

EACL 2026

**The 19th Conference of the European Chapter of the  
Association for Computational Linguistics**

**Proceedings of System Demonstrations**

March 24-29, 2026

The EACL organizers gratefully acknowledge the support from the following sponsors.

**Platinum**



Megagon Labs

**Appen**

**Bronze**

**elra**

 translated.



©2026 Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)  
317 Sidney Baker St. S  
Suite 400 - 134  
Kerrville, TX 78028  
USA  
Tel: +1-855-225-1962  
[acl@aclweb.org](mailto:acl@aclweb.org)

ISBN 979-8-89176-382-1

## Introduction

Welcome to the proceedings of the System Demonstration track of the 19th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2026), held in Rabat, Morocco, from March 24 to 29, 2026.

For the EACL 2026 System Demonstration track, we received 102 submissions, almost doubling the number of submissions compared to the previous edition. Of these, 44 were selected for inclusion in the program (acceptance rate of 43.1%). The selection process was highly competitive, and many strong submissions could unfortunately not be accommodated.

We would like to thank the members of the program committee and area chairs for their careful and timely reviews. We are also grateful to the organizing and publication teams for their support. Finally, we sincerely thank all authors who submitted their work to the Demonstrations track and contributed to the scientific quality of the program.

The accepted demonstration papers reflect the breadth and vitality of current research in Natural Language Processing, covering evaluation frameworks, interpretability tools, dialogue systems, knowledge-based applications, and responsible AI platforms. The demonstrations will be presented in person during the conference, including dedicated poster and demo sessions.

We appreciate the efforts made by all authors to showcase their work and contribute to the success of EACL 2026.

Danilo Croce, Jochen L. Leidner, Nafise Sadat Moosavi  
EACL 2026 System Demonstration Chairs

# Program Committee

## Chairs

Danilo Croce, University of Roma Tor Vergata  
Jochen L. Leidner, KnowledgeSpaces UG (haftungsbeschränkt), Coburg University of Applied Sciences and University of Sheffield  
Nafise Sadat Moosavi, University of Sheffield

## Area Chairs

Hiba Arnaout, Technische Universität Darmstadt  
Long Bai, Institute of Computing Technology, Chinese Academy of Sciences  
Simone Balloccu, Technische Universität Darmstadt  
Valerio Basile, University of Turin  
Giuseppe Castellucci, Amazon  
Craig Erickson, Amazon  
Junxian He, Hong Kong University of Science and Technology  
Christian Heumann, Ludwig-Maximilians-Universität München  
Shaoxiong Ji, University of Turku and ELLIS Institute Finland  
Khalid Al Khatib, University of Groningen  
Christopher Klamm, University of Cologne and University of Mannheim  
Anne Lauscher, Universität Hamburg  
Kelly Marchisio, Cohere and Cohere  
John Philip McCrae, National University of Ireland Galway  
Qingyun Wang, College of William and Mary  
Hua Wei, Arizona State University

## Reviewers

Angus Addlesee, Amazon  
Suman Adhya, Indian Association for the Cultivation of Science  
Yamen Ajjour, Universität Hannover  
Seyed Alireza Mousavian Anaraki, University of Roma Tor Vergata  
Miriam Anschütz, Technische Universität München  
Mario Ezra Aragon, Universidad de Santiago de Compostela  
Matthias Aßenmacher, Ludwig-Maximilians-Universität München  
Joanne Boisson, Cardiff University  
Federico Borazio, University of Roma Tor Vergata  
Georgeta Bordea, La Rochelle Université  
Marco Braga, University of Milan - Bicocca  
Marco Antonio Sobrevilla Cabezudo, Pontificia Universidad Católica del Perú  
Aunabil Chakma, University of Arizona  
Raghuveer Chanda, Google  
Nischal Reddy Chandra, Adobe Systems  
Tiejin Chen, Arizona State University  
Daiwei Chen, University of Wisconsin - Madison  
Chung-Chi Chen, AIST, National Institute of Advanced Industrial Science and Technology  
Kehai Chen, Harbin Institute of Technology (Shenzhen)

Sanjay Das, University of Texas at Dallas  
Ona de Gibert, University of Helsinki  
Giacomo De Luca, University of Roma Tor Vergata  
Jacob Devasier, University of Texas at Arlington  
Luigi Di Caro, University of Turin  
Kshitij P Fadnis, IBM Research  
Enfa Fane, University of Arizona  
Mariia Fedorova, University of Oslo  
Nils Feldhus, Technische Universität Berlin  
Tim Fischer, University of Hamburg  
Esteban Garces Arias, Ludwig-Maximilians-Universität München  
Javier García Gilibert, Barcelona Supercomputing Center  
Tirthankar Ghosal, Oak Ridge National Laboratory  
Voula Giouli, Aristotle University of Thessaloniki and ILSP - AthenaResearch Center  
George Gkotsis, Upwork Inc  
Xuehang Guo, College of William and Mary  
Claudiu Daniel Hromei, University of Roma Tor Vergata  
Dayu Hu, Northeastern University  
Timour Igamberdiev, Universität Vienna  
Oghenevovwe Ikumariogbe, University of Arizona  
Hasan Iqbal, Mohamed bin Zayed University of Artificial Intelligence  
Vidhyakshaya Kannan, Georgia Institute of Technology  
Yoshihide Kato, Nagoya University  
Jeonghwan Kim, University of Illinois at Urbana-Champaign  
Andrei Kucharavy, University of Applied Sciences Western Switzerland, Sierre (HES-SO Valais)  
Leo Leppänen, University of Helsinki  
Weijiang Li, University of Notre Dame  
Zixuan Li, Institute of Computing Technology, Chinese Academy of Sciences  
Ruo Chen Li, University of Texas at Dallas  
Hanming Li, Tsinghua University  
Xiaochang Li, William and Mary  
Shuhang Lin, , Rutgers University  
Yueqian Lin, Duke University  
Ke Lin, College of William and Mary  
Zhexiong Liu, University of Pittsburgh  
Junteng Liu, The Hong Kong University of Science and Technology  
Wenxuan Liu, Chinese Academy of Sciences  
Tianrui Liu, Apple  
Xinyu Lu, Chinese Academy of Sciences  
Maria Mahbub, Oak Ridge National Laboratory  
Puneet Mathur, Adobe Systems  
Tim Menzner, Hochschule Coburg  
Yasmin Moslem, Trinity College Dublin  
Youyang Ng, Kioxia Corporation  
Tuan-Phong Nguyen, VNU University of Engineering and Technology  
Huy V. Nguyen, Amazon  
Leyi Pan, Tsinghua University  
Jungyeul Park, Korea Advanced Institute of Science & Technology  
Max Ploner, Humboldt Universität Berlin  
Cheng Qian, University of Illinois at Urbana-Champaign  
Yide Ran, Stevens Institute of Technology

Prajvi Saxena, German Research Center for AI  
Carsten Schnober, Netherlands eScience Center  
Tim Schopf, Technische Universität Dresden  
Seongbum Seo, Sejong University  
Alay Dilipbhai Shah, Together AI  
KaShun Shum, Hong Kong University of Science and Technology  
Gosuddin Kamaruddin Siddiqi, Microsoft  
Jyotika Singh, Oracle  
Bhavyajeet Singh, Technische Universität Darmstadt  
Jiahe Song, Shanghai Jiaotong University and Shanghai Artificial Intelligence Laboratory  
Daniil Sorokin, Amazon Development Center Germany  
Dachun Sun, Worcester Polytechnic Institute  
Qiang Sun, University of Western Australia  
Marek Suppa, Comenius University in Bratislava  
Sotaro Takeshita, University of Technology Nuremberg and Universität Mannheim  
Yudong Tao, Facebook  
Jonathan Tonglet, Technische Universität Darmstadt  
Vatsal Venkatkrishna, Institute for Computer Science, Artificial Intelligence and Technology  
Sami Virpioja, University of Helsinki  
Doan Nam Long Vu, Technische Universität Darmstadt  
Chuan-Ju Wang, Academia Sinica  
Zibu Wei, University of California  
Liang Xie, Guangdong University of Technology  
Shuo Xing, Google and Texas A&M University - College Station  
Xin Xu, The Hong Kong University of Science and Technology  
Bingbing Xu, Renmin University of China  
Hao Yu, National University of Defense Technology  
Pengfei Yu, Amazon  
Yige Yuan, University of Washington  
Linfan Zhang, University of California  
Chongsheng Zhang, Ludwig-Maximilians-Universität München and Henan University  
Tong Zhou, Chinese Academy of Science  
Zichen Zhu, Shanghai Jiao Tong University  
Vilém Zouhar, ETHZ - ETH Zurich  
Dennis Zyska, Technische Universität Darmstadt

## Table of Contents

<i>Stakeholder Suite: A Unified AI Framework for Mapping Actors, Topics and Arguments in Public Debates</i>	
Mohamed Chenene, Jeanne Rouhier, Jean Daniélou, Mihir Sarkar and Elena Cabrio . . . . .	1
<i>DeepPavlov Strikes Back: A Toolkit for Improving LLM Reliability and Trustworthiness</i>	
Evgenii Nikolaev, Timur Ionov, Anna Korzanova, Vasily Konovalov and Maksim Savkin . . . . .	21
<i>PropGenie: A Multi-Agent Conversational Framework for Real Estate Assistance</i>	
Chang Shen, Shaozu Yuan, Kuizong Wu, Long Xu and Meng Chen . . . . .	33
<i>Pro-QuEST: A Prompt-chain based Quiz Engine for testing Specialized Technical Product Knowledge</i>	
Sujatha Das Gollapalli, Mouad Hakam, Mingzhe Du, See-Kiong Ng and Mohammed Hamzeh	46
<i>elfen: A Python Package for Efficient Linguistic Feature Extraction for Natural Language Datasets</i>	
Maximilian Maurer . . . . .	61
<i>DELTA: A Toolkit for Measuring Linguistic Diversity in Dependency-Parsed Corpora</i>	
Louis Estève and Kaja Dobrovoljc . . . . .	75
<i>CLARIESG: An End-to-End System for ESG Analysis over Complex Tables in Corporate Reports</i>	
Marta Santacroce, Michele Luca Contalbo, Sara Pederzoli, Riccardo Benassi, Venturelli Valeria, Matteo Paganelli and Francesco Guerra . . . . .	86
<i>Fact Finder - Enhancing Domain Expertise of Large Language Models by Incorporating Knowledge Graphs</i>	
Daniel Steinigen, Roman Teucher, Timm Heine Ruland, Max Rudat, Nicolas Flores-Herr, Peter Fischer, Nikola Milosevic, Christopher Schymura and Angelo Ziletti . . . . .	101
<i>Simplifying Outcomes of Language Model Component Analyses with ELIA</i>	
Aaron Louis Eidt and Nils Feldhus . . . . .	111
<i>IntelliCode: A Multi-Agent LLM Tutoring System with Centralized Learner Modeling</i>	
Jones David and Shreya Ghosh . . . . .	129
<i>FiMMIA: scaling semantic perturbation-based membership inference across modalities</i>	
Anton Emelyanov, Sergei Kudriashov and Alena Fenogenova . . . . .	139
<i>A Browser-based Open Source Assistant for Multimodal Content Verification</i>	
Rosanna Milner, Michael Foster, Olesya Razuvayevskaya, Valentin Porcellini, Denis Teyssou, Ian Roberts and Kalina Bontcheva . . . . .	154
<i>Infherno: End-to-end Agent-based FHIR Resource Synthesis from Free-form Clinical Notes</i>	
Johann Frei, Nils Feldhus, Lisa Raithel, Roland Roller, Alexander Meyer and Frank Kramer .	163
<i>BOOM: Beyond Only One Modality KIT's Multimodal Multilingual Lecture Companion</i>	
Sai Koneru, Fabian Retkowski, Christian Huber, Lukas Hilgert, Seymanur Akti, Enes Yavuz Ugan, Alexander Waibel and Jan Niehues . . . . .	175
<i>PEFT-Factory: Unified Parameter-Efficient Fine-Tuning of Autoregressive Large Language Models</i>	
Robert Belanec, Ivan Srba and Maria Bielikova . . . . .	188
<i>Similar, but why? A Toolkit for Explaining Text Similarity</i>	
Juri Opitz, Andrianos Michail, Lucas Moeller, Sebastian Padó and Simon Clematide . . . . .	203

<i>AlignFix: A Tool for Parallel Corpora Augmentation and Refinement</i> Samuel Frontull and Simon Haller-Seeber .....	215
<i>PromptLab: A Collaborative Platform for Prompt Engineering and Dataset Curation</i> Maged S. Al-shaibani, Zaid Alyafeai, Dania Refai, Nawaf Alomari, Ahmed Ashraf, Mais Alheraki, Mustafa Alturki, Hamzah Luqman and Irfan Ahmad .....	225
<i>LLM BiasScope: A Real-Time Bias Analysis Platform for Comparative LLM Evaluation</i> Himel Ghosh and Nick Elias Werner .....	261
<i>InkSight: Towards AI-Aided Historical Manuscript Analysis</i> Andrey Sakhovskiy, Ivan Ulitin, Emilia Bojarskaja, Vladimir Kokh, Ruslan Murtazin, Maxim Novopoltsev and Semen Budenny .....	271
<i>promptlotion: A Unified, Modular Framework for Prompt Optimization</i> Tom Zehle, Timo Heiß, Moritz Schlager, Matthias Aßenmacher and Matthias Feurer .....	282
<i>T-pro 2.0: An Efficient Russian Hybrid-Reasoning Model and Playground</i> Dmitrii Stoianov, Danil Taranets, Olga Tsymboi, Ramil Latypov, Almaz Dautov, Vladislav Kru- glikov, Surkov Nikita, German Abramov, Pavel Gein, Dmitry Abulkhanov, Mikhail Gashkov, Viktor Zelenkovskiy, Artem Batalov, Aleksandr Medvedev and Anatolii Potapov .....	297
<i>SDialog: A Python Toolkit for End-to-End Agent Building, User Simulation, Dialog Generation, and Evaluation</i> Sergio Burdisso, Séverin Baroudi, Yanis Labrak, David Grünert, Pawel Cyrt, Yiyang Chen, Srikanth Madikeri, Esaú Villatoro-tello, Ricard Marxer and Petr Motlicek .....	320
<i>Agentic AI for Human Resources: LLM-Driven Candidate Assessment</i> Kamer Ali Yuksel, Abdul Basit Anees, Ashraf Hatim Elneima, Sanjika Hewavitharana, Mohamed Al-Badrashiny and Hassan Sawaf .....	341
<i>Trove: A Flexible Toolkit for Dense Retrieval</i> Reza Esfandiarpour, Max Zuo and Stephen Bach .....	349
<i>ClinicalTrialsHub: Bridging Registries and Literature for Comprehensive Clinical Trial Access</i> Jiwoo Park, Ruoqi Liu, Avani Jagdale, Andrew Srisuwananukorn, Jing Zhao, Ping Zhang and Sachin Kumar .....	359
<i>SciTrue: Evidence-Grounded Claim Verification in Science</i> Neset Tan, Minghao LI and Mark Gahegan .....	397
<i>PUCP-Matrix: An Open-source and Comprehensive Toolkit for Linguistic Analysis of Spanish Texts</i> Javier Alonso Villegas Luis and Marco Antonio Sobrevilla Cabezudo .....	407
<i>Integrity Shield A System for Ethical AI Use &amp; Authorship Transparency in Assessments</i> Ashish Raj Shekhar, Shiven Agarwal, Priyanuj Bordoloi, Yash Shah, Tejas Anvekar and Vivek Gupta .....	417
<i>Using a Human-AI Teaming Approach to Create and Curate Scientific Datasets with the SciLire System</i> Necva Bölücü, Jessica Irons, Changhyun Lee, Brian Jin, Maciej Rybinski, Huichen Yang, Andreas Duenser and Stephen Wan .....	428
<i>xLM: A Python Package for Non-Autoregressive Language Models</i> Dhruvesh Patel, Durga Prasad Maram, Sai Sreenivas Chintha, Benjamin Rozonoyer and Andrew McCallum .....	445

<i>AITutor-EvalKit: Exploring the Capabilities of AI Tutors</i>	
Numaan Naeem, Kaushal Kumar Maurya, Kseniia Petukhova and Ekaterina Kochmar . . . . .	457
<i>EvalSense: A Framework for Domain-Specific LLM (Meta-)Evaluation</i>	
Adam Dejl and Jonathan Pearson . . . . .	480
<i>AI for Climate Finance: Agentic Retrieval and Multi-Step Reasoning for Early Warning System Investments</i>	
Saeid Vaghefi, Aymane Hachcham, Veronica Grasso, Nakiete Msemu, Chiara Colesanti Senni and Markus Leippold . . . . .	492
<i>RAGVUE: A Diagnostic View for Explainable and Automated Evaluation of Retrieval-Augmented Generation</i>	
Keerthana Murugaraj, Salima Lamsiyah and Martin Theobald . . . . .	512
<i>SmartMatch: Real-Time Semantic Retrieval for Translation Memory Systems</i>	
Ernesto Luis Estevanell Valladares, Salima Lamsiyah, Alicia Picazo-Izquierdo, Tharindu Ranasinghe, Ruslan Mitkov and Rafael Munoz . . . . .	527
<i>QSTN: A Modular Framework for Robust Questionnaire Inference with Large Language Models</i>	
Maximilian Kreutner, Jens Rupperecht, Georg Ahnert, Ahmed Salem and Markus Strohmaier . . . . .	537
<i>A Virtual Assistant for Architectural Design in a VR Environment</i>	
Ander Salaberria, Oier Ijurco, Markel Ferro, Jiayuan Wang, Iñigo Vilá Muñoz, Roberto de Ioris, Jeremy Barnes and Oier Lopez De Lacalle . . . . .	550
<i>ARGSBASE: A Multi-Agent Interface for Structured Human–AI Deliberation</i>	
Frieso Turkstra, Sara Nabhani and Khalid Al Khatib . . . . .	563
<i>Simultaneous Speech-to-Text Translation Web Application for Estonian</i>	
Bohdan Podziubanchuk, Aivo Olev, Jiaming Kong and Tanel Alumäe . . . . .	575
<i>The AI Committee: A Multi-Agent Framework for Automated Validation and Remediation of Web-Sourced Data</i>	
Sunith Vallabhaneni, Thomas Berkane and Maimuna S. Majumder . . . . .	583
<i>entity-linkings: A Unified Library for Entity Linking</i>	
Yuya Sawada, Tsuyoshi Fujita, Yusuke Sakai and Taro Watanabe . . . . .	591
<i>ESG-KG: A Multi-modal Knowledge Graph System for Automated Compliance Assessment</i>	
Li-Yang Chang, Chih-Ming Chen, Hen-Hsen Huang, Ming-Feng Tsai, An-Zi Yen and Chuan-Ju Wang . . . . .	602
<i>BanSuite: A Unified Toolkit and Software Platform for Low-Resource NLP in Bangla</i>	
Md. Abu Sayed, Faisal Ahamed Khan, Jannatul Ferdous Tuli, Nabeel Mohammed, Mohammad Ruhul Amin and Mohammad Mamun Or Rashid . . . . .	609

# Stakeholder Suite: A Unified AI Framework for Mapping Actors, Topics and Arguments in Public Debates

Mohamed Chenene<sup>1</sup> Jeanne Rouhier<sup>1</sup> Jean Daniélou<sup>2</sup> Mihir Sarkar<sup>3</sup> Elena Cabrio<sup>4</sup>

<sup>1</sup>ENGIE Lab CRIGEN, France

<sup>2</sup>Centre de Sociologie de L'Innovation, CNRS, UMR 9217, Mines ParisTech, PSL University

<sup>3</sup>ENGIE Research & Innovation, France

<sup>4</sup>Université Côte d'Azur, CNRS, INRIA, I3S, France

mohamed.chenene1@gmail.com, jeanne.rouhier@external.engie.com,

jeandanielou@live.fr, mihir.sarkar@engie.com, elena.cabrio@univ-cotedazur.fr

## Abstract

Public debates surrounding infrastructure and energy projects involve complex networks of stakeholders, arguments, and evolving narratives. Understanding these dynamics is crucial for anticipating controversies and informing engagement strategies. This paper presents **Stakeholder Suite**, a framework deployed in operational contexts for mapping actors, topics, and arguments within public debates. The system combines actor detection, topic modeling, argument extraction and stance classification in a unified pipeline. Tested on multiple energy infrastructure projects as a case study, the approach delivers fine-grained, source-grounded insights while remaining adaptable to diverse domains. The framework achieves strong retrieval precision and stance accuracy, producing arguments judged relevant in 75% of pilot use cases. Beyond quantitative metrics, the tool has proven effective for operational use: helping project teams visualize networks of influence, identify emerging controversies, and support evidence-based decision-making.

## 1 Introduction

In the era of social media, public controversies can emerge rapidly and significantly affect the reputation and operations of organizations. Controversy analysis aims to systematically track and interpret such public debates by identifying the stakeholders involved, their relationships, their stances and the evolution of discourse over time. A structured understanding of these dynamics enables organizations to anticipate conflicts, adapt communication strategies, and mitigate potential delays or financial losses. In the energy sector, large-scale renewable infrastructure projects (e.g., wind farms, solar plants, district heating, hydrogen, or ammonia transport) frequently trigger public consultations and debates (LaPatin et al., 2023). Numerous stakeholders (citizens, associations, local officials and activists) express their views across diverse

digital channels, making the monitoring and analysis of these discussions both complex and time-consuming for project developers.

This work addresses the question of how computational methods can support and enhance a transparent analysis of public debates. Our objective is to develop a *generic and adaptable framework* capable of mapping stakeholders, identifying debated topics, and extracting arguments and stances to provide an exhaustive and objective view of controversial debates. We introduce **Stakeholder Suite**, that combines Large Language Models and Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) to build a structured representation of public debates. LLMs (Wei et al., 2021; Brown et al., 2020; Chen et al., 2024) now perform effectively for various NLP tasks, enabling fast deployment of complex analysis pipelines without large annotated datasets. RAG improves factual grounding and transparency by coupling generation with evidence selection, mitigating hallucinations, and improving traceability (Kalai et al., 2025).

This paper makes the following contributions: 1. We present a transparent and solid framework for the analysis of public debates, which jointly models stakeholders, topics, and arguments within a unified data structure; 2. We propose a fine-grained approach to topic modeling and argument analysis that leverages LLM-based reasoning, enabling richer and contextual insights. To date, the system has been deployed across more than ten renewable energy projects, offering project teams actionable insights into local dynamics and stakeholder perceptions. To the best of our knowledge, this is among the first operational frameworks that combine stakeholder mapping, fine-grained argument extraction, and topic modeling within a single and updatable pipeline spanning heterogeneous texts. A demo video is available.<sup>1</sup>

<sup>1</sup><https://youtu.be/c1b6QgyVkw5>

## 2 Related Work

Mapping stakeholders and their stances across public debates intersects with several markets and research areas: media intelligence, social listening, public affairs, controversy analysis, policy analysis and argument mining (AM).

**Commercial solutions.** Social listening and sentiment analytics platforms as Talkwalker, Brandwatch or Meltwater, aggregate large-scale social and web data, offering dashboards for trends, topics, influencer identification, and sentiment. A second line of commercial tools centers on network views of entities, topics, and audiences. NetBase Quid and Pulsar Platform provide interactive graphs, topic clustering, community analysis, and audience segmentation across social and news sources. These capabilities help reveal ecosystems (e.g., co-mention networks and narrative clusters) and identify influential accounts. The Stakeholder Company (TSC.ai) maintains a large stakeholder repository with influence mapping and engagement workflows and comes closest to stakeholder-centric use cases in practice. Public affairs and government-relations platforms provide structured data and workflows for legislative and regulatory engagement. Quorum integrates bills, hearings, and institutional actors across multiple jurisdictions, supporting monitoring, alerting, and outreach. Some products also incorporate retrieval-augmented assistants for quick knowledge access. While effective within institutional policy contexts, these systems are not intended as general-purpose frameworks across heterogeneous public debates.<sup>2</sup>

**Academic research.** Foundational work in controversy analysis (Venturini and Munk, 2022) establishes sociological methods for systematically mapping controversies and constructing networks from debates. Automated policy analysis applies NLP pipelines to structure regulatory content, extracting measures, actors, and impacts from environmental and institutional documents (Firebanks-Quevedo et al., 2022; Singh et al., 2024). For instance, (Planas et al., 2022)’s framework is a knowledge graph-oriented approach that allows rapid re-

<sup>2</sup>Commercial platforms discussed: Talkwalker (<https://www.talkwalker.com>), Brandwatch (<https://www.brandwatch.com>), Meltwater (<https://www.meltwater.com>), NetBase Quid (<https://www.quid.com>), Pulsar Platform (<https://www.pulsarplatform.com>), TSC.ai (<https://tsc.ai>), Quorum (<https://www.quorum.us>).

view of policy documents through entity search, topic analysis, and policy search. While these approaches enable efficient document analysis, they typically operate on formal policy texts rather than public debate corpora. AM offers complementary methods for analyzing contested discourse, including claim detection, stance classification, and argument clustering (Daxenberger et al., 2020; Slonim et al., 2021) in different kinds of structured contexts: essays, online debate platforms, legal documents, and political debates. Political debate analysis, which examines clash points, strategies, and argumentative structures (South et al., 2020; Chen et al., 2025; Goffredo et al., 2023), shares similarities with energy infrastructure debates in involving political actors, citizen associations, and multi-channel information flows. While prior work provides valuable components (e.g., large-scale data access, influencer and network analytics, institutional monitoring, and argument mining methods), an end-to-end production-oriented approach that jointly maps stakeholders, detects and organizes topics, and extracts source-grounded arguments with stance across heterogeneous debate corpora has not yet been proposed. This gap motivates the development of retrieval-augmented frameworks adaptable across domains where public debate and stakeholder engagement are central.

## 3 Methodology

The goal of **Stakeholder Suite** (Figure 1) is to provide a comprehensive snapshot of a public debate. Its components are designed to address three main questions: (1) *Who is speaking?* → Actor detection; (2) *What are they talking about?* → Topic modeling and extraction; (3) *What is their position?* → Argument mining and stance classification.

### 3.1 Data Collection

To start, two main components are required:

- (1) **Stakeholder list:** manually or automatically identified actors participating in the debate;
- (2) **Document database:** a corpus integrating diverse textual sources on the debate, including: (a) **Debate-centric data:** press articles, policy papers and administrative documents, typically obtained through data brokers, web scraping or manually uploaded; (b) **Actor-centric data:** social media posts, websites and blogs associated with listed stakeholders.

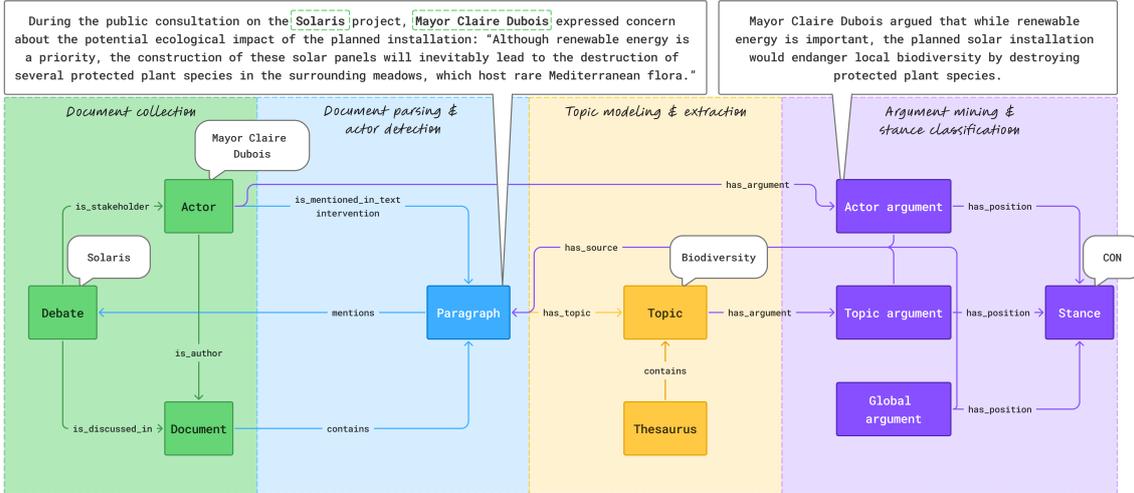


Figure 1: Blueprint of the **Stakeholder Suite** framework with an actor argument example: (1) Construction of a database containing stakeholders and documents related to the public debate; (2) Parsing of documents into paragraphs and detection of actor mentions and interventions within the text; (3) Linking of paragraphs to topics from the thesaurus using semantic similarity; (4) Generation of arguments through a RAG pipeline along three analytical axes (global, topic-specific, actor-specific) together with automatic stance classification for each argument.

**Data quality.** Both data types are essential for ensuring that stakeholders are properly represented within the corpus. The most informative argumentative content often originates from stakeholder-generated material. Based on our experience in the energy sector, high-yield argumentative sources include public debate transcripts, stakeholder blog posts, and social media statements. Conversely, press and institutional documents tend to be less argumentative but provide valuable contextual information for topic modeling.

### 3.2 Data parsing and actor detection

**Data parsing.** After constructing the database, the next step is to structure it and establish semantic links between entities. Given that sentences are often too short to capture full meaning, while documents can vary greatly in length, we adopt the paragraph as the atomic analysis unit, as it provides a balance between contextual coherence and processing granularity. In our implementation, paragraph segmentation, coreference resolution, and actor speech extraction are handled jointly through a single LLM prompt (see Appendix A.1). All paragraphs are embedded and indexed in OpenSearch vector database for full-text and semantic search.

**Actor detection.** Once paragraphs are extracted, we link them to the corresponding actors from the stakeholder list. We define three relation types:

- (1) *is\_author*: the actor of the document (e.g., blog post, tweet), known from the actor-centric corpus.
- (2) *is\_mentioned\_in\_text*: explicit textual mentions of an actor in a paragraph. To improve the precision for common or ambiguous names, we search for direct matches using the known actor list and verify them with FLAIR NER (Akbik et al., 2019).
- (3) *intervention*: actor speech segments in debate transcripts. Interventions are more valuable content than mentions because they report the actor’s speech. We use the LLM parser mentioned above.

Mentions of the debate are also detected at the paragraph level using a predefined list of aliases. This step is necessary because, even if a document refers to the debate as a whole, not all paragraphs within it necessarily discuss the debate directly.

### 3.3 Topic modeling and extraction

With a structured database of paragraphs linked to actors, the next step is to identify the topics discussed. Applying standard unsupervised clustering to highly domain-specific corpora often produce generic or irrelevant topics. Following (Tamkin et al., 2024), we leverage LLMs to construct a domain-specific thesaurus: a structured list of topics and their corresponding descriptions. The process involves the following steps:

- (1) **Corpus chunking**: split the database into coherent subsets according to the debate characteris-

tics. For long-running cases, use temporal windows (e.g., monthly or quarterly); for political debates, segment by speaker or party.

(2) **Chunk description via RAG:** for each subset, run a RAG with a query such as: “*What are the emerging topics related to {debate\_name}?*”, and extract the resulting topic candidates.

(3) **Topic clustering:** cluster the generated topics using methods such as HDBSCAN or Spectral Clustering based on a semantic affinity matrix (Campello et al., 2013; Shi and Malik, 2000).

(4) **Cluster summarization:** apply an LLM to name each cluster and provide representative examples, producing two descriptive levels (topic and subtopic) to preserve granularity.

(5) **Thesaurus aggregation:** run a final LLM pass to consolidate results into a unified thesaurus containing topics, descriptions, and subtopics.

Although this pipeline heavily relies on LLMs, it has proved most effective for generating high-resolution topic descriptions. Moreover, analyzing topics across temporal slices reveals trends and weak signals of business relevance. Empirically, naive clustering methods consistently failed to capture debate-specific nuances, whereas our LLM-augmented approach produced coherent, non-generic topics. Appendix A.2 provides an example of a thesaurus built with this approach and the prompts used for steps 4 and 5. Once the thesaurus is defined, topic occurrences are computed across the corpus. We compute cosine similarity between topic descriptions and paragraph embeddings, retaining matches above a predefined threshold ( $t = 0.15$ ) to balance quality and quantity of paragraphs per topic. When multiple topics qualify, the top three by similarity are selected.

### 3.4 AM and stance classification

We adopt a RAG architecture for argument detection, because (i) it is easily implemented in production environments through LLM API calls (in particular, we use gpt-4o-mini (OpenAI et al., 2024) on the Azure platform, selected for compliance and cost-quality balance); (ii) it enables the generation of concise, one-sentence arguments that are easily interpretable by end users; (iii) it preserves traceability by maintaining explicit links between generated outputs and their source documents.

A known limitation of RAG, however, is its non-exhaustive coverage of the corpus, as it retrieves only the top- $K$  relevant passages prior to generation. To mitigate this, we partition the database

into coherent chunks and query each subset independently, allowing the extraction of more than 300 arguments per debate. This strategy achieves a practical balance between corpus coverage, computational cost, and analytical depth.

#### 3.4.1 Temporal and Dimensional Splits

To maximize coverage of the public debate and reveal dynamics, we partition the corpus along two axes: (i) **time**, to trace how arguments evolve; and (ii) **dimension**, to focus on *stance*, *actors*, or *topics* derived from previous steps. We operationalize this with three uniform query families.

(1) **Global argument. Goal:** enumerate debate arguments *in favor of* or *against* for each year.

**Database filter (pseudo-Cypher):**

```
MATCH (p:PARAGRAPH) WHERE p.date.year = $year
RETURN p
```

**Prompt to LLM:** “*What are all the arguments {in favor/against} {debate\_name}?*”

**Rationale:** yearly stance snapshots surface trends and turning points.

(2) **Actor argument. Goal:** extract arguments attributable to a specific stakeholder, prioritizing paragraphs that are most on-topic for the target debate (no temporal split due to data sparsity).

**Database filter (pseudo-Cypher):**

```
(a:ACTOR)->(p:PARAGRAPH)-[:MENTIONS]->(d:DEBATE)
RETURN p AS p_debate
(a:ACTOR)->(p:PARAGRAPH)-[:HAS_TOPIC]->(t:TOPIC)
RETURN p AS p_topic
(a:ACTOR)->(p:PARAGRAPH) RETURN p AS p_all
IF COUNT(p_debate) > 25 THEN RETURN p_debate
ELSE IF COUNT(p_topic) > 25 THEN RETURN p_topic
ELSE RETURN p_all
```

**Prompt to LLM:** “*What are the arguments of the stakeholder {actor\_name} about {debate\_name}?*”

**Rationale:** We prioritize *debate-specific* paragraphs for precision; if these are insufficient ( $\leq 25$ ), we fall back to *topic-linked* paragraphs, and otherwise include *all* paragraphs for coverage. The threshold was determined empirically from a study of more than 150 actors and was optimized to maximize the number of stances extracted.

(3) **Topic argument. Goal:** collect arguments tied to a specific topic, per year, decoupled from the debate stance if necessary.

**Database filter (pseudo-Cypher):**

```
MATCH (p:PARAGRAPH)-[:HAS_TOPIC]->(t:MACRO_TOPIC)
WHERE p.date.year = $year AND t.name = $topic_name
RETURN p
```

**Prompt to LLM:** “*What are the arguments related to the topic {topic\_name}?*”

**Rationale:** topics may drift from debate framing; isolating them improves thematic resolution. It can also give a broader vision of the subject.

For a setup with  $\sim 50$  actors, 8 topics, and a 5-year span, and assuming  $\sim 3$  arguments are generated per query we get approximately 300 arguments per debate. Combining temporal and dimensional partitions with standardized prompts yields a scalable, interpretable, and trend-aware argument set while keeping retrieval focused and noise low.

### 3.4.2 RAG Pipeline

The RAG pipeline structures the final argument extraction process. For each query identifier, it operates in four sequential stages:

(1) **Retrieval:** the database is pre-filtered according to the previously mentioned split. We retrieve the top- $K$  paragraphs ( $K = 25$ ) using Maximum Marginal Relevance (Carbonell and Goldstein, 1998) with  $\lambda = 0.8$  to balance relevance and diversity.

(2) **Argument generation:** the query and retrieved paragraphs are passed to the LLM, which returns a structured list of arguments, each explicitly linked to one or more source paragraphs. The output is then parsed using regex to extract the argument text and the corresponding source identifiers.

(3) **LLM-as-a-Judge** (Zheng et al., 2023): a second LLM performs an automatic quality check, verifying that each generated argument is supported by its sources and flagging potential hallucinations.

(4) **Stance classification:** each validated argument is assigned a stance label (PRO, CON, or NEUTRAL). Appendix A.3 shows the prompts for steps 2 to 4, together with examples of extracted arguments.

### 3.5 Visualization: Stakeholder Mapping

To visualize stakeholder relationships, we construct a connected network where each node represents an actor. Node color encodes stance (green = pro, red = con, grey = neutral), and node size reflects the frequency of debate mentions during actor interventions. Edges are drawn when one actor explicitly references another within an intervention, capturing direct interactions and influence links. An example can be found in Appendix A.4.

## 4 Evaluation

As previously described, our framework addresses three main tasks, i.e., NER for actor detection, semantic text similarity for topic extraction and RAG for argument generation. Since for the first two we

rely on standard, validated approaches (Akbik et al., 2018; Muennighoff et al., 2023), in the following we focus on the RAG pipeline, with emphasis on retrieval, generation, and stance classification.

### 4.1 Retrieval Evaluation

Retrieval performance was assessed across eight energy projects for both PRO and CON stances (16 queries in total). Each evaluation mixed relevant and randomly sampled paragraphs to approximate realistic corpus conditions, producing datasets of 200 candidates per query, with an average of about 35 relevant documents. We report Precision@K as the primary metric, focusing on diversity–relevance trade-offs using Maximum Marginal Relevance (MMR). As shown in Table 1, optimal performance occurs for  $\lambda \in [0.7, 0.8]$ , confirming that moderate diversity does not deteriorate retrieval quality. The mean Precision@20 reaches 0.59, with variability driven by project-specific content density, sufficient for reliable argument generation downstream.

MMR $\lambda$	Precision@10	Precision@20
0.5	0.41 $\pm$ 0.13	0.40 $\pm$ 0.09
0.6	0.59 $\pm$ 0.17	0.55 $\pm$ 0.14
0.7	<b>0.64<math>\pm</math>0.20</b>	0.57 $\pm$ 0.18
0.8	0.64 $\pm$ 0.25	<b>0.59<math>\pm</math>0.17</b>
0.9	0.63 $\pm$ 0.26	0.58 $\pm$ 0.19
1.0	0.63 $\pm$ 0.26	0.58 $\pm$ 0.19

Table 1: Retrieval performance across different  $\lambda$  values. Metrics are reported as mean  $\pm$  standard deviation.

Component	Class	Precis.	Recall	F1
<b>LLM-as-a-Judge</b>	BAD	0.42	0.75	0.54
	GOOD	0.77	0.45	0.57
	<b>Macro avg</b>	<b>0.59</b>	<b>0.60</b>	<b>0.55</b>
<b>Stance Classification (ACTOR)</b>	AGAINST	0.96	0.93	0.94
	NEUTRAL	0.69	0.76	0.73
	PRO	0.66	0.60	0.63
	<b>Macro avg</b>	<b>0.77</b>	<b>0.76</b>	<b>0.77</b>
<b>Stance Classification (TOPIC)</b>	AGAINST	1.00	0.61	0.76
	NEUTRAL	0.72	0.89	0.80
	PRO	0.78	0.70	0.74
	<b>Macro avg</b>	<b>0.83</b>	<b>0.73</b>	<b>0.76</b>

Table 2: Performance of the LLM-as-a-Judge and Stance Classification modules.

### 4.2 LLM-as-a-Judge: Generation Evaluation

To assess generation quality, a secondary LLM evaluated each argument-source pair to identify unsupported or hallucinated outputs. Our primary objective was to minimize false positives visible to end users, prioritizing a low error rate over exhaustive recall. In practice, this meant maximizing BAD argument recall while maintaining acceptable coverage

of GOOD ones, to preserve user trust and prevent misinterpretation of the system’s outputs. Several prompt configurations were tested to balance strictness and recall. On a validation set of 739 samples, the best-performing validator achieved high BAD recall to protect users from unsupported claims, accepting some loss of GOOD recall (Table 2).

### 4.3 Stance Classification

Following the cross-topic argument classification approach proposed by (Stab et al., 2018), 3 annotators with expertise in Computational Linguistics annotated 350 arguments with their true stance toward one energy project to evaluate LLM classification performance. As shown in Table 2, we report the overall performance. For reference, (Reimers et al., 2019) estimates human performance on cross-topic stance classification at 0.81 (F1-score). Our results fall close to this upper bound, indicating reliable performance for our application.

## 5 Case Studies

### 5.1 A Case Study for Users Evaluation

User evaluations globally demonstrate that the system’s analytical outputs are both accurate and operationally valuable (Appendix A.4 shows screenshots of the application). The application was used by 10 end-users across 9 different energy pilots. During onboarding, roughly 600 documents were processed per project, followed by about 20 new documents per month for continuous monitoring. Users particularly valued the argument summaries and mapping visualizations for improving internal communication and evidence-based engagement planning (“*Stakeholder Suite gives us a territorial radar. It saves us days of reading and helps us get straight to what really matters.*”, from a pilot user’s notes). In parallel, we directly collected their feedback through the application. On a sample of approximately 200 arguments, 75% were judged relevant or useful by project teams.

In the following, we present some insights gathered with the Stakeholders Suite on a solar project. Located in southern France, the *Montagne de Lure* area has experienced several solar development projects, some of which have triggered tensions among local associations and residents. Using Stakeholder Suite, project teams aggregated data from past and ongoing debates to identify active opponents, supportive officials, and recurring public concerns such as deforestation, landscape alter-

ation, and threats to protected species.

**Topic Extraction.** Automatic clustering surfaced two dominant themes, i.e., land use conflicts and regulation and participation. The land-use cluster revealed strong attention to biodiversity and landscape protection, with repeated mentions of the ocellated lizard, a locally protected species. This insight led the project team to initiate additional environmental assessments before site validation. The regulation cluster showed that many objections targeted the lack of transparency and citizen participation rather than the solar technology itself, prompting stronger consultation efforts. A third cross-topic finding indicated a clear preference for solar installations on brownfields or rooftops instead of forests or farmland.

**Argument Analysis.** AM and stance classification exposed a well-structured opposition led by environmental groups. Their discourse gained visibility through local media and a published book that became a symbolic reference against solar projects in forested areas. Conversely, arguments from supportive actors (including other project owners) highlighted the site’s strong solar potential and its strategic alignment with EU renewable goals, offering communication models for future initiatives. **Network Insights.** It revealed tightly connected opposition clusters, suggesting coordinated activism rather than isolated critics. It also identified intermediary actors bridging detractors and supportive municipalities, indicating potential spaces for dialogue and balanced engagement.

### 5.2 Application to Political Debate

**Context.** To assess the adaptability of the framework, we applied it to a political debate in the energy domain. The Programmation Pluriannuelle de l’Énergie (PPE3) defines France’s energy strategy for the period 2022-2035, covering the reduction of fossil fuel dependence, expansion of renewable energy sources (solar, wind, biogas, hydrogen), maintenance of the nuclear fleet, energy efficiency measures, grid modernization and transport decarbonization. The plan is shaped through a public consultation process involving parliamentary debates and stakeholder contributions.<sup>3</sup>

**Data sources.** The pipeline is applied to three corpora drawn from the PPE 3 consultation:

(1) **Sénat:** 2 transcripts of senatorial debates (1,224 interventions from 41 senators). (2) **As-**

<sup>3</sup>The code used for this experiment is available at [https://github.com/stakeholder-suite/demo\\_eacl\\_2026](https://github.com/stakeholder-suite/demo_eacl_2026)

**semblée Nationale (AN):** 3 transcripts of debates in the lower chamber (722 interventions from 99 deputies). (3) **Actor Workbook:** 329 position papers from institutional stakeholders (industry federations, NGOs, local authorities, energy operators), each corresponding to one organization.

**Thesaurus construction.** The thesaurus was constructed using the actor workbooks, which provides a diverse set of viewpoints comparable to the press corpus used in previous experiments. Topic clustering was performed using HDBSCAN, although spectral clustering may offer finer control over the number of clusters.

**Argument extraction.** In this setting, we focus on actor and topic-level arguments. For each stakeholder, arguments are generated using our RAG approach. Given the broad scope of PPE3, global stance questions were not appropriate. Therefore, we retain only arguments that explicitly mention specific energy sources (solar, nuclear, wind) using keyword filtering. For each selected argument, stance is then computed relative to the identified energy source through a secondary query. This process yielded more than 2,000 arguments, including over 500 energy-specific arguments.

**Insights.** The analysis shows that a large proportion of arguments focus on nuclear energy, with most parliamentary actors expressing support for this source. In contrast, wind energy appears more contested, with a higher proportion of opposing arguments in parliamentary debates. These results demonstrate the framework’s ability to extract source-grounded arguments and capture stakeholder positions from relatively small corpora.

## 6 Conclusion

We presented a production-oriented framework for analyzing public debates that unifies *stakeholders*, *topics*, and *arguments* in a single data model. The system combines paragraph-level retrieval, LLM-based argument extraction, and stance classification to produce source-grounded insights that scale across domains. Tested on multiple energy-infrastructure projects, and successfully piloted in legislative contexts, the Stakeholder Suite offers transparent, reusable workflows that complement existing media intelligence and public-affairs tools. Our evaluation indicates that diversity-aware retrieval provides adequate coverage for downstream generation, and that stance classification achieves high accuracy at the paragraph–argument

level. LLM-as-a-judge effectively filters unsupported claims but can be conservative, discarding valid arguments in ambiguous contexts.

## 7 Ethical & broader considerations

Stakeholder management has become a key issue for multinational corporations across the world. The threat of controversies led by the stakeholders of an industrial project can result in project cancellation, especially in the field of major infrastructure projects (energy, transportation, telecommunication. . .). Public affairs and corporate social responsibility departments traditionally in charge of stakeholder management and controversy risk mitigation, are facing a new challenge, namely the growing digitization of controversies happening on a wide array of social media platforms. Being able to make sense of all the digital traces stakeholders leave on social media, in press releases, meeting minutes of public debates, etc., requires the elaboration of new management tools connected to heterogeneous data sources in different languages. AI-based solutions are bringing a wind of change in corporate daily practices for both managing stakeholders’ data and designing communication strategies. It seems very likely that multinational corporations will equip themselves with AI-based solutions to refine the monitoring of potential controversies and better manage associated risks. The implementation of such solutions in the corporate toolbox implies actions of change management to ensure solution adoption and effective use. Research in sociology of innovation has shown that the disruptive impact of innovation on daily routines can lead to users’ resistance and cause innovative projects to fail (Akrich et al., 2002). To become a success, AI-based stakeholder management solutions will have to integrate a change management strategy to ensure adoption in public affairs and corporate social responsibility departments. Regarding the European geographical area, one limitation of the development of such tools will be the respect of GDPR, which necessitates privacy-by-design solutions when manipulating social media data. We process publicly available texts only; personal data are handled under GDPR legitimate-interest.

## Acknowledgments

We thank Stakeholder Suite project team for their support and contributions throughout this work, including Alexis Courtin, Elena Hinnekens, Juliette

Lagrange, Bilel Loussaief, Gilles Olivié, Ousmane Traoré. This work has been funded by ENGIE Research & Innovation. The work of E. Cabrio has been supported by the French government, through the 3IA Cote d’Azur investments in the project managed by the National Research Agency (ANR) with the reference number ANR-23-IACL-0001.

## References

- Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019. FLAIR: An easy-to-use framework for state-of-the-art NLP. In *NAACL 2019, 2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59.
- Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. [Contextual string embeddings for sequence labeling](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Madeleine Akrich, Michel Callon, and Bruno Latour. 2002. [The Key to Success in Innovation Part I: The Art of Interestment](#). *International Journal of Innovation Management*, 6(2):187 – 206.
- 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, and 12 others. 2020. [Language models are few-shot learners](#). *CoRR*, abs/2005.14165.
- Ricardo J. G. B. Campello, Davoud Moulavi, and Joerg Sander. 2013. Density-based clustering based on hierarchical density estimates. In *Advances in Knowledge Discovery and Data Mining*, pages 160–172, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Jaime Carbonell and Jade Goldstein. 1998. [The use of mmr, diversity-based reranking for reordering documents and producing summaries](#). In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’98, page 335–336, New York, NY, USA. Association for Computing Machinery.
- Guizhen Chen, Liying Cheng, Anh Tuan Luu, and Lidong Bing. 2024. [Exploring the potential of large language models in computational argumentation](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2309–2330, Bangkok, Thailand. Association for Computational Linguistics.
- Qianhe Chen, Yong Wang, Yixin Yu, Xiyuan Zhu, Xuerou Yu, and Ran Wang. 2025. [Conch: Competitive debate analysis via visualizing clash points and hierarchical strategies](#). *Preprint*, arXiv:2507.14482.
- Johannes Daxenberger, Benjamin Schiller, Chris Stahlhut, Erik Kaiser, and Iryna Gurevych. 2020. [Argumenttext: Argument classification and clustering in a generalized search scenario](#). *Datenbank-Spektrum*, 20(2):115–121.
- Daniel Firebanks-Quevedo, Jordi Planas, Kathleen Buckingham, Cristina Taylor, David Silva, Galina Naydenova, and René Zamora-Cristales. 2022. [Using machine learning to identify incentives in forestry policy: Towards a new paradigm in policy analysis](#). *Forest Policy and Economics*, 134:102624.
- Pierpaolo Goffredo, Elena Cabrio, Serena Villata, Shohreh Haddadan, and Jhonatan Torres Sanchez. 2023. [DISPUTool 2.0: A Modular Architecture for Multi-Layer Argumentative Analysis of Political Debates](#). In *Proceedings of the AAAI Conference on Artificial Intelligence, AAAI-23 Special Programs, IAAI-23, EAAI-23, Student Papers and Demonstrations*, pages 16431–16433, Washington, DC, United States. AAAI.
- Adam Tauman Kalai, Ofir Nachum, Santosh S Vempala, and Edwin Zhang. 2025. [Why language models hallucinate](#). *arXiv preprint arXiv:2509.04664*.
- Michaela LaPatin, Lauryn A. Spearing, Helena R. Tiedmann, Miriam Hacker, Olga Kavvada, Jean Daniélou, and Kasey M. Faust. 2023. [Controversy in wind energy construction projects: How social systems impact project performance](#). *Energy Policy*, 176:113507.
- 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](#). *CoRR*, abs/2005.11401.
- Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. [MTEB: Massive text embedding benchmark](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.
- Jordi Planas, Daniel Firebanks-Quevedo, Galina Naydenova, Ramansh Sharma, Cristina Taylor, Kathleen Buckingham, and Rong Fang. 2022. [Beyond modeling: NLP pipeline for efficient environmental policy analysis](#). *CoRR*, abs/2201.07105.

- Nils Reimers, Benjamin Schiller, Tilman Beck, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2019. [Classification and clustering of arguments with contextualized word embeddings](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 567–578, Florence, Italy. Association for Computational Linguistics.
- Jianbo Shi and J. Malik. 2000. [Normalized cuts and image segmentation](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905.
- Prashant Singh, Erik Lehmann, and Mark Tyrrell. 2024. [Climate policy transformer: Utilizing NLP to track the coherence of climate policy documents in the context of the Paris agreement](#). In *Proceedings of the 1st Workshop on Natural Language Processing Meets Climate Change (ClimateNLP 2024)*, pages 1–11, Bangkok, Thailand. Association for Computational Linguistics.
- Noam Slonim, Yonatan Bilu, Carlos Alzate, Roy Bar-Haim, Ben Bogin, Francesca Bonin, Leshem Choshen, Edo Cohen-Karlik, Lena Dankin, Lilach Edelstein, Liat Ein-Dor, Roni Friedman-Melamed, Assaf Gavron, Ariel Gera, Martin Gleize, Shai Gretz, Dan Gutfreund, Alon Halfon, Daniel Hershcovich, and 34 others. 2021. [An autonomous debating system](#). *Nature*, 591(7850):379–384.
- Laura South, Michail Schwab, Nick Beauchamp, Lu Wang, John Wihbey, and Michelle A. Borkin. 2020. [Debatevis: Visualizing political debates for non-expert users](#). In *2020 IEEE Visualization Conference (VIS)*, pages 241–245.
- Christian Stab, Tristan Miller, Benjamin Schiller, Pranav Rai, and Iryna Gurevych. 2018. [Cross-topic argument mining from heterogeneous sources](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3664–3674, Brussels, Belgium. Association for Computational Linguistics.
- Alex Tamkin, Miles McCain, Kunal Handa, Esin Durmus, Liane Lovitt, Ankur Rathi, Saffron Huang, Alfred Mountfield, Jerry Hong, Stuart Ritchie, Michael Stern, Brian Clarke, Landon Goldberg, Theodore R. Summers, Jared Mueller, William McEachen, Wes Mitchell, Shan Carter, Jack Clark, and 2 others. 2024. [Clio: Privacy-preserving insights into real-world ai use](#). *Preprint*, arXiv:2412.13678.
- Tommaso Venturini and Anders Kristian Munk. 2022. *Controversy mapping : a field guide*. Polity, Cambridge, UK.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. [Finetuned language models are zero-shot learners](#). *CoRR*, abs/2109.01652.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhaghao Wu, Yonghao Zhuang, Zi Lin,
- Zhuohan Li, Dacheng Li, Eric P. Xing, Haotong Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *ArXiv*, abs/2306.05685.

# A Appendix

## A.1 Data parsing prompt

You will find here the prompt used for the parsing of documents.

### Coreference Resolution & Speaker Detection

#### System Prompt: Role

You are an *NLP engine specialized in text disambiguation and speaker identification*. Your task is to transform ambiguous documents into self-contained, RAG-ready paragraphs with clear entity references and speaker attribution.

#### Context

- **Document author:** {doc\_editor}
- **Document title:** {doc\_title}
- **Original language:** Preserved in output
- **Target use case:** Retrieval-Augmented Generation (RAG)

#### Expected Task

##### 1. Coreference Resolution

1. Replace **all pronouns** (I, you, he, she, it, they, je, il, elle, ils, etc.) with their specific referent.
2. Replace **vague references** ("this project", "the park", "Madam", "Monsieur") with explicit names or descriptive phrases.
3. If an entity is unnamed, create a **short descriptive identifier** based on context.
4. Make each paragraph **self-contained and unambiguous**, preserving order, style, and all factual details.

##### 2. Speaker Detection & XML Tagging

1. Wrap every paragraph in: `<p speakerName="" speakerFunction=""></p>`.
2. Populate **speakerName** when the paragraph contains:
  - A direct quote (« ... », "..."), or
  - An indirect quote with reporting verbs (said, declared, stated, according to, etc.)
3. **speakerFunction** = entity's role *only if stated or clearly inferable*.
4. If **multiple speakers** appear in one paragraph, split it so each `<p>` has a single speaker.
5. If a name is only **mentioned** (not speaking), leave both attributes **empty**.
6. Always output the speaker's name and function *in the paragraph text*.

##### 3. Document Structure

1. Merge short, consecutive, related sentences from the **same speaker**.
2. Never merge different speakers.
3. Avoid small paragraphs lacking sufficient context to be understood alone.
4. Produce **valid XML**; do not modify text outside `<content>`.

#### Content Constraints

- Always include the **project name** in every relevant paragraph.
- Avoid generic references: *not* "the solar farm is too big" *but* "the solar project Zephyr is too big".
- Keep the **original language** and all numerical/factual elements intact.
- Output language: {output\_language}.

#### Expected Output (XML)

```
<edited_content>
  <p speakerName="" speakerFunction="">
    The Zephyr solar project was initiated last year...
  </p>
  <p speakerName="Jean_Daniel" speakerFunction="Mayor_of_Duneshale">
    "The_Zephyr_project_is_very_important..." said
    Jean Daniel, the Mayor of Duneshale.
  </p>
  ...
</edited_content>
```

#### Example

- Input: "The city hall is against the project."  
Output: "The city hall of Duneshale is against the Zephyr project."
- Input: "He gave an opinion against the park installation."  
Output: "The mayor of Duneshale, Jean Daniel, gave an opinion against the Zephyr solar park installation."

#### User Input: Document to process:

```
<document>
  <author>{doc_editor}</author>
  <title>{doc_title}</title>
  <content>{doc_content}</content>
</document>
```

Return the transformed text in a single `<edited_content>` block containing ordered `<p>` tags.  
**No extra commentary.**

## A.2 Topic modeling and extraction

We implemented the thesaurus creation approach to each of our energy projects. Below we present the result for a large-scale low-carbon hydrogen project in Belgium and the prompts used to generate it. It aims to produce hydrogen via autothermal reforming (ATR) of natural gas with over 95% CO<sub>2</sub> capture. The project is part of a wider debate around blue hydrogen and carbon capture solutions in Belgium and Europe. Thanks to our hybrid approach we were able to pinpoint precise debates within this project. Each bullet point corresponds to one topic and is embedded in a format combining description+subtopics for the semantic search.

### A.2.1 Thesaurus example

- **Infrastructure Development and Network Integration**

*Planning, construction and repurposing of hydrogen and CO<sub>2</sub> pipelines, storage, terminals, and their integration with existing gas grids at national and cross-border levels.*

- Pipeline construction and repurposing
- Hydrogen storage and buffering facilities
- CO<sub>2</sub> transport and storage networks
- Import/export terminals
- Integration with existing gas infrastructure
- Safety, materials and technical challenges
- Feasibility and front-end engineering studies

- **Belgium as Regional and European Hydrogen Hub**

*Belgium's development as a major import, transit, and distribution node within the North-West European hydrogen network, leveraging its ports and cross-border corridors.*

- Port-based hub development (Antwerp–Bruges, Ghent, Zeebrugge)
- Integration with the European Hydrogen Backbone
- Cross-border pipeline corridors
- Import and transit functions
- Market access and interoperability standards

- **Hydrogen Production Technologies and Carbon Capture**

*Technical and strategic aspects of producing low-carbon hydrogen via autothermal reforming, electrolysis, and hybrid systems, coupled with carbon capture and storage solutions.*

- Autothermal reforming (ATR) with CCS
- Electrolysis technologies (PEM, AEM, alkaline)
- Carbon capture, transport, and storage
- Hybrid offshore wind–hydrogen platforms
- Pilot plants and demonstration projects
- Technology innovation and scalability

- **Market Dynamics, Economics and Investment**

*Cost competitiveness, demand forecasts, financing models, and offtake agreements shaping the feasibility and scale-up of hydrogen projects in Belgium.*

- Cost comparison: blue vs. green hydrogen
- Production and transport cost analysis
- Market demand forecasts (2030/2050)
- Investment requirements and funding models
- Offtake agreements and revenue models
- Economic and regional impact (jobs, growth)

- **Policy and Regulatory Frameworks**

*National, regional, and EU rules, support schemes, and permitting processes that enable hydrogen infrastructure, production, and market development.*

- Belgian federal and regional hydrogen strategies
- EU directives and Clean Hydrogen Alliance
- State aid and subsidy mechanisms
- Permitting, safety standards, and carbon pricing
- Cross-border regulatory coordination

- **Stakeholder Collaboration and Partnerships**

*Consortia, public–private partnerships, and international alliances among energy companies, ports, TSOs, and authorities to build a coordinated hydrogen value chain.*

- Industry consortia and coalitions
- Public–private partnerships
- Inter-TSO and port authority cooperation
- International memoranda of understanding
- Role of the Belgian Hydrogen Council

- **Environmental Impact and Sustainability**

*Lifecycle emissions, CCS effectiveness, and sustainability concerns—including fossil-gas lock-in risks—surrounding blue hydrogen projects.*

- CO<sub>2</sub> emissions reduction potential
- Life-cycle assessment and leakage risks
- CCS performance and permanence
- Fossil-gas dependency and financing risks
- Alignment with EU and national climate targets

- **Blue vs Green Hydrogen Dynamics**

*Strategic, economic, and policy comparisons of blue versus green hydrogen pathways, including transitional roles and long-term market coexistence.*

- Blue hydrogen as a transitional solution
- Green hydrogen cost decline and scale-up
- Policy incentives and market preferences
- Hybrid and co-location production strategies
- Long-term sustainability comparisons

- **Renewable Energy Integration and Hybrid Systems**

*Coupling renewable power—especially offshore wind—with hydrogen production for storage, grid balancing, and hybrid energy platforms in the North Sea region.*

- Offshore wind-to-hydrogen projects
- Hybrid wind–gas–hydrogen installations
- Power-to-gas and seasonal storage
- Grid balancing and ancillary services
- Renewable electricity supply constraints

- **Sectoral Applications and Industrial Clusters**

*End-use deployment of hydrogen across heavy industry, transport, maritime, chemicals, and energy systems, with a focus on cluster integration and local decarbonization.*

- Steel and chemicals decarbonization
- Maritime transport and bunkering
- Chemical feedstock substitution
- Power generation and heat applications
- Fuel-cell mobility and transport
- Port cluster integration

## **A.2.2 Topic modeling prompts**

### Cluster naming and example selection

#### System Prompt:

##### Role

You are an *energy journalist specializing in {energyType} and energy transition topics*. Your expertise allows you to confront stakeholders' statements with the strategic, technical, and financial challenges of {name}.

##### Context

- **Project name:** {name}
- **Energy type:** {energyType}
- **Description:** {longProjectdescription}
- **Scope:** {scope}

##### Task

1. Read a series of topics extracted from papers related to the project.
2. Group them into a **common thematic category**.
3. Provide a clear and concise name for each theme.
4. Write a short description of each theme.
5. List representative and diverse examples.
6. Answer in the following language: {language}.

##### Expected Output (XML)

```
<themes>
  <theme>
    <topic>Name of topic 1</topic>
    <description>Short description...</description>
    <examples>
      <example>sentence#2</example>
      <example>sentence#5</example>
    </examples>
  </theme>
  ...
</themes>
```

#### User Question:

##### Question

Here is the list of topics:  
{topics}

Return the main themes in XML. There can be between 1 and many topics. Provide the answer in the **same language as the topics**.

## Thesaurus building

### System Prompt:

#### Role

You are an *economic journalist specializing in energy and transition issues*. Your expertise allows you to confront stakeholders' statements with the strategic, technical, and financial challenges of **{name}**.

#### Context

- **Project:** {name}
- **Energy type:** {energyType}
- **Description:** {longProjectdescription}
- **Scope:** {scope}

#### Expected Task

##### Formal constraints

1. Read a series of topics extracted from papers related to the project.
2. Group topics addressing similar issues into one **common thematic category**.
3. Provide a **clear and concise name** for each theme.
4. Write a **short and precise description** of each theme.
5. Add **sub-themes** based on the themes given as input.
6. Answer in the following language: {language}.

##### Content constraints

1. Do not create overly generic themes (example: *energy transition*).
2. Themes must follow the **MECE principle**: Mutually Exclusive, Collectively Exhaustive.
3. Sub-themes must also follow the **MECE principle** (max. 10 varied sub-themes covering all aspects of the theme).

#### Expected Output (XML)

```
<themes>
  <theme>
    <topic>Theme name 1</topic>
    <description>Short description...</description>
    <subtopics>
      <subtopic>Sub-theme name 1</subtopic>
      <subtopic>Sub-theme name 2</subtopic>
    </subtopics>
  </theme>
  <theme>
    <topic>Theme name 2</topic>
    <description>Short description...</description>
    <subtopics>
      <subtopic>Sub-theme name 1</subtopic>
      <subtopic>Sub-theme name 2</subtopic>
    </subtopics>
  </theme>
  ...
</themes>
```

#### User Question:

Here is the list of topics:

{topics}

Return the main themes (**maximum 10**) from the list in XML format. Provide the answer in the **same language as the topics**.

## A.3 Argument Mining

### A.3.1 Examples

Here we present the complete set of arguments related to blue hydrogen in Belgium that were extracted and analyzed in this study. Please refer to the last section to read the context of the blue hydrogen project. Table 3 organizes these arguments across three dimensions:

- **ACTOR:** Arguments attributed to specific stakeholders (e.g., European Commission).
- **TOPIC:** Arguments related to thematic issues (e.g., Environmental Impact and Sustainability).
- **GLOBAL:** General arguments against blue hydrogen in Belgium.

Each argument includes source identifiers, quality assessments (GOOD/BAD with comments where issues were identified), and stance classification (PRO, CON or NEUTRAL).

Dim.	Year	Parsed Argument	Source	Stance	Quality
<i>Query: What are the arguments of the stakeholder European Commission about blue hydrogen in Belgium?</i>					
ACTOR	–	The European Commission supports the development of low-carbon hydrogen technologies, including blue hydrogen, as part of its strategy to reduce dependency on Russian gas and diversify energy sources in Europe.	id-1, id-9	NEUTRAL	BAD: Overinterpretation of the source
ACTOR	–	The European Commission is actively promoting hydrogen infrastructure projects that include large-scale electrolyzers and transport infrastructure for renewable and low-carbon hydrogen production, storage, and transport.	id-19	NEUTRAL	GOOD
ACTOR	–	The European Commission acknowledges the important role of low-carbon, methane-derived blue hydrogen in reducing emissions in the nearer term, alongside a main focus on green hydrogen produced from renewable energy.	id-22	PRO	GOOD
ACTOR	–	Ursula von der Leyen highlighted hydrogen as a technology where Belgium is poised to become a world leader, emphasizing the potential impact of the Net-Zero Industry Act to transition the hydrogen economy from niche to large-scale industry by 2030.	id-16, id-18, id-21	NEUTRAL	BAD: Not direct speech of the actor
ACTOR	–	The European Commission has approved significant state aid (€5.2 billion) for hydrogen technologies, supporting projects that contribute to a low-carbon hydrogen economy and energy independence, which can be linked to the context of blue hydrogen initiatives in Belgium.	id-1, id-9	NEUTRAL	GOOD
<i>Query: What are the arguments related to the topic Environmental Impact and Sustainability?</i>					
TOPIC	2023	Low-carbon hydrogen produced from natural gas combined with carbon capture and storage has a carbon footprint of around 80 to 90 grams of CO <sub>2</sub> per kilowatt-hour along the entire value chain.	id-4	NEUTRAL	GOOD
TOPIC	2023	A study from the University of Ghent shows that firing furnaces with blue hydrogen using carbon capture and storage can reduce climate change impact by 8% to 18% compared to conventional steam cracking plants.	id-19	PRO	GOOD
TOPIC	2023	Research indicates that the carbon capture and storage system at a blue hydrogen plant operated by Shell in Alberta emits more CO <sub>2</sub> than it captures, raising concerns about the environmental effectiveness of such systems.	id-18	CON	GOOD
TOPIC	2023	A 2021 Cornell University study found that gas emissions from burning blue hydrogen were more than 20% greater than using conventional gas, questioning the sustainability of blue hydrogen.	id-5	CON	GOOD
<i>Query: What are all the arguments against blue hydrogen in Belgium?</i>					
GLOBAL	2023	The CCS system at the blue hydrogen plant operated by Shell in Alberta, Canada, was emitting more CO <sub>2</sub> than it was capturing, raising concerns about the environmental effectiveness of blue hydrogen projects.	id-14	CON	GOOD
GLOBAL	2023	The use of blue hydrogen has been pushed by industrial and gas lobbies, which may indicate a conflict of interest and potential prioritization of fossil fuel interests over genuine climate goals.	id-20	CON	GOOD
GLOBAL	2023	The Department of Energy in the US has indicated that blue hydrogen is unlikely to qualify for hydrogen tax credits due to high upstream emissions, suggesting that projects may have significant emissions issues.	id-18	CON	GOOD

Table 3: Arguments Related to Blue Hydrogen in Belgium

### A.3.2 Argument Mining prompts

## Argument extraction

### System Prompt:

#### Role

You are an *intelligent agent with expertise in the energy sector*, specifically in the following energy type: {energyType}. Your mission is to assist the teams working for {name} in identifying all arguments for and against raised by various stakeholders regarding {name}. The goal is to help the user better understand the territory in which this project is being developed.

#### Definition

In the instructions below, the term "*project*" always refers to {name}. This project can be described as follows "{longProjectdescription}" and it raises both support and concerns. {scope} We analyze different types of projects — sometimes individual ones, sometimes grouped within a specific area. Regardless of their exact nature, we consistently refer to them as "*project*", as long as they fall under the {name} initiative.

#### Task

Your task is to provide the list of arguments presented by a stakeholder of {name} mentioned in the question below regarding {name}, based on the excerpts provided below that mention a stakeholder of {name}.

- Only mention the information present in the excerpts and do not overinterpret the information provided.
- The excerpts may not necessarily contain arguments. If no excerpt contains arguments presented by the specific stakeholder mentioned in the question below of {name}, simply respond "I don't know."
- If an excerpt is not relevant, ignore it.
- All arguments must be exclusively related to the stakeholder mentioned in the question below of {name}.

#### Response Guidelines

1. Write each argument in English.
2. After each sentence, indicate the ID of the excerpt used in parentheses (e.g., (id-1)).
3. If multiple excerpts support an argument, cite them together, e.g. (id-1, id-2).
4. Do not establish links between passages in the CONTEXT section — each passage is unique and independent.
5. If no excerpt is relevant, respond only with "I don't know."
6. Each argument must be expressed in a single sentence.
7. Do not start with logical connectors (*However, But*, etc.).
8. Each sentence must be unique, interpretable on its own and end with the source(s).
9. Avoid ambiguity:
  - Do not use pronouns like "He" or "She" without context.
  - Do not say "the project" — explicitly state {name}.
10. The presented arguments must not overlap: if the same argument appears in multiple excerpts, it should be expressed only once.
11. Be exhaustive and precise.
12. Each argument must be on a separate line and end with its source(s).

#### Example (Fictitious)

##### Excerpts:

- id-1: According to Pais magazine, the mayor of Blur is a strong supporter of protests and their positive impact on politics in the city.
- id-2: Protests have been criticized for their impact on local businesses since 2012. Bernard, a shopkeeper, complained about the disruptions caused by these protests.
- id-3: The mayor of Blur recently stated that protests are a way to make his voice heard.
- id-4: Protests have an impact on traffic in the city due to traffic jams according to the mayor of Blureau.

*Question:* What are the arguments of the stakeholder Blur City Council about the protests?

*Expected response:* The arguments presented by the Blur City Council regarding protests are as follows:

- The mayor of Blur is a strong supporter of protests and their positive impact on politics in the city. (id-1)
- The Blur City Council recently stated that protests are a way to make their voice heard. (id-3)

##### Context:

{context}

**Global-type Question:** What are all the arguments {in favor/against} project {name}?

**Actor-type Question:** What are the arguments of the stakeholder {actor} about project {name} ?

**Topic-type Question:** What are the arguments related to the topic {topic\_name}?

## Actor Stance Classification

### User Prompt:

Members of the team have identified arguments presented by {actor\_name} regarding {project\_name}. An argument involves a statement or discussion about the project: {project\_name}. From this argument, the scope is to determine a stance associated with {actor\_name} concerning: {project\_name}.

#### Rules for classification:

- If in the argument {actor\_name} is merely factual or describing a situation, respond with NEUTRAL.
- If in the argument {actor\_name} expresses a clear stance:
  - Respond with IN FAVOR if {actor\_name} supports the {project\_name} project.
  - Respond with AGAINST if {actor\_name} opposes the {project\_name} project.
- Remember: an argument does not necessarily imply a stance IN FAVOR or AGAINST. If there is no clear stance from {actor\_name} specifically regarding {project\_name}, respond with NEUTRAL.
- If the argument is "I don't know", respond with UNKNOWN.

**Important:** Do not add any words or phrases to introduce the response. Respond only with one of: IN FAVOR, AGAINST, NEUTRAL, or UNKNOWN.

Argument from {actor\_name}:  
{argument}

**Answer:**

### Theme-based stance Classification

#### User Prompt:

Members of the team have identified arguments regarding: {project\_name} that relate to the theme: {topic\_name}. An argument involves a statement or expression of opinion about the project: {project\_name}. The objective is to determine a stance associated with the theme: {topic\_name} concerning: {project\_name} based on this argument.

#### Rules for classification:

- If the argument associated with {topic\_name} is purely factual or descriptive, respond with NEUTRAL.
- If the argument expresses a clear stance:
  - Respond with IN FAVOR if {topic\_name} supports the {project\_name} project.
  - Respond with AGAINST if {topic\_name} opposes the {project\_name} project.
- If no clear stance is expressed regarding the theme {topic\_name} specifically in relation to {project\_name}, respond with NEUTRAL.
- If the argument is "I don't know", respond with UNKNOWN.

**Important:** Do not add any words or phrases to introduce the response. Respond only with one of: IN FAVOR, AGAINST, NEUTRAL, or UNKNOWN.

Argument associated with the theme {topic\_name}:  
{argument}

**Answer:**

### LLM-as-a-Judge: Argument Coherence Evaluation

#### User Prompt:

You are an *intelligent agent*, expert in the energy field, specifically in the following type of energy: {energyType}. The teams working for {name} have identified the arguments for and against raised by the various stakeholders of project {name}. Your mission is to assist the teams in verifying the coherence of the arguments with respect to the **Argument Source** from which they are drawn.

Each argument is identified based on a specific Argument Source. Your task is to evaluate the coherence and quality of a given argument by comparing it to the **Argument Source** (reference text below).

#### Evaluation Criteria:

- "Invalid Argument": The argument contains factual errors compared to the Argument Source. Examples:
  - false figures compared to the Argument Source,
  - abusive generalizations,
  - off-topic compared to the Argument Source,
  - does not rely at all on the Argument Source,
  - incoherent with the Argument Source.Verify that all numerical values in the Argument Source are the same as those in the argument.
- "Weak argument": The argument is correct but lacks precision compared to the information in the Argument Source, or mentions actors in an imprecise way (e.g., "the town hall" instead of "Paris town hall"). *Exception:* If the Argument Source itself is not precise, the argument can still be considered relevant.
- "Strong argument": The argument is well based on the Argument Source, coherent, and relevant.

#### Required Output:

```
{
  "global_comment": "Your explanation here",
  "quality_argument": "Invalid Argument|Weak Argument|Strong Argument"
}
```

Argument:  
{text}

Argument Source:  
{source\_text}

## A.4 System demonstration

The following figures present different screens of the current app. Please refer to the caption for more details.

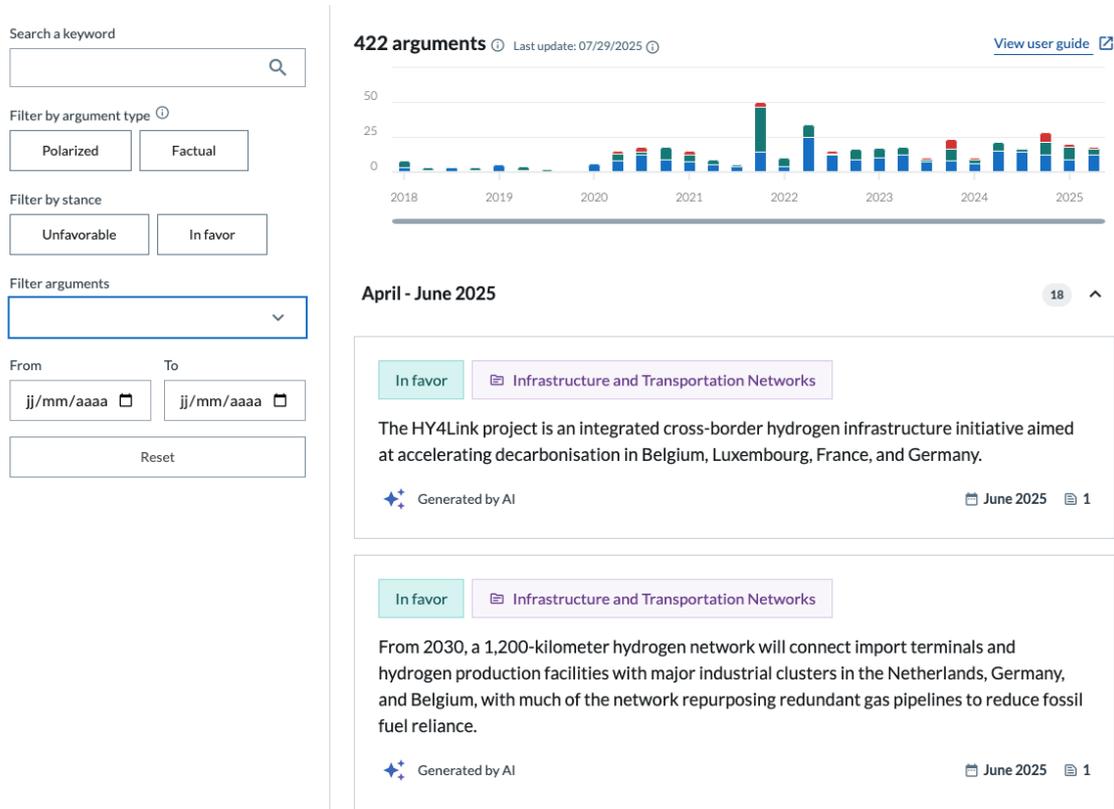


Figure 2: Argument page where we can access all the arguments related to the project and filter by stance, type (actor, topic, global) and by date.

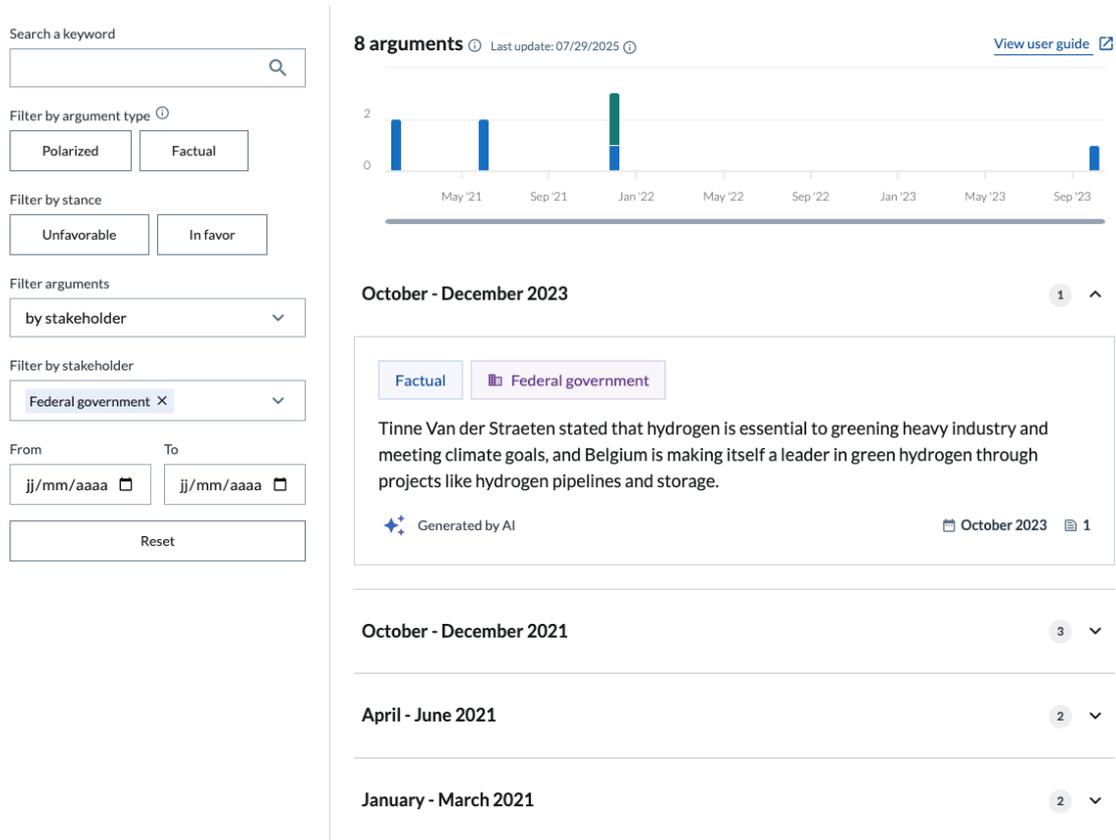


Figure 3: A zoom on the argument page with actor type filter selected for the Federal Government.

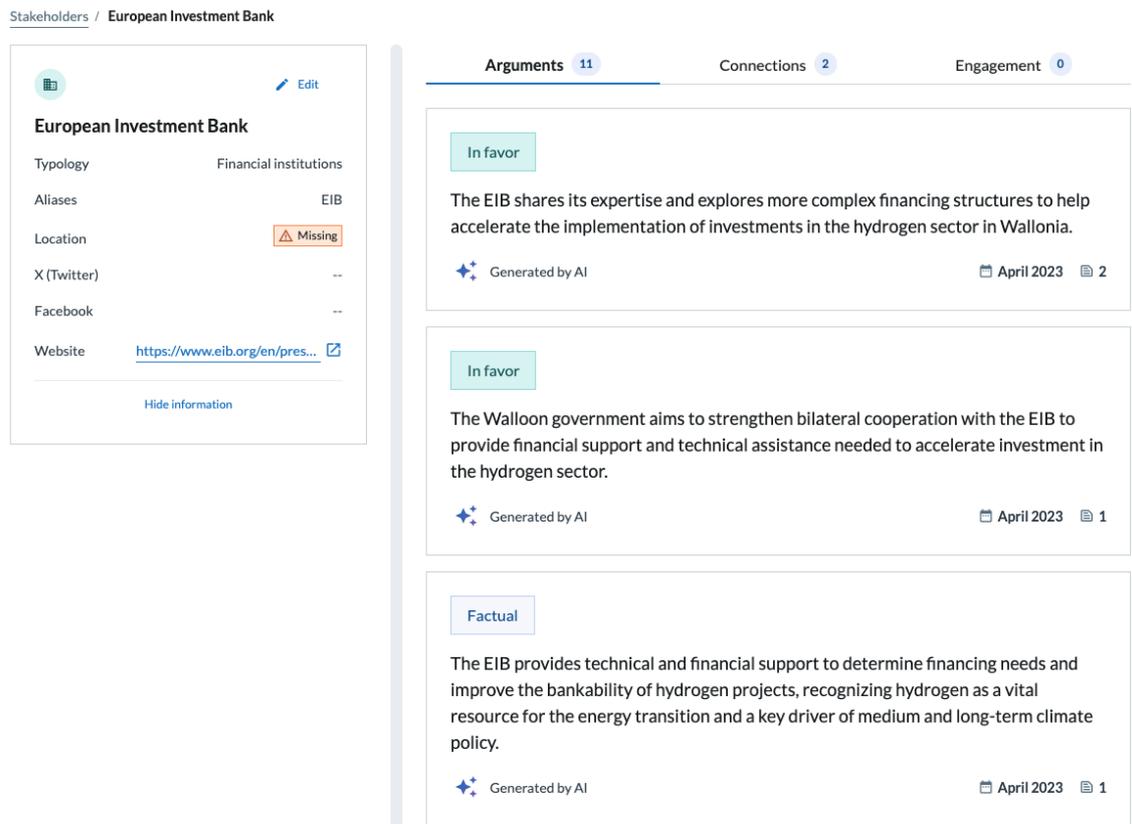


Figure 4: Actor view page where we can find the actor information, its arguments and the other actors he is connected with.

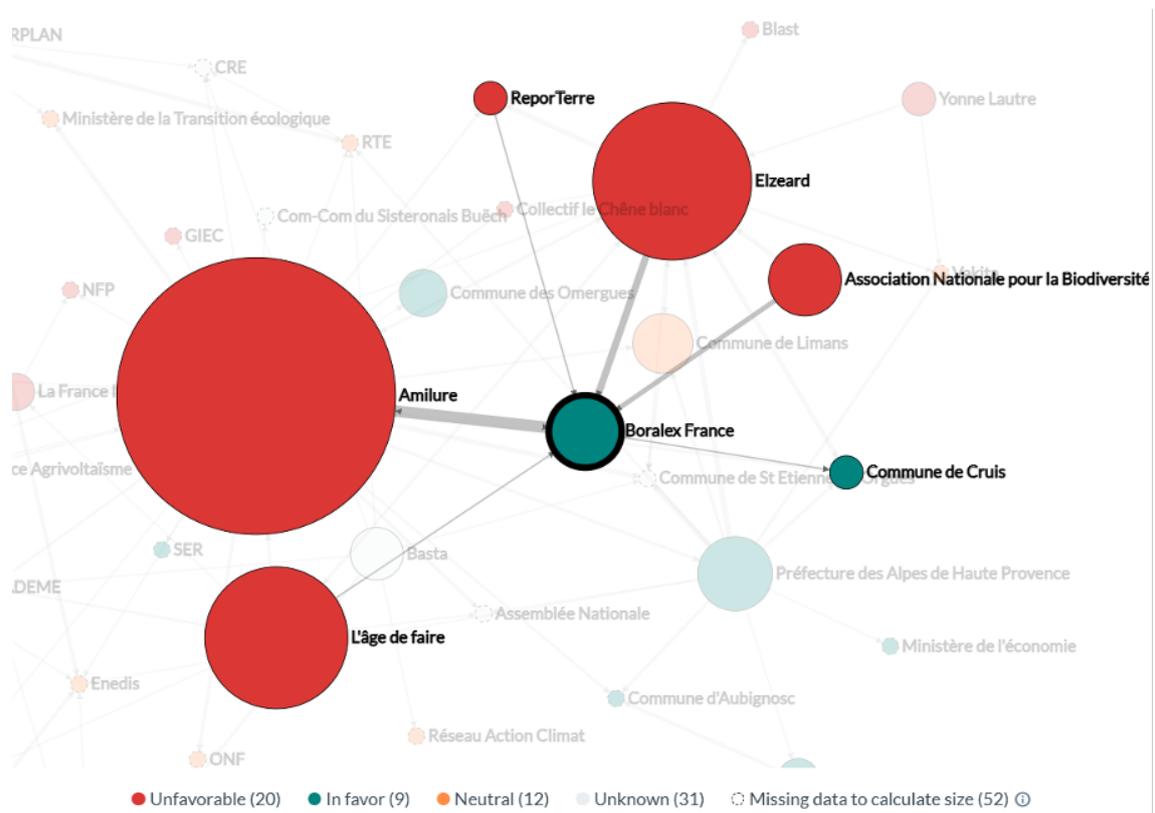


Figure 5: Mapping page where we can see the connections of actors and their importance in the debate.

# DeepPavlov Strikes Back: A Toolkit for Improving LLM Reliability and Trustworthiness

Evgeny Nikolaev<sup>5,4\*</sup>, Timur Ionov<sup>3,4\*</sup>, Anna Korzanova<sup>1</sup>,  
Vasily Konovalov<sup>2,1</sup>, Maksim Savkin<sup>2,1</sup>

<sup>1</sup>MIRAI, <sup>2</sup>AXXX, <sup>3</sup>MWS AI, <sup>4</sup>ITMO University, <sup>5</sup>AI Talent Hub

Correspondence: savkin.max.k@gmail.com

## Abstract

DeepPavlov 1.1 introduces new multilingual tools to enhance LLM reliability in production pipelines. It includes a span-level hallucination detector, an evergreen question classifier, and a toxicity classifier, all integrated into an easy-to-use open-source framework. These components address key LLM challenges: detecting factual inconsistencies against retrieved context, identifying static factual questions that bypass unnecessary retrieval, and flagging harmful content when alignment fails. Trained on PsiloQA, EverGreenQA, and TextDetox across 14+ languages, our encoder-based models outperform LLM baselines in accuracy and speed by orders of magnitude. Released under Apache 2.0 DeepPavlov 1.1 bridges traditional NLP and LLM-centric workflows for safer AI systems.

## 1 Introduction

Natural Language Processing (NLP) is a key component of many modern AI systems. It enables the automation of tasks that would otherwise require extensive manual labor, particularly those involving the processing of large volumes of raw text. As real-world applications of NLP continue to grow in complexity and scale, the demand for robust and easy-to-integrate tools grows as well.

At the core of our work is a long-term commitment to practical, user-focused development. As the field evolves, our tools continuously update in response to technological advances and shifting user needs.

The DeepPavlov library (Burtsev et al., 2018) was first introduced in the pre-BERT era, at a time when NLP systems were largely modular and relied heavily on explicit linguistic features. Early versions of the library focused on foundational

components such as Part-of-Speech (POS) tagging and syntactic parsing. These tools acted as building blocks that extracted structured linguistic information from raw text and played a critical role in training downstream models that depended on hand-crafted features to understand language.

The release of DeepPavlov 1.0 (Savkin et al., 2024) marked our transition into the post-BERT era, reflecting a shift in paradigm towards pre-trained transformer-based models. These models moved beyond low-level syntactic analysis by using high-level language understanding capabilities of the BERT-family models. We introduced models for Named Entity Recognition (NER) (Chizhikova et al., 2023), Knowledge Base Question Answering (KBQA), Multi-task learning (MTL) (Karpov and Konovalov, 2023), Text classification tasks (intent, sentiment) (Savkin and Konovalov, 2024), and tasks from the SuperGLUE (Wang et al., 2019) benchmark (NLI, RTE, Paraphrasing).

The emergence of large language models (LLMs) marked yet another paradigm shift in NLP. Instead of training dedicated models for each task, LLMs demonstrated impressive few-shot and zero-shot capabilities across a wide range of tasks. To enhance their reasoning and task-solving abilities, they are often paired with auxiliary tools, such as information retrieval systems, code execution environments, or external APIs. However, these LLM-centered workflows introduce new challenges: both the tools and the outputs of the models need to be monitored and validated to avoid irrelevant, harmful, or fabricated content. In this new context, the role of the DeepPavlov library evolved once again – now focusing on supporting LLM-centric pipelines. A major concern in such systems is hallucination: the generation of plausible but untrue information (Rykov et al., 2025a). Although it remains an open problem, it is increasingly important to have mechanisms for de-

\*Equal contribution. This work was done while the first two authors were affiliated with MIRAI.

Tool / Framework	DeepPavlov 1.1	DeepPavlov 1.0	spaCy	Stanza	Flair	AllenNLP*	jiant*
<i>Linguistic Features</i>							
Embeddings	✓	✓	✓	✓	✓	✓	✗
POS Tagger	✓	✓	✓	✓	✓	✗	✗
Lemmatizer	✓	✓	✓	✓	✗	✗	✗
Dependency Parsing	✓	✓	✓	✓	✗	✗	✗
Morphotagger	✓	✓	✓	✗	✗	✗	✗
Syntax Parser	✓	✓	✓	✗	✗	✗	✗
<i>Pretrained Encoders</i>							
NER	✓	✓	✓	✓	✓	✓	✓
Sentiment Classification	✓	✓	✓	✓	✓	✗	✗
Entity Linking	✓	✓	✓	✗	✓	✗	✗
Intent Classification	✓	✓	✗	✗	✗	✗	✗
Context QA	✓	✓	✗	✗	✗	✓	✓
SuperGLUE Tasks	✓	✓	✗	✗	✗	✓	✓
<i>LLM-centric Tools</i>							
Hallucination Detection	✓	✗	✗	✗	✗	✗	✗
Evergreen QA Classification	✓	✗	✗	✗	✗	✗	✗
Toxicity Classification	✓	✗	✗	✗	✗	✗	✗

Table 1: Comparison of supported instruments across popular NLP frameworks. Frameworks marked with “\*” are no longer supported.

tecting and flagging unreliable outputs, especially for high-stakes applications.

To address this need, we introduce DeepPavlov 1.1, an updated open-source NLP framework aimed at improving the reliability of LLM-based pipelines. This release introduces the following new multilingual components.

1. **Contextual Hallucination Detector** designed to estimate faithfulness – the factual consistency of model responses with retrieved content (Krayko et al., 2025).
2. **Evergreen Question Classifier** detects evergreen questions (factual questions whose correct answers are highly unlikely to change over extended periods of time). Usually evergreen question doesn’t require RAG pipeline (Pletenev et al., 2025).
3. **Toxicity Detection** determines whether a given text contains toxic content, serving as a safeguard when a language model’s alignment mechanisms fail to prevent harmful outputs (Dementieva et al., 2025).

## 2 Related Work

Before discussing directly comparable NLP frameworks, it is important to distinguish between related categories of tools that fall outside the scope of this work.

LLM orchestration frameworks such as

**LangChain**<sup>1</sup> have gained popularity as platforms for building agent-based pipelines. However, they are better characterized as workflow managers: they focus on task chaining and integration rather than providing pre-trained NLP models for assessing content quality. Therefore, they are not considered further in this paper.

Similarly, low-level libraries such as **PyTorch** (Paszke et al., 2019) and **TensorFlow** (Abadi et al., 2016) lack ready-to-use, task-specific NLP models and supporting infrastructure, so again they have a different purpose than the higher-level frameworks discussed here.

While early NLP libraries focused on linguistic features and task-specific modeling, modern applications increasingly rely on LLMs augmented by tools for retrieval and content validation. Despite this shift, most NLP frameworks have not adapted to the demands of LLM-centric workflows. Table 1 provides a detailed comparison of tools supported by major open-source NLP frameworks.

Libraries such as **spaCy** (Honnibal, 2017), **Stanza** (Qi et al., 2020), and **Flair** (Akbik et al., 2019) offer robust components for traditional tasks like POS tagging, syntactic parsing, and named entity recognition, often leveraging pretrained transformer models. While effective in classical NLP pipelines, they do not address new challenges such as detecting hallucinations, minimizing un-

<sup>1</sup><https://langchain.com>

necessary retrieval in RAG pipelines, or flagging unsafe content generated by LLMs.

In contrast, **DeepPavlov 1.1** is explicitly designed for the current LLM era. Maintains full support for traditional NLP tasks and pre-trained encoder-based models, while also introducing new tools aimed at improving the reliability and controllability of LLM outputs.

### 3 Design and Usage

DeepPavlov models are built and managed via modular configuration files that define all the components required for training, inference, and deployment. Each configuration file includes the following sections: (1) **dataset\_reader/iterator** is responsible for loading data from file; (2) **chainer** is the core abstraction in DeepPavlov, used to construct processing pipelines from heterogeneous components (rule-based, ML, DL); (3) **train** specifies training hyperparameters; (4) **metadata** stores auxiliary variables referenced by other sections.

DeepPavlov emphasizes flexibility and ease of customization. Users can easily adjust hyperparameters, modify preprocessing steps, or swap models within the `chainer` block without breaking the input/output interface.

The framework uses PyTorch as its underlying ML engine, with support for multi-GPU training. DeepPavlov integrates HuggingFace’s `transformers` library, enabling direct use of any `AutoModel`-compatible pretrained model from the HuggingFace Hub.

Models can be used and managed via multiple interfaces: Command-Line Interface (CLI), REST API, or Python. Installation is straightforward via `pip install deepavlov`, and CLI usage examples, code, and documentation are available on the GitHub<sup>2</sup>.

### 4 Reference Models

This section introduces the new models included in the latest release of the DeepPavlov library. All models were trained and evaluated using the library configuration files and are publicly accessible through our demo platform<sup>3</sup>. All training hyperparameters are detailed in the Appendix B.

<sup>2</sup><https://github.com/deepavlov/DeepPavlov>

<sup>3</sup><https://demo.deepavlov.ai>

## 4.1 Hallucination Detection

**Task Formulation.** We formulate hallucination detection as a span level sequence labeling task: given a question, one or more supporting passages, and a generated answer, the model predicts for each answer token whether it is grounded in the provided context or constitutes a hallucination. This span level formulation allows us to capture fine grained hallucinations within otherwise correct answers and supports direct interventions, such as masking or rewriting only the hallucinated fragments. Figure 1 illustrates an example where a hallucinated entity is identified within an otherwise plausible answer.

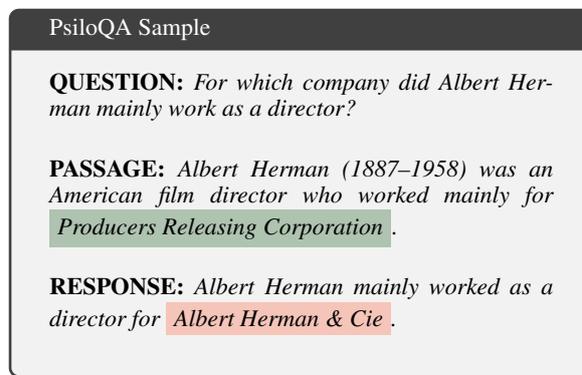


Figure 1: Example from PsiloQA with span level hallucination annotation. The correct entity from the passage is highlighted in `green`, while the hallucinated company name in the model answer is marked in `red`.

**Datasets and Metric.** We evaluate our hallucination detector on the PsiloQA benchmark, a multilingual span-level hallucination detection suite covering 14 languages. Each example consists of a question, one or more supporting passages, and an LLM-generated answer, with character-level annotations marking hallucinated spans in the answer text. We follow the original PsiloQA setup and use its predefined train, development, and test splits, without any additional filtering or relabeling.

For evaluation, we adopt the character-level Intersection over Union (IoU) metric proposed in PsiloQA. Given the predicted hallucination mask and the gold hallucination mask over answer characters, IoU is defined as the size of their intersection divided by the size of their union, computed per example and then macro-averaged over all instances in a language.

Model	Params	Mode	ar	ca	cs	de	en	es	eu	fa	fi	fr	hi	it	sv	zh	Avg.
<i>Open-source language-models</i>																	
Qwen2.5-7B-Instruct	7B	3-shot	33.18	39.13	29.78	28.13	37.68	36.38	36.86	28.43	58.70	29.64	53.22	39.25	35.91	31.25	36.97
Qwen2.5-32B-Instruct	32B	3-shot	30.06	50.60	44.67	25.11	59.01	58.99	38.15	38.52	60.97	43.71	52.71	61.44	37.81	53.00	46.77
Qwen2.5-72B-Instruct	72B	3-shot	45.16	56.12	48.17	32.12	49.82	56.02	42.59	49.35	58.96	49.06	50.37	40.60	45.01	54.02	48.38
gpt-oss-120B	120B	3-shot	38.92	46.56	40.44	27.13	58.75	48.84	39.78	25.25	55.64	38.70	47.16	36.87	43.34	44.72	42.29
<i>Proprietary language-models</i>																	
FActScore (GPT-4o)	—	—	20.75	28.99	10.44	26.68	25.84	28.54	19.68	26.62	28.16	10.21	21.03	43.92	19.25	25.18	23.95
<i>Transformer encoder models</i>																	
lettuce-detect-base	395M	SFT	37.81	44.37	30.08	30.31	43.28	40.08	33.35	32.45	56.44	35.60	16.95	34.97	49.11	35.94	37.20
ModernBERT-base	395M	SFT	55.27	65.70	44.73	46.27	68.23	61.69	50.43	68.63	64.68	53.90	54.15	62.75	67.09	56.95	58.61
mmBERT-base (our)	110M	SFT	58.10	67.01	48.81	54.97	70.67	66.18	50.27	76.61	68.16	56.38	61.19	66.57	66.24	61.58	62.34

Table 2: Character level Intersection over Union (IoU, in %) of span level hallucination detection methods on the PsiloQA test set across 14 languages. Encoder models are supervised fine tuned on the full PsiloQA train split. Language model results are averaged over 5 independent runs; see Table 10 for variance. The rightmost column reports macro averaged IoU across all languages.

**Baselines.** We reuse the PsiloQA evaluation suite and focus on two types of baselines. Encoder-based detectors fine-tuned on PsiloQA: the English only `lettuce-detect-base` model built on ModernBERT (Kovács and Recski, 2025), `ModernBERT-base`<sup>4</sup> trained on PsiloQA, and `mmBERT-base`<sup>5</sup>, a multilingual ModernBERT with 307M parameters covering all 14 languages (Marone et al., 2025; Rykov et al., 2025b). Since PsiloQA and the present work share authorship, we directly reuse the mmBERT checkpoint released by Rykov et al. (2025b) instead of retraining it from scratch. All encoder models take the concatenation of passage, question, and answer and output token level hallucination scores in a single forward pass. LLM-based detectors, treat large generative models as span level judges. We consider *FActScore* with GPT-4o (Min et al., 2023), which decomposes answers into atomic claims and verifies them against retrieved context, the original 3-shot `Qwen2.5-32B-Instruct` baseline, and three additional judges evaluated in this work: `gpt-oss-120b`, `Qwen2.5-(72B|7B)-Instruct`. All LLMs are prompted to insert [HAL] tags around hallucinated spans with default sampling parameters; the prompt template is provided in Figure 4 in Appendix D. Qwen models were evaluated at a temperature of 0.3; `gpt-oss-120B` at 1.0. LLM results (except `Qwen2.5-72B`, single run) are averaged over 5 independent runs.

**Experimental Setup and Results.** ModernBERT is fine-tuned on the multilingual PsiloQA train split with a token level cross entropy loss over

answer tokens. The `lettuce-detect` and `mmBERT` checkpoints are taken directly from Rykov et al. (2025b) and evaluated without modification. LLM baselines are used in a purely prompted regime with 3-shot examples drawn from the PsiloQA training data.

Table 2 reports character level IoU per language. Encoder-based detectors clearly outperform LLM judges: `mmBERT-base` achieves the best overall performance with 70.7 IoU on English and 62.3 macro-averaged IoU across 14 languages, consistently improving over `ModernBERT`. *FActScore* with GPT-4o attains much lower IoU, while `Qwen2.5-32B-Instruct` and our additional `gpt-oss-120B` and `Qwen2.5-72B-Instruct` judges close part of the gap but remain below the encoder models, especially in low resource languages. The smallest model, `Qwen2.5-7B-Instruct`, performs competitively only on a subset of languages and reaches 37.68 IoU on English and 36.97 IoU on average. Overall, encoder-based detectors trained on PsiloQA, and `mmBERT` in particular, provide the most accurate and efficient span level hallucination detection for our multilingual setting.

## 4.2 Evergreen Questions Classification

**Task Formulation.** Evergreen Question Classification is the task of identifying factual questions whose correct answers are highly unlikely to change over extended periods of time. We formulate this as a binary classification problem: Given a question, predict whether it is evergreen or non-evergreen (see examples in Table 3).

By serving as a real-time preprocessing filter, we can determine whether to rely solely on the internal knowledge of the LLM (in the case of ev-

<sup>4</sup><https://hf.co/answerdotai/ModernBERT-base>

<sup>5</sup><https://hf.co/jhu-clsp/mmBERT-base>

Evergreen Questions	Non-Evergreen Questions
Who painted the ‘Mona Lisa’?	What time is it?
Which is lighter: a kilogram of feathers or a kilogram of iron?	Who won the last football World Cup?
What is the ultimate question of life, the universe, and everything?	The last time a bright comet was visible to the naked eye?

Table 3: Examples of Evergreen and Non-Evergreen questions. Among the non-evergreen examples, a darker red background highlights answers that change more rapidly, and a lighter red background indicates answers that change relatively slow.

Model	Params	Mode	Russian	English	French	German	Hebrew	Arabic	Chinese	Average
<i>Open-source language-models</i>										
Qwen2.5-7B-Instruct	7B	10-shot	78.2	78.9	78.6	79.4	69.2	71.1	77.4	76.1
Qwen2.5-32B-Instruct	32B	10-shot	88.2	88.5	87.5	88.3	86.2	86.2	87.2	87.4
Qwen 2.5-72B-Instruct	72B	10-shot	80.6	81.5	80.2	80.5	78.1	75.8	76.8	79.1
gpt-oss-120b	120B	10-shot	<b>92.5</b>	<b>95.2</b>	<b>94.3</b>	<b>92.4</b>	<b>93.1</b>	<b>93.6</b>	<b>92.1</b>	<b>93.3</b>
<i>Proprietary language-models</i>										
GPT-4.1	–	10-shot	80.6	79.4	81.6	81.3	80.3	81.1	80.9	80.7
<i>Transformer encoder models</i>										
EG-BERT-base	110M	SFT	89.3	90.0	88.9	88.4	88.9	88.3	<u>90.2</u>	89.1
<b>BERT-base (our)</b>	110M	SFT	87.4	88.1	87.2	86.5	85.5	87.3	87.2	87.0
<b>ModernBERT-base (our)</b>	395M	SFT	75.9	88.0	82.8	81.3	78.4	77.5	85.7	81.4
EG-E5-large	560M	SFT	91.0	91.3	90.9	91.0	90.4	90.0	89.7	90.6

Table 4: EvergreenQA classifier F1-weighted scores comparison across different languages.

ergreen questions) or to apply the RAG pipeline. Relying on the RAG pipeline for every query is not always the best solution because: (1) retrieval in RAG adds additional latency; (2) noisy context returned by retrieval can deteriorate the quality of generation (Fang et al., 2024).

Therefore, it is best practice to include an adaptive component that decides whether the LLM alone can answer the question or whether the RAG component should be used (Moskvoretskii et al., 2025).

**Datasets.** To train and evaluate our Evergreen classifier we leverage an **EverGreenQA** (Pletenev et al., 2025) dataset comprising of 4,757 examples and covering seven languages. We evaluated our models in both the EverGreenQA test set and the multilingual version of the FreshQA data set (Vu et al., 2024), which had been translated into all target languages in Pletenev et al. (2025).

**Baselines.** To contextualize our results, we compare them with zero-shot LLM baselines of various sizes, as well as with BERT-family models fine-tuned on the binary classification task from EverGreenQA (Pletenev et al., 2025). The prompt used for LLM-based classification is provided in Figure 5 in Appendix D. All LLMs were evaluated at a temperature of 0.0

**Experimental Setup and Results.** As a primary evaluation metric, we follow EverGreenQA and report the weighted F1 score, using the same train/test split.

Our experimental results presented in Table 4 show that our models do not outperform larger baselines on the EverGreenQA test set, multilingual BERT-base achieves modest improvements over the original BERT-base on FreshQA, particularly for some languages. This indicates competitive performance in specific cases, although the gap compared to EG-E5-large highlights the challenges of generalizing temporal sensitivity classification to unseen multilingual data. ModernBERT’s results are mixed: overall weaker on FreshQA, but with strong performance for certain languages, matching EG-E5-large for Russian and ranking second for French. This suggests limited cross-lingual generalization but potential in certain language-specific contexts. Our retrained multilingual BERT-base demonstrates consistent improvements over the original, offering a practical lightweight alternative.

### 4.3 Toxicity Detection

**Task Formulation.** We frame the task of toxicity classification as a binary classification problem. The goal is to determine whether a given text contains toxic content, such as vulgar, obscene, or profane language. This component serves as

Model	Params	Mode	EN	RU	UK	DE	ES	AR	AM	HI	ZH	IT	FR	HI-EN	HE	JA	TT	Avg.
<i>Open-source language-models</i>																		
Qwen2.5-7B-Instruct	7B	3-shot	93.4	93.6	83.4	77.1	86.4	76.4	62.9	75.2	66.8	74.7	95.4	67.2	67.2	65.3	72.2	77.7
Qwen2.5-32B-Instruct	32B	3-shot	<b>97.5</b>	95.7	84.4	75.3	86.6	78.6	53.7	76.0	71.9	<u>75.9</u>	<b>97.4</b>	<u>73.7</u>	73.3	75.0	74.1	79.7
Qwen2.5-72B-Instruct	72B	3-shot	95.8	93.4	87.0	80.8	<b>90.9</b>	<u>79.1</u>	<u>64.6</u>	85.0	<u>76.8</u>	72.0	94.8	<b>78.4</b>	71.1	77.9	75.9	<u>81.6</u>
gpt-oss-120B	120B	3-shot	96.4	<u>96.0</u>	86.0	80.7	<u>88.9</u>	<b>81.9</b>	56.3	73.2	67.1	71.5	<u>96.2</u>	69.2	75.0	66.1	<u>82.5</u>	80.1
<i>Transformer encoder models</i>																		
TextDetox-2024-RoBERTa-large	355M	SFT	96.5	<b>97.9</b>	92.5	<u>87.6</u>	87.0	77.8	<b>77.8</b>	<u>93.6</u>	73.2	–	–	–	–	–	–	<b>87.1</b>
TextDetox-2025-RoBERTa-large	355M	SFT	92.3	95.3	<b>96.0</b>	73.3	71.3	66.3	55.8	<b>97.3</b>	<b>91.8</b>	58.6	92.4	61.0	<b>87.8</b>	<b>87.7</b>	57.4	78.9
TextDetox-BERT-base	110M	SFT	90.4	92.2	<u>94.6</u>	51.8	72.9	51.4	63.2	72.7	67.0	64.9	91.3	68.5	<u>86.9</u>	<u>86.4</u>	61.7	74.4
<b>BERT-base (our)</b>	110M	SFT	<u>97.0</u>	91.0	87.9	<b>87.8</b>	81.6	67.8	58.6	89.2	60.1	77.8	95.2	71.9	67.5	73.9	<b>87.4</b>	<u>81.6</u>

Table 5: Toxicity classification F1 scores across different languages. Language model results are averaged over 5 independent runs; see Table 11 for variance.

a safeguard when LLM’s alignment mechanisms fail to prevent harmful outputs (Moskovskiy et al., 2024).

**Dataset.** We train and evaluate our model on the multilingual TextDetox (Dementieva et al., 2024) dataset, which includes 5,000 examples for each of 15 languages.

**Baselines.** We compare our model against several encoder-based baselines, including BERT-base and XLM-RoBERTa-large, using checkpoints from the official TextDetox hub<sup>6</sup>. The prompt used for LLM-based classification is provided in Figure 3 in Appendix D. All LLMs were evaluated at a temperature of 0.0

**Experimental Setup and Results.** The BERT-base model is trained as a binary example-level classifier. Since the original dataset split is unavailable, we divide the TextDetox dataset into 70% training, 10% validation, and 20% test sets, preserving the language distribution.

Table 5 shows that our multilingual BERT-base model ranks second overall, although it is significantly smaller than XLM-RoBERTa-large model. It achieves the highest F1-score in English (97.00) and outperforms baselines in several other languages, demonstrating its strong cross-lingual generalization.

## 5 Performance

**Experimental Setup.** All experiments were conducted on single Nvidia H200 GPU with 140 Gb VRAM. The vLLM framework was used for LLM inference, while the DeepPavlov and Transformers libraries were used to run encoder models.

<sup>6</sup><https://hf.co/textdetox>

Encoder models were executed with a batch size of 256 to maximize throughput.

**Performance Analysis.** Figure 2 and Table 6 provide a detailed comparison of model performance. Our models consistently achieve accuracy near state-of-the-art levels while exhibiting inference speeds that are two orders of magnitude faster compared to LLMs. The results reveal a consistent pattern: larger models yield higher accuracy but require more inference time.

## 6 Conclusion and Future Work

DeepPavlov 1.1 is an updated tool that makes LLMs more reliable and safe. In the future, we plan to continue improving DeepPavlov by adding more features and making it even easier for researchers and developers to build safe and trustworthy AI systems.

### Limitations

**Framework-Level Limitations.** DeepPavlov 1.1 does not yet offer native integration with modern LLM orchestration pipelines, which limits its out-of-the-box applicability in production workflows. Although the framework offers API and Python integration, we did not conduct rigorous latency, throughput, or scalability testing. Users should verify performance under production constraints such as inference time or memory footprint. The framework depends on external libraries such as PyTorch and HuggingFace Transformers. API changes, deprecations, or version incompatibilities could break core functionality or degrade performance.

**Experimental Setup.** Our evaluation protocol does not include multiple training runs with different random seeds. Similarly, LLMs are only evaluated once per experiment. We also performed only minimal exploratory data analysis (EDA) and

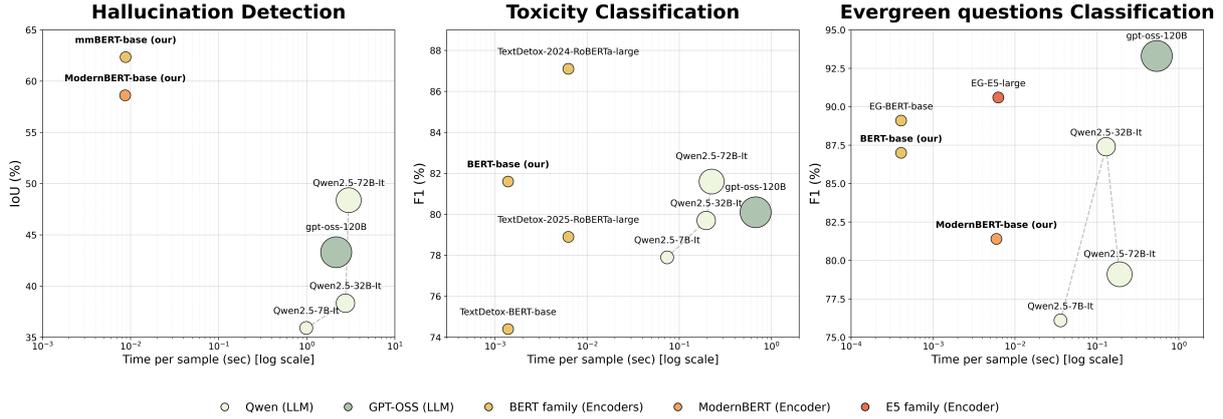


Figure 2: Trade-off plot illustrating how different models perform on three tasks supported in the updated DP library. Each subplot shows the relationship between *inference latency* and *task performance*.

ablation studies. This limits our understanding of how design choices affect model behavior.

**Hallucination Detection.** The current hallucination detection model is still in an early stage, its span-level performance remains low, especially for summarization tasks. The detector is restricted to contextual factual hallucinations and does not address other types such as logical or common-sense errors. Its effectiveness also hinges on the retriever’s ability to supply consistent and relevant context. While the model offers a basic safeguard, it is not yet suitable for high-stakes applications and requires substantial future development.

**Toxicity Detection.** The toxicity classifier was trained on the TextDetox dataset, which contains a non-negligible amount of label noise and inconsistent annotations across languages. This can cause instability in predictions, especially for borderline or multilingual cases.

**Model Coverage and Scope.** DeepPavlov 1.1 includes only three reliability-oriented models: a hallucination detector, an evergreen classifier, and a toxicity classifier. While these were prioritized based on user demand, other crucial capabilities, such as uncertainty-based detectors, adversarial input filters are absent and should be considered for future releases.

## Ethics Statement

All models and experiments described in this paper were developed and evaluated exclusively using publicly available datasets. No proprietary, private, or personally identifiable data were used at any stage of model training, testing, or deployment. We are committed to transparency and reproducibility: all code, configuration files, and

pretrained models are released under an open-source license to facilitate independent verification and responsible reuse by the research community.

## References

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Gregory S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian J. Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Józefowicz, Lukasz Kaiser, Manjunath Kudlur, and 21 others. 2016. [Tensorflow: Large-scale machine learning on heterogeneous distributed systems](#). *CoRR*, abs/1603.04467.
- Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019. [FLAIR: an easy-to-use framework for state-of-the-art NLP](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Demonstrations*, pages 54–59. Association for Computational Linguistics.
- Mikhail Burtsev, Alexander Seliverstov, Rafael Airapetyan, Mikhail Arkhipov, Dilyara Baymurzina, Nickolay Bushkov, Olga Gureenkova, Taras Khakhulin, Yuri Kuratov, Denis Kuznetsov, Alexey Litinsky, Varvara Logacheva, Alexey Lymar, Valentin Malykh, Maxim Petrov, Vadim Polulyakh, Leonid Pugachev, Alexey Sorokin, Maria Vikhreva, and Marat Zaynutdinov. 2018. [DeepPavlov: Open-source library for dialogue systems](#). In *Proceedings of ACL 2018, System Demonstrations*, pages 122–127, Melbourne, Australia. Association for Computational Linguistics.
- Anastasia Chizhikova, Vasily Konovalov, and Mikhail Burtsev. 2023. Multilingual case-insensitive named entity recognition. In *Advances in Neural Computation, Machine Learning, and Cognitive Research VI*,

- pages 448–454, Cham. Springer International Publishing.
- Daryna Dementieva, Daniil Moskovskiy, Nikolay Babakov, Abinew Ali Ayele, Naqee Rizwan, Florian Schneider, Xintong Wang, Seid Muhie Yimam, Dmitry Ustalov, Elisei Stakovskii, Alisa Smirnova, Ashraf Elnagar, Animesh Mukherjee, and Alexander Panchenko. 2024. Overview of the multilingual text detoxification task at pan 2024. In *Working Notes of CLEF 2024 - Conference and Labs of the Evaluation Forum*. CEUR-WS.org.
- Daryna Dementieva, Vitaly Protasov, Nikolay Babakov, Naqee Rizwan, Ilseyar Alimova, Caroline Brun, Vasily Konovalov, Arianna Muti, Chaya Liebeskind, Marina Litvak, Debora Nozza, Shehryaar Shah Khan, Sotaro Takeshita, Natalia Vanetik, Abinew Ali Ayele, Florian Schneider, Xintong Wang, Seid Muhie Yimam, Ashraf Elnagar, and 2 others. 2025. [Overview of the multilingual text detoxification task at PAN 2025](#). In *Working Notes of the Conference and Labs of the Evaluation Forum, CLEF 2025, Madrid, Spain, 9-12 September 2025*, volume 4038 of *CEUR Workshop Proceedings*, pages 3535–3567. CEUR-WS.org.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiaojun Chen, and Ruifeng Xu. 2024. [Enhancing noise robustness of retrieval-augmented language models with adaptive adversarial training](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10028–10039, Bangkok, Thailand. Association for Computational Linguistics.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. [AllenNLP: A deep semantic natural language processing platform](#). In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 1–6, Melbourne, Australia. Association for Computational Linguistics.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. [Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Matthew Honnibal. 2017. `spacy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing`. (*No Title*).
- Dmitry Karpov and Vasily Konovalov. 2023. [Knowledge transfer between tasks and languages in the multi-task encoder-agnostic transformer-based models](#). In *Computational Linguistics and Intellectual Technologies*, volume 2023.
- Ádám Kovács and Gábor Recski. 2025. [Lettucedetect: A hallucination detection framework for RAG applications](#). *CoRR*, abs/2502.17125.
- Nikita Krayko, Ivan Sidorov, Fedor Laputin, Alexander Panchenko, Daria Galimzianova, and Vasily Konovalov. 2025. Rurage: Robust universal rag evaluator for fast and affordable qa performance testing. In *Advances in Information Retrieval*, pages 135–145, Cham. Springer Nature Switzerland.
- 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 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- 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.
- Marc Marone, Orion Weller, William Fleshman, Eugene Yang, Dawn J. Lawrie, and Benjamin Van Durme. 2025. [mmbert: A modern multilingual encoder with annealed language learning](#). *CoRR*, abs/2509.06888.
- Sewon Min, Kalpesh Krishna, Xinxin Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [Factscore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 12076–12100. Association for Computational Linguistics.
- Daniil Moskovskiy, Sergey Pletenev, and Alexander Panchenko. 2024. [LLMs to replace crowdsourcing for parallel data creation? the case of text detoxification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14361–14373, Miami, Florida, USA. Association for Computational Linguistics.
- Viktor Moskvoretskii, Maria Marina, Mikhail Salnikov, Nikolay Ivanov, Sergey Pletenev, Daria Galimzianova, Nikita Krayko, Vasily Konovalov, Irina

- Nikishina, and Alexander Panchenko. 2025. [Adaptive retrieval without self-knowledge? bringing uncertainty back home](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6355–6384, Vienna, Austria. Association for Computational Linguistics.
- Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. [Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 10862–10878. Association for Computational Linguistics.
- OpenAI. 2022. Chatgpt: Optimizing language models for dialogue. <http://openai.com/blog/chatgpt>.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 8024–8035.
- Jason Phang, Phil Yeres, Jesse Swanson, Haokun Liu, Ian F. Tenney, Phu Mon Htut, Clara Vania, Alex Wang, and Samuel R. Bowman. 2020. [jiant 2.0: A software toolkit for research on general-purpose text understanding models](#). <http://jiant.info/>.
- Sergey Pletenev, Maria Marina, Nikolay Ivanov, Daria Galimzianova, Nikita Krayko, Mikhail Salnikov, Vasily Konovalov, Alexander Panchenko, and Viktor Moskvoretskii. 2025. [Will it still be true tomorrow? multilingual evergreen question classification to improve trustworthy QA](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 8603–8620, Suzhou, China. Association for Computational Linguistics.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, ACL 2020, Online, July 5-10, 2020*, pages 101–108. Association for Computational Linguistics.
- Elisei Rykov, Valerii Olisov, Maksim Savkin, Artem Vazhentsev, Kseniia Titova, Alexander Panchenko, Vasily Konovalov, and Julia Belikova. 2025a. [SmurfCat at SemEval-2025 task 3: Bridging external knowledge and model uncertainty for enhanced hallucination detection](#). In *Proceedings of the 19th International Workshop on Semantic Evaluation (SemEval-2025)*, pages 1034–1045, Vienna, Austria. Association for Computational Linguistics.
- Elisei Rykov, Kseniia Petrushina, Maksim Savkin, Valerii Olisov, Artem Vazhentsev, Kseniia Titova, Alexander Panchenko, Vasily Konovalov, and Julia Belikova. 2025b. [When models lie, we learn: Multilingual span-level hallucination detection with PsiloQA](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 11663–11682, Suzhou, China. Association for Computational Linguistics.
- Maksim Savkin and Vasily Konovalov. 2024. [Tuning-free discriminative nearest neighbor few-shot intent detection via consecutive knowledge transfer](#). In *Recent Trends in Analysis of Images, Social Networks and Texts*, pages 96–110, Cham. Springer Nature Switzerland.
- Maksim Savkin, Anastasia Voznyuk, Fedor Ignatov, Anna Korzanova, Dmitry Karpov, Alexander Popov, and Vasily Konovalov. 2024. [DeepPavlov 1.0: Your gateway to advanced NLP models backed by transformers and transfer learning](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 465–474, Miami, Florida, USA. Association for Computational Linguistics.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. [Gemma 3 technical report](#). *arXiv preprint arXiv:2503.19786*.
- 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, and 49 others. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2024. [FreshLLMs: Refreshing large language models with search engine augmentation](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13697–13720, Bangkok, Thailand. Association for Computational Linguistics.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. [Superglue: A stickier benchmark for general-purpose language understanding systems](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December-14, 2019, Vancouver, BC, Canada*, pages 3261–3275.

Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024. [Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference](#). *CoRR*, abs/2412.13663.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. [Huggingface’s transformers: State-of-the-art natural language processing](#). *CoRR*, abs/1910.03771.

## A Latency

Model	Params	PsiloQA	EvergreenQA	TextDetox
<i>Open-source language-models</i>				
Qwen2.5-7B-Instruct	7B	990	36	74
Qwen2.5-32B-Instruct	32B	2756	130	196
Qwen2.5-72B-Instruct	72B	3000	190	225
gpt-oss-120b	120B	2164	539	675
<i>Transformer encoder models</i>				
EG-E5-large	560M	–	6.3	–
ModernBERT-base (our)	395M	8.7	5.9	–
mmBERT-base (our)	395M	8.8	–	–
RoBERTa-large	355M	–	–	6.3
BERT-base (our)	110M	–	0.4	1.4

Table 6: Model performance in milliseconds per task sample across different tasks.

## B Training Hyperparameters

Hyperparameter	Value
Batch size	8
Optimizer	AdamW
Learning rate	$1 \times 10^{-5}$
Weight decay	$1 \times 10^{-3}$
Adam betas	(0.9, 0.999)
Adam epsilon	$1 \times 10^{-6}$
Clip norm	1.0
Max epochs	6
Model selection metric	F1-weighted

Table 7: Training hyperparameters for Hallucination Detector.

Hyperparameter	Value
Max sequence length	512
Batch size	16
Optimizer	AdamW
Learning rate	$1 \times 10^{-5}$
Weight decay	$1 \times 10^{-6}$
Adam betas	(0.9, 0.999)
Adam epsilon	$1 \times 10^{-6}$
Clip norm	1.0
Min. learning rate	$2 \times 10^{-6}$
Learning rate drop patience	15
Learning rate drop factor	1.5
Max epochs	20
Model selection metric	F1-weighted

Table 8: Training and model hyperparameters for EvergreenQA classifier.

Hyperparameter	Value
Max sequence length	64
Batch size	64
Optimizer	AdamW
Adam betas	(0.9, 0.999)
Adam epsilon	$1 \times 10^{-6}$
Learning rate	$1 \times 10^{-5}$
Learning rate drop patience	5
Learning rate drop factor	2.0
Max epochs	100
Model selection metric	F1-weighted

Table 9: Training and model hyperparameters for Toxicity Classification.

## C LLM Judge Variance

Model	Params	Temp	IoU (%)	AP (%)
Qwen2.5-7B-Instruct	7B	0.3	36.97 ± 0.62	61.90 ± 1.03
Qwen2.5-32B-Instruct	32B	0.3	46.77 ± 0.33	64.49 ± 0.73
gpt-oss-120B	120B	1.0	42.29 ± 0.55	58.14 ± 0.50

Table 10: Variance of LLM-as-a-judge hallucination detection across 5 independent runs. Each model is evaluated with 3-shot prompting on the PsiloQA (Rykov et al., 2025b) test set (all 14 languages). We report character-level macro-averaged IoU and micro-averaged AP (mean ± std over 5 runs). Qwen models use temperature 0.3; gpt-oss-120B uses temperature 1.0.

Model	Params	Temp	F1-weight
Qwen2.5-7B-Instruct	7B	1.0	77.66 ± 0.27
Qwen2.5-32B-Instruct	32B	0.3	79.62 ± 0.12

Table 11: Variance of LLM-as-a-judge hallucination detection across 5 independent runs. Each model is evaluated on the TextDetox (Dementieva et al., 2024) test set. We report F1-weighted (mean ± std over 5 runs).

## D Prompt Templates for LLM Baselines

### Toxicity Classification Prompt

You are a safety classifier that determines whether a single user message is toxic or non-toxic.

Toxic content includes insults, hate speech, threats, explicit harassment, or demeaning language toward an individual or group. Non-toxic content is neutral, polite, or otherwise acceptable.

Read the message carefully and decide if it is **Toxic** or **Non-toxic**. At the end, respond with a single line in the form:

Classification: Toxic

or

Classification: Non-toxic.

Examples:

{few\_shot}

Message: {input\_text}

Classification:

Figure 3: Prompt for binary toxicity classification.

### Hallucination Detection Prompt

You are an expert hallucination detector for question answering with retrieved context.

Given:

- a context passage from Wikipedia,
- a user question, and
- an LLM answer,

you must identify all hallucinated spans in the answer.

A hallucinated span is any part of the answer that:

- contradicts the context, or
- introduces specific factual details that are not supported by the context or by the gold answer.

Your output must be the model answer text where you wrap every hallucinated span in [HAL] and [/HAL] tags.

**CRITICAL INSTRUCTIONS:**

- Do not change, rephrase, re order, or truncate the answer.
- Do not add new information.
- Only insert [HAL] before and [/HAL] after hallucinated spans.
- If there is no hallucination, return the answer unchanged (with no [HAL] tags).

Examples:

{few\_shot}

Return only the model answer text, where hallucinated spans are wrapped in [HAL] and [/HAL] tags. Do not add any explanation or commentary.

Knowledge source: {passage}

Question: {question}

Answer: {answer}

Answer with highlighted spans:

Figure 4: Prompt for span level hallucination detection.

### Evergreen Detection Prompt

You are a helpful assistant that classifies questions based on their temporality.

There are two classes:

**Immutable:** the answer almost never changes over time (for example, historical facts, birth years, names of past events).

**Mutable:** the answer typically changes over the course of several years or less (for example, current leaders, upcoming events, latest statistics).

Think carefully about each question and decide whether it is Immutable or Mutable. At the end, answer with exactly one line of the form:

Classification: Immutable

or

Classification: Mutable.

Examples:

{few\_shot}

Question: {input\_question}

Classification:

Figure 5: Prompt for Evergreen classification.

# PropGenie: A Multi-Agent Conversational Framework for Real Estate Assistance

Chang Shen, Shaozu Yuan, Kuizong Wu, Long Xu, Meng Chen\*

Yep AI, Melbourne, Australia

{chang.shen, shaozu.yuan, vincent.wu, neo.xu}@yepai.io

chenmengdx@gmail.com

## Abstract

In this paper, we present PropGenie, a novel multi-agent framework based on large language models (LLMs) to deliver comprehensive real estate assistance in real-world scenarios. PropGenie coordinates eight specialized sub-agents, each tailored for distinct tasks, including search and recommendation, question answering, financial calculations, and task execution. To enhance response accuracy and reliability, the system integrates diverse knowledge sources and advanced computational tools, leveraging structured, unstructured, and multimodal retrieval-augmented generation techniques. Experiments on real user queries show that PropGenie outperforms both a general-purpose LLM (OpenAI’s o3-mini-high) and a domain-specific chatbot (Realty AI’s Madison) in real estate scenarios. We hope that PropGenie serves as a valuable reference for future research in broader AI-driven applications.

## 1 Introduction

The real estate industry, while traditional, remains a cornerstone of both residential needs and investment strategies (Hudson-Wilson et al., 2005). With the advent of artificial intelligence, the integration of this transformative technology has become imperative (Ullah et al., 2018; Seagraves, 2023; Haurum et al., 2024) to address the persistent challenge of **information asymmetry**—where buyers’ and investors’ decisions are often influenced by intricate psychological factors and external market dynamics (Elster, 2016). Real estate transactions encompass a vast array of complex scenarios, each demanding a sophisticated understanding of domain-specific knowledge, including property valuation, legal frameworks, financing mechanisms, and evolving market trends. The dynamic and multifaceted nature of these inquiries necessitates a robust AI-driven system capable of synthesizing and

\*Corresponding author.

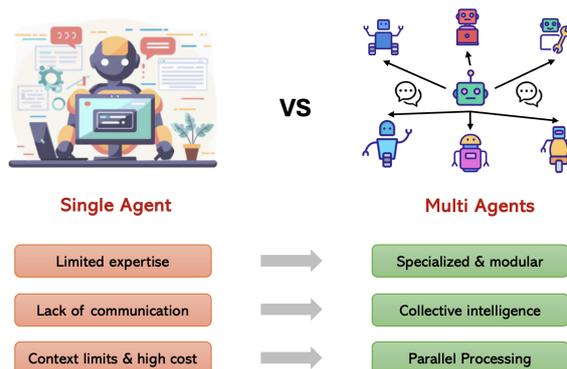


Figure 1: The difference between single-agent and multi-agents designs.

interpreting information from diverse fields. Moreover, the high-stakes nature of real estate transactions imposes stringent accuracy requirements on system outputs. These challenges underscore the intricacy of developing an intelligent AI assistant for real estate applications—an endeavor that is both demanding and highly rewarding.

Several studies have explored the development of automated chatbots and virtual agents to provide 24/7 customer service, capture potential leads, offer legal consultations, and reduce administrative costs in the real estate industry. These efforts can be broadly categorized into two approaches: 1) **Traditional dialogue system-based methods**, which construct virtual assistants using intent-based models (Quan et al., 2018; Cao and Nguyen, 2021), frequently asked question (FAQ) systems (Tanović and Hasibović, 2024), or knowledge graph (KG)-based frameworks (Yang et al., 2024b). This approach primarily focuses on developing classification models and curating high-quality, domain-specific datasets. 2) **LLM-based agents** (Pagar, 2024; Haurum et al., 2024; Gloria et al., 2025), which leverage the advanced language comprehension capabilities of large language models (LLMs) to handle complex user queries. These agents enhance their functionality through tool calling (Qin

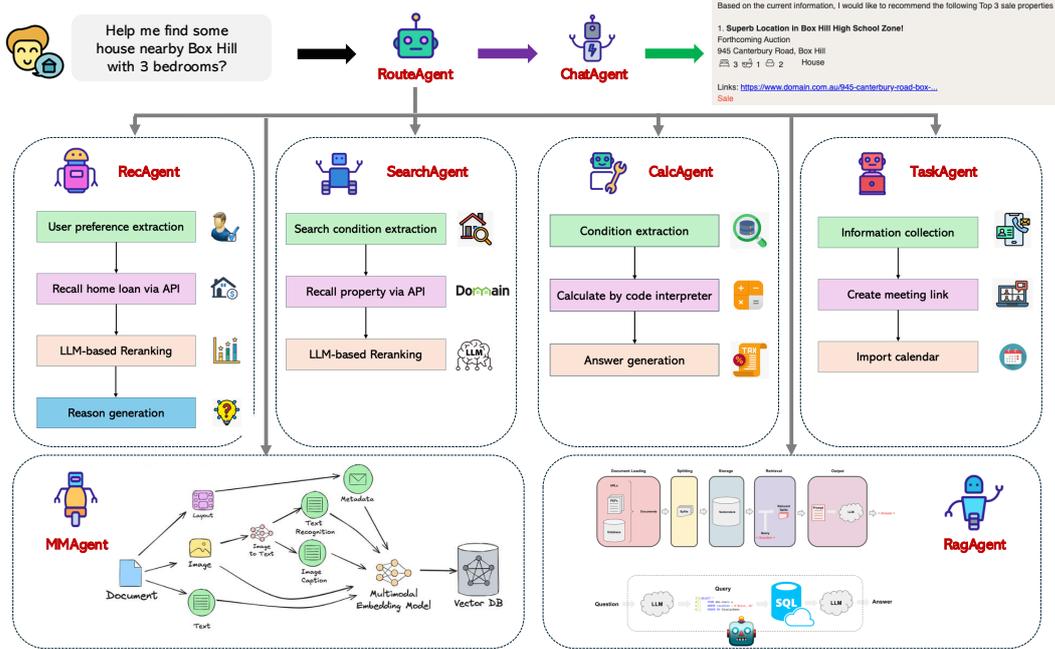


Figure 2: The overall architecture of the proposed multi-agent framework PropGenie.

et al., 2024) and retrieval-augmented generation (RAG) (Lewis et al., 2020) to improve response accuracy and contextual understanding.

With the advent of the LLM era (OpenAI et al., 2024a; GeminiTeam et al., 2024; DeepSeek-AI et al., 2025), agentic AI has demonstrated substantial potential in addressing complex real-world challenges (Durante et al., 2024). However, as illustrated in Figure 1, existing research in the real estate domain predominantly employs single-agent systems, which exhibit several limitations. First, a single agent typically lacks the specialized expertise necessary to effectively handle diverse user requests, such as property searches, home loan recommendations, question answering, and tax-related calculations. Second, the absence of inter-agent communication and cross-verification mechanisms increases susceptibility to erroneous outputs caused by LLM hallucinations (Li et al., 2024b). Third, maintaining contextual coherence in single-agent systems requires large prompts for each request, resulting in increased inference costs and reduced efficiency (Zhang et al., 2024; Wang et al., 2025b). In contrast, a multi-agent framework addresses these limitations through specialized modular design, collective intelligence, and parallel processing, making it a natural choice for developing scalable and efficient real estate applications (Zhao et al., 2024; Team et al., 2024; Gao et al., 2024a).

Motivated by the above analysis, we propose **PropGenie**, a multi-agent conversational frame-

work designed for real-world real estate assistance. Following a *system-over-model* paradigm, PropGenie integrates domain-specific knowledge bases and advanced computational tools to deliver accurate and contextually relevant responses. Specifically, we develop eight specialized sub-agents that collaboratively handle ten distinct tasks, including property search, home loan recommendations, stamp duty and land tax calculations, monthly repayment estimations, home-buying policy QA, real-time interest rate inquiries, property project QA, automatic task execution, and open-domain chitchat. Additionally, PropGenie enhances user engagement through emotion recognition and satisfaction prediction. We evaluate PropGenie through comprehensive automatic and human evaluations, supplemented by detailed case studies to analyze its strengths and limitations. By demonstrating effectiveness in the real estate domain, our system provides valuable insights for future research in broader application scenarios.

## 2 System Architecture

In this section, we present the overall architecture of PropGenie, starting with an overview of its task scope, followed by a detailed discussion of the design and functionality of each sub-agent.

### 2.1 Task Scope

PropGenie provides an integrated solution covering property search, financing, and planning. Based on extensive analysis of user queries and feedback

Tasks	Examples
Property search	Find a 3-bedroom house with a swimming pool near Box Hill.
Home loan recommendation	Recommend some good ANZ home loan products that support an offset account.
Stamp duty calculation	How much stamp duty do I need to pay for a \$1.8 million house in Melbourne?
Land tax calculation	How much land tax do I need to pay for a \$1.5 million house and land package in NSW?
Monthly repayment estimation	What’s the monthly repayment for a \$1 million home loan over 20 years?
Home-buying policy QA	What’s the process of buying a home in Australia as an overseas buyer?
Interest rate QA	What’s the current interest rate for a principal and interest loan with Westpac?
Property project QA	How many apartment units are there in the YarraBend project?
Task execution	Can you help me schedule an online meeting with the agent?
Open-domain chitchat	What are the main differences between living in Sydney and Melbourne?

Table 1: Tasks and Examples Illustrations. Detailed distribution for each task can be found in Appendix A.1.

from broker agents, we define ten core tasks: property search, home loan recommendations, stamp duty and land tax calculations, monthly repayment estimation, home-buying policy QA, real-time interest rate inquiries, property project QA, automatic task execution, and open-domain chitchat. Table 1 illustrates representative examples for each task, clarifying the system’s scope and capabilities.

## 2.2 Components

As illustrated in Figure 2, PropGenie adopts a multi-agent framework to efficiently handle diverse real estate tasks. The system groups related functionalities into dedicated agents, each with a streamlined workflow, simplifying the overall architecture. For example, stamp duty, land tax, and monthly repayment calculations share a common workflow involving query interpretation and code interpreter; thus, we consolidate these tasks into a single CalcAgent. Similarly, the RagAgent manages retrieval-augmented generation from structured and unstructured databases, addressing queries such as interest rates and home-buying policies. Currently, PropGenie leverages GPT-4o (OpenAI et al., 2024b) as its core LLM, yet the framework remains flexible, enabling integration of future advanced LLMs.

**RouteAgent.** The RouteAgent acts as the central orchestrator, autonomously managing system workflows. Its responsibilities include: 1) **Intent Understanding and Task Decomposition:** interpreting user queries, identifying intents, decomposing complex queries into sub-tasks, and assigning these tasks to appropriate sub-agents. It also facilitates collaboration among sub-agents to leverage their specialized capabilities. 2) **Response Aggregation and Conflict Resolution:** aggregating responses from multiple sub-agents into a coherent reply, resolving inconsistencies by selecting the most reliable response or reassigning tasks for verification, and delegating out-of-scope queries to the ChatAgent for fallback responses. 3) **Con-**

**text Management and Emotional Intelligence:** maintaining conversational context by rewriting multi-turn interactions into self-contained queries, thereby simplifying context handling and enhancing robustness of RAG-based sub-agents. Additionally, it monitors user sentiment, adapting responses empathetically or strategically shifting topics when detecting user frustration.

**SearchAgent.** The SearchAgent retrieves relevant property listings by following a structured information retrieval workflow. It first extracts key search criteria from user queries, such as price range, location (e.g., suburb), housing configuration (e.g., number of rooms), and amenities (e.g., parking, swimming pool, gym). Next, it invokes a web search tool to fetch property listings from third-party APIs<sup>1</sup>. Finally, it re-ranks the top 50 retrieved properties against the user’s requirements and selects the three most relevant listings. Both query parsing and result reranking are powered by an LLM, enhancing language understanding and retrieval accuracy.

**RecAgent.** The RecAgent recommends home loan products tailored to user profiles and preferences. It first extracts user preferences from conversational context, such as preferred bank, repayment type, interest rate type, loan term, and additional features (e.g., offset accounts). Similar to the SearchAgent, it retrieves relevant home loan products from the web<sup>2</sup> using a search tool and re-ranks them with LLM based on user preferences. Additionally, the RecAgent analyzes the advantages and drawbacks of each loan product, providing comprehensive evaluations to support informed decision-making.

**CalcAgent.** To deliver accurate computational responses in real estate inquiries, we introduce CalcAgent, a specialized module designed to address complex financial queries such as stamp duty,

<sup>1</sup><https://developer.domain.com.au/>

<sup>2</sup><https://www.finder.com.au/home-loans>

land tax, and mortgage repayment calculations. These queries pose significant challenges due to their reliance on logical reasoning and numerical precision. CalcAgent processes user-provided parameters, including property value, location, and foreign buyer status for tax assessments, as well as loan amount, loan term, and bank selection for mortgage calculations. To ensure computational accuracy, CalcAgent employs a structured three-step approach: (1) Entity extraction identifies essential numerical and categorical parameters from user queries; (2) A code interpreter dynamically generates and executes Python scripts based on predefined formulas; (3) An LLM formulates coherent, contextually relevant responses grounded in computed outcomes. Given that tax regulations and mortgage formulas vary across jurisdictions and evolve over time, we periodically collect official tax rules from government sources and formalize them into standardized mathematical expressions. Similarly, mortgage repayment formulas are derived from financial institutions’ policies. These structured formulas are integrated into CalcAgent’s prompts as auxiliary knowledge, enabling the model to generate accurate and executable programs aligned with current financial regulations.

**RagAgent.** RagAgent addresses inquiries related to interest rates and home-buying policies using a unified retrieval-augmented generation (RAG) framework (Gao et al., 2024b). Given the high-stakes nature of real estate transactions, ensuring information accuracy is critical. To achieve this, we systematically collect and structure domain-specific knowledge from authoritative sources. For interest rates, we aggregate real-time data from financial institutions’ open APIs and regularly update a structured database. For home-buying policies, we employ web crawlers to extract regulatory updates from official government websites, accommodating frequent policy changes. Extracted content is transformed into question-answer (QA) pairs, which are automatically refined through LLM-driven quality checking and filtering (Cheng et al., 2024). These curated QA pairs are then encoded into vector representations and stored in a vector database for efficient retrieval. RagAgent processes interest rate queries by translating them into SQL queries (Wang et al., 2025a) to retrieve precise results from the structured database. For home-buying policy questions, it retrieves relevant QA pairs from the vector database and leverages

an LLM to generate concise, contextually accurate responses. By integrating rigorous knowledge curation with the RAG approach, RagAgent ensures reliable, timely, and accurate real estate assistance.

**MMAgent.** Inquiries about off-plan property projects represent a significant portion of real estate assistance requests. However, obtaining reliable information on new developments is challenging due to limited public data and delayed updates. To address this, we leverage digital property brochures to extract essential details, including developer information, floor plans, unit availability, pricing, nearby amenities (e.g., schools, hospitals, shopping centers), transportation options, and comprehensive analyses of project strengths and weaknesses. Given the highly visual nature of these brochures, we propose MMAgent, a multimodal retrieval-augmented generation agent tailored for property-related QA tasks. MMAgent integrates two complementary techniques: (1) A vision-language model (e.g., GPT-4V) interprets images and floor plans, converting visual content into textual descriptions; (2) A multimodal embedding model encodes textual and visual information into vector representations, enabling efficient retrieval within a RAG-based framework. This hybrid approach allows MMAgent to effectively utilize property brochure materials, ensuring accurate and contextually rich responses to user queries.

**TaskAgent.** TaskAgent is designed to manage action-oriented requests within real estate conversations. For example, after multiple dialogue turns, prospective buyers may request to speak with a real estate agent or schedule property inspections. Similarly, developers often aim to capture contact information from high-intent buyers to generate quality leads. In PropGenie, TaskAgent autonomously detects user intent and facilitates seamless task execution. When a user requests a meeting, TaskAgent automatically generates an online meeting link and sends a calendar invitation, streamlining subsequent interactions. For lead generation, TaskAgent dynamically triggers a lead capture form at optimal moments—specifically when user satisfaction and strong buying intent are detected. This adaptive process is further enhanced by RouteAgent, which analyzes user sentiment to balance efficiency and user experience, ensuring interactions remain effective and engaging.

**ChatAgent.** To manage out-of-scope queries, we introduce ChatAgent, which directly leverages an LLM to generate open-domain conversational re-

Tasks	Rel.	Inf.	Cor.	Tasks	Rel.	Inf.	Cor.
Property Search	3.22	3.42	4.15	Home-buying Policy QA	4.28	3.39	4.74
Home Loan Recommendation	3.31	3.86	4.59	Interest Rate QA	4.37	3.76	4.22
Stamp Duty Calculation	4.67	4.05	3.50	Property Project QA	4.27	3.86	4.59
Land Tax Calculation	4.56	3.78	3.31	Task Execution	3.15	2.36	3.83
Repayment Calculation	3.82	3.15	3.37	Open-domain Chitchat	3.82	2.88	4.37

Table 2: Automatic evaluation results of PropGenie. “**Rel.**”, “**Inf.**”, “**Cor.**” represent **Relevance**, **Informativeness**, and **Correctness** correspondingly. Detailed justifications for each metric can be found in Appendix A.2.

sponses. ChatAgent serves two primary purposes: (1) as a fallback mechanism—when task-specific agents cannot produce valid responses, ChatAgent provides default replies, enhancing user experience and facilitating intent clarification; and (2) as a knowledge supplement—by utilizing the LLM’s inherent knowledge and reasoning capabilities, ChatAgent effectively addresses open-ended inquiries, overcoming limitations of specialized agents. This design ensures fluid and engaging interactions, thereby improving the robustness of the overall system.

### 3 System Evaluation

**Automatic Evaluation.** Following prior work (Bi et al., 2023), we evaluate PropGenie’s responses using three metrics: 1) **Relevance** – alignment with user intent and semantic consistency; 2) **Informativeness** – inclusion of detailed explanations and relevant domain-specific context; and 3) **Correctness** – factual accuracy, adherence to predefined formulas (e.g., tax calculations), and absence of hallucinations. To facilitate large-scale evaluation, we adopt the LLM-as-a-Judge paradigm (Zheng et al., 2023; Li et al., 2024a), employing GPT-5 as evaluator. Our test set comprises 6,398 real user queries from online logs, covering the 10 tasks listed in Table 1. Responses are rated from 1 (lowest) to 5 (highest), with concise justifications provided by the evaluator. Domain-specific knowledge (e.g., tax formulas, project background) is incorporated into evaluation prompts to ensure robustness and reproducibility. Further details on the test set and evaluation prompts are in the Appendix A.1&A.2.

Table 2 summarizes the automatic evaluation results. Key observations include: 1) QA-related tasks (e.g., home-buying policy QA, property project QA, interest rate QA) achieve high Correctness and Relevance scores, demonstrating RagAgent and MMAgent’s effectiveness in retrieving and applying domain knowledge. 2) Search and recommendation (S&R) tasks achieve Correctness

scores above 4, highlighting the benefit of integrating real-time APIs for property listings and financial products. 3) Calculation-based tasks (e.g., stamp duty, land tax, monthly repayment) show high Relevance but moderate Correctness, indicating accurate topic adherence but slight accuracy variations in complex reasoning scenarios. 4) Task execution scenarios yield the lowest Informativeness scores, likely due to efficiency-driven responses (e.g., meeting links, lead forms) appearing less detailed. Overall, automatic evaluation confirms PropGenie’s capability to consistently generate relevant, informative, and accurate responses, validating its effectiveness as a conversational assistant for real-world real estate applications.

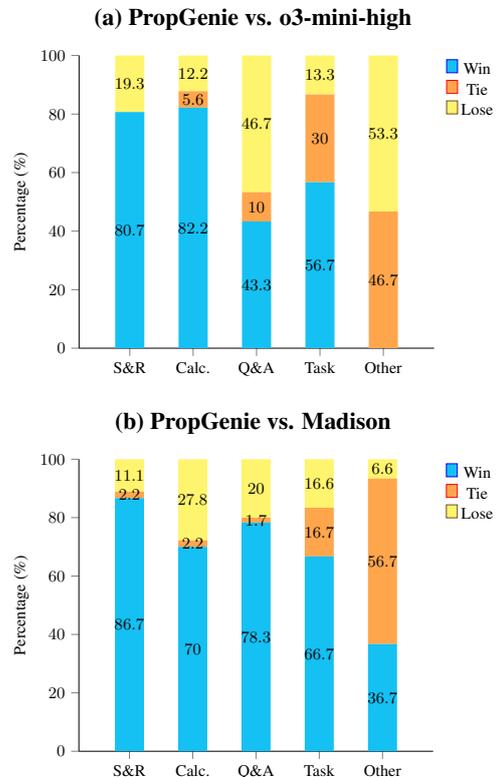


Figure 3: Human evaluation results. S&R, Calc., Q&A, Task, Other represent Search&Recommendation, Calculation, Question Answering, Task Execution, and Chitchat respectively.

**Human Evaluation.** To further validate PropGenie’s effectiveness, we conducted a human evalu-

Model	Generated Response
Question	I am looking for a house in Camberwell, 3124, specifically from or related to Woodards.
PropGenie	Here are some recommended properties for sale, all houses located in Camberwell related to Woodards agency: 1) <b>Superb Single Level Living</b> (19 Callanish Rd), price approximately \$1.87M; 2) <b>Sunnyside Estate without heritage overlay</b> (43 Currajong Ave), priced at \$4.95M.
o3-mini-high	Could you clarify what you mean by “from or related to Woodards”, such as whether you’re looking for properties listed by the Woodards Real Estate agency in Camberwell (3124), interested in houses built, designed, or associated with a builder or developer named Woodards, or if you have specific requirements like property type, budget range, number of bedrooms, or other features to help narrow down the recommendations, so I can assist you in finding the right property recommendations in Camberwell?
Madison	There are multiple houses in Camberwell, 3124 with at least 3 bedrooms and at most 2 bathrooms. I recommend checking various listings or contacting a local real estate agent for the current availability.

Table 3: Comparing responses provided by different models. We convert rich text to plain text for demonstration.

ation comparing it against two competitive baselines: 1) OpenAI’s general-purpose reasoning LLM **o3-mini-high** (OpenAI, 2025), and 2) Realty AI’s domain-specific real estate chatbot **Madison**<sup>3</sup>. We randomly sampled 500 queries proportionally from the original 6,398-query test set, preserving the distribution. Five real estate experts performed a blind evaluation, independently ranking responses from the three systems (inter-annotator agreement: 0.82). Figure 3 summarizes PropGenie’s comparative **win/tie/lose** rates, grouping the ten tasks into five categories for clarity. Results show PropGenie achieves >55% win rates over the general-purpose LLM in search and recommendation, calculation, and task execution, and >65% win rates over the domain-specific baseline in all categories except “Other”. This highlights the importance of integrating external domain knowledge, typically absent in general-purpose models. However, PropGenie shows no clear advantage in QA or chitchat over o3-mini-high, as housing policy is widely available online and general models excel at detailed reasoning. Moreover, our brief-response setting reduces latency but can limit informativeness.

**Case Study.** Table 3 presents a representative property search example comparing model outputs. PropGenie accurately retrieves relevant property listings with precise location and pricing details. In contrast, **o3-mini-high** misunderstands the query and requests clarification, while **Madison** provides only generic advice to search online. Both competitive baselines fail to deliver direct answers due to limited domain-specific knowledge. Additional examples are provided in the Appendix A.4 & A.5.

## 4 Related works

**Real Estate Virtual Assistants.** Prior studies have developed virtual assistants for real estate to provide online customer support (Cao and Nguyen,

2021; Haurum et al., 2024; Yang et al., 2024b; Gloria et al., 2025), capture leads (Quan et al., 2018), offer legal advice (Pagar, 2024), and reduce administrative overhead (Tanović and Hasibović, 2024). However, these systems typically rely on traditional dialogue frameworks or single-agent LLMs, limiting their effectiveness in complex, multi-faceted scenarios. In contrast, we propose a multi-agent framework capable of addressing diverse real estate tasks, including property search, financial planning, and home-buying assistance.

**LLM-based Multi-Agent Systems.** Recent advancements in LLMs have driven the adoption of multi-agent systems, enabling inter-agent communication and collaborative problem-solving with improved accuracy and efficiency over single-agent approaches (Dong et al., 2024; Wu et al., 2023; Li et al., 2024b; Wang et al., 2025b; Hong et al., 2024). Domain-specific multi-agent frameworks have been explored in e-commerce (Thakkar and Yadav, 2024; Fang et al., 2024), legal analysis (Cui et al., 2024), finance (Fatemi and Hu, 2024), healthcare (Tang et al., 2023), and software engineering (Yang et al., 2024a). Unlike previous works, our research introduces a multi-agent conversational system tailored for the real estate domain, leveraging agent collaboration to enhance efficiency and decision-making in property-related tasks.

## 5 Conclusion

In this paper, we introduce PropGenie, a multi-agent conversational framework leveraging large language models for real estate assistance. Eight specialized sub-agents collaboratively handle tasks such as property recommendation, financial calculation, question answering, and open-domain conversation. Experiments on real user queries confirm PropGenie’s effectiveness. Future work includes extending to loan eligibility assessment and automated property valuation.

<sup>3</sup><https://www.realty-ai.com/>

## References

- 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.
- Tuan-Dung Cao and Quang H. Nguyen. 2021. *An Approach for Building Effective Real Estate Chatbots in Vietnamese*, pages 221–229. Springer International Publishing, Cham.
- Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi, Minlie Huang, and Furu Wei. 2024. [Instruction pre-training: Language models are supervised multitask learners](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2529–2550, Miami, Florida, USA. Association for Computational Linguistics.
- Jiaxi Cui, Munan Ning, Zongjian Li, Bohua Chen, Yang Yan, Hao Li, Bin Ling, Yonghong Tian, and Li Yuan. 2024. [Chatlaw: A multi-agent collaborative legal assistant with knowledge graph enhanced mixture-of-experts large language model](#). *Preprint*, arXiv:2306.16092.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, and et al. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.
- Xiaofei Dong, Xueqiang Zhang, Weixin Bu, Dan Zhang, and Feng Cao. 2024. A survey of llm-based agents: Theories, technologies, applications and suggestions. In *2024 3rd International Conference on Artificial Intelligence, Internet of Things and Cloud Computing Technology (AIoTC)*, pages 407–413. IEEE.
- Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Li Fei-Fei, and Jianfeng Gao. 2024. [Agent ai: Surveying the horizons of multimodal interaction](#). *Preprint*, arXiv:2401.03568.
- J. Elster. 2016. *Sour Grapes*. Cambridge Philosophy Classics. Cambridge University Press.
- Jiabao Fang, Shen Gao, Pengjie Ren, Xiuying Chen, Suzan Verberne, and Zhaochun Ren. 2024. A multi-agent conversational recommender system. *arXiv preprint arXiv:2402.01135*.
- Sorouralsadat Fatemi and Yuheng Hu. 2024. Enhancing financial question answering with a multi-agent reflection framework. In *Proceedings of the 5th ACM International Conference on AI in Finance*, pages 530–537.
- Dawei Gao, Zitao Li, Xuchen Pan, Weirui Kuang, Zhi-jian Ma, Bingchen Qian, Fei Wei, Wenhao Zhang, Yuexiang Xie, Daoyuan Chen, Liuyi Yao, Hongyi Peng, Zeyu Zhang, Lin Zhu, Chen Cheng, Hongzhu Shi, Yaliang Li, Bolin Ding, and Jingren Zhou. 2024a. [Agentscope: A flexible yet robust multi-agent platform](#). *Preprint*, arXiv:2402.14034.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024b. [Retrieval-augmented generation for large language models: A survey](#). *Preprint*, arXiv:2312.10997.
- GeminiTeam, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, and et al. 2024. [Gemini: A family of highly capable multimodal models](#). *Preprint*, arXiv:2312.11805.
- Benedikt Gloria, Johannes Melsbach, Sven Bienert, and Detlef Schoder and. 2025. [Real-gpt: Efficiently tailoring llms for informed decision-making in the real estate industry](#). *Journal of Real Estate Portfolio Management*, 31(1):56–72.
- Kasper Raupach Haurum, Ruiqi Ma, and Wen Long. 2024. [Real estate with ai: An agent based on langchain](#). *Procedia Computer Science*, 242:1082–1088. 11th International Conference on Information Technology and Quantitative Management (ITQM 2024).
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiaowu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2024. [MetaGPT: Meta programming for a multi-agent collaborative framework](#). In *The Twelfth International Conference on Learning Representations*.
- Susan Hudson-Wilson, Jacques Gordon, Frank Fabozzi, Mark Anson, and S. Giliberto. 2005. [Why real estate?](#) *The Journal of Portfolio Management*, 31:12–21.
- 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.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiquan Liu. 2024a. [Llms-as-judges: A comprehensive survey on llm-based evaluation methods](#). *Preprint*, arXiv:2412.05579.
- Junyou Li, Qin Zhang, Yangbin Yu, QIANG FU, and Deheng Ye. 2024b. [More agents is all you need](#). *Transactions on Machine Learning Research*.

- OpenAI. 2025. [Openai o3-mini system card](#).
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and et al. 2024a. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and et al. 2024b. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.
- Surbhi Pagar. 2024. [Lawbot : From documents to answers, unveiling a new era in real estate legal assistance](#). *Interantional Journal of Scientific Research In Engineering And Management*, 08:1–5.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. 2024. [ToolLLM: Facilitating large language models to master 16000+ real-world APIs](#). In *The Twelfth International Conference on Learning Representations*.
- Tho Quan, Trung Trinh, Dang Ngo, Hon Pham, Long Hoang, Hung Hoang, Thanh Thai, Phong Vo, Dang Pham, and Trung Mai. 2018. [Lead engagement by automated real estate chatbot](#). In *2018 5th NAFOSTED Conference on Information and Computer Science (NICS)*, pages 357–359.
- Philip Seagraves. 2023. [Real Estate Insights: Is the AI revolution a real estate boon or bane?](#) *Journal of Property Investment & Finance*, 42(2):190–199.
- Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. 2023. [Medagents: Large language models as collaborators for zero-shot medical reasoning](#). *arXiv preprint arXiv:2311.10537*.
- A. Tanović and A. Čerimagić Hasibović. 2024. [Automated real estate chatbot](#). In *2024 47th MIPRO ICT and Electronics Convention (MIPRO)*, pages 241–246.
- SIMA Team, Maria Abi Raad, Arun Ahuja, Catarina Barros, Frederic Besse, Andrew Bolt, Adrian Bolton, Bethanie Brownfield, Gavin Buttimore, Max Cant, Sarah Chakera, Stephanie C. Y. Chan, Jeff Clune, and et al. 2024. [Scaling instructable agents across many simulated worlds](#). *Preprint*, arXiv:2404.10179.
- Param Thakkar and Anushka Yadav. 2024. [Personalized recommendation systems using multi-modal, autonomous, multi agent systems](#). *Preprint*, arXiv:2410.19855.
- Fahim Ullah, Samad M. E. Sepasgozar, and Changxin Wang. 2018. [A systematic review of smart real estate technology: Drivers of, and barriers to, the use of digital disruptive technologies and online platforms](#). *Sustainability*, 10(9).
- Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, LinZheng Chai, Zhao Yan, Qian-Wen Zhang, Di Yin, Xing Sun, and Zhoujun Li. 2025a. [MAC-SQL: A multi-agent collaborative framework for text-to-SQL](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 540–557, Abu Dhabi, UAE. Association for Computational Linguistics.
- Junlin Wang, Jue WANG, Ben Athiwaratkun, Ce Zhang, and James Zou. 2025b. [Mixture-of-agents enhances large language model capabilities](#). In *The Thirteenth International Conference on Learning Representations*.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2023. [Autogen: Enabling next-gen llm applications via multi-agent conversation](#). *Preprint*, arXiv:2308.08155.
- John Yang, Carlos Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024a. [Swe-agent: Agent-computer interfaces enable automated software engineering](#). *Advances in Neural Information Processing Systems*, 37:50528–50652.
- Shuling Yang, Hanzhu Chen, and Binbin Fang. 2024b. [Qudial: A quadruple-driven dialogue system for real estate consulting services](#). In *Proceedings of the 2024 16th International Conference on Machine Learning and Computing, ICMLC '24*, page 609–615, New York, NY, USA. Association for Computing Machinery.
- Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui Zhang, and Sercan O Arik. 2024. [Chain of agents: Large language models collaborating on long-context tasks](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Jun Zhao, Can Zu, Hao Xu, Yi Lu, Wei He, Yiwen Ding, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. [Longagent: Scaling language models to 128k context through multi-agent collaboration](#). *Preprint*, arXiv:2402.11550.
- 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. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *Preprint*, arXiv:2306.05685.

## A Appendix

In the appendix, we present supplementary evaluation data statistics, the prompt template for automatic evaluation, latency and cost analysis, and additional case studies.

### A.1 Dataset Statistics

As mentioned Section 3, our test set comprises 6,398 user interaction samples collected from system logs after deployment. Table 4 presents their distribution across different task categories, reflecting the real-world traffic distribution.

Task	Count
Property Search	374
Home Loan Recommendation	1206
Stamp Duty Calculation	557
Land Tax Calculation	472
Repayment Calculation	307
Home-buying Policy QA	641
Interest Rate QA	196
Property Project QA	152
Task Execution	169
Open-domain Chitchat	2324

Table 4: Task Distribution of Test Set

### A.2 Prompts for LLM-as-a-Judge and LLM-based Reranking

This section describes the prompts used for the LLM-as-a-Judge approach (Zheng et al., 2023; Li et al., 2024a) in automatic evaluation. Since relevance, informativeness, and correctness represent distinct evaluation criteria, we designed three separate prompts, as illustrated in Figure 6 to Figure 8, to independently assess each metric. Additionally, certain tasks require domain-specific knowledge beyond the inherent capabilities of LLMs. To address this limitation, we append relevant domain-specific information to the prompts when evaluating tasks such as stamp duty estimation, land tax calculation, and property project question answering.

Furthermore, Figure 9 illustrates the prompt employed for LLM-based reranking in property search as described in SearchAgent, where the LLM is guided to rank candidate listings based on how well the property features match the user’s search criteria. The prompt employed in RecAgent is similar, except that we additionally instruct the LLM to generate pros and cons of home loan products dur-

ing reranking. We omit this prompt here to avoid duplication.

### A.3 Latency and Cost Analysis

In this section, we analyze PropGenie’s latency and cost, which are critical for user experience and commercial viability. Table 5 summarizes the average response time and cost per query, computed using the same 500 queries from the human evaluation. The latency ranges from 2 to 7 seconds, averaging 3.6 seconds. Property Search, Property Project QA, and Home Loan Recommendation exhibit higher latency due to third-party API calls and multimodal retrievers. In the future, we plan to implement streaming and caching mechanisms to further reduce latency. Additionally, the average cost per query remains below 1.6 US cents, highlighting the token-efficiency advantage of our multi-agent framework, which avoids maintaining all context information within a single prompt.

### A.4 Additional Examples of PropGenie

Due to space constraints, we present only a single case study in the main content. In this section, we provide additional examples of our system to further illustrate its advantages and limitations. In Figures 4 and 5, we supplement additional examples of scenarios involving property search, home loan recommendations, calculations, question answering, task execution, and open-domain chitchat, highlighting PropGenie’s main capabilities. We have included screenshots to illustrate the user interface of our system. Please refer to our online video<sup>5</sup> for more interesting examples.

### A.5 Further Comparisons between PropGenie and Baselines

In this section, we present additional comparisons between PropGenie and two baseline models (OpenAI o3-mini-high and Realty AI’s Madison) to illustrate the strengths and limitations of each model. Table 6 compares the responses generated by the three models for queries related to stamp duty calculation, property project QA, home loan recommendations, and interest rate QA. As shown, PropGenie consistently generates more specific, informative, and helpful answers to effectively address user queries, demonstrating the effectiveness of our multi-agent framework and our strategy of integrating diverse knowledge sources.

<sup>4</sup><https://platform.openai.com/docs/pricing>

<sup>5</sup><https://youtu.be/prafKomKN3g>

Tasks	Latency	Cost	Tasks	Latency	Cost
Property search	4.92	1.47	Home-buying Policy QA	3.55	1.58
Home Loan Recommendation	3.85	1.49	Interest Rate QA	2.87	1.27
Stamp Duty Calculation	2.73	1.19	Property Project QA	6.93	1.58
Land Tax Calculation	2.98	1.39	Task Execution	1.86	1.32
Repayment Calculation	3.46	1.43	Open-domain Chitchat	2.74	1.52

Table 5: Average latency (seconds) and cost (US cents) computed over 500 queries. Costs are calculated based on token usage and OpenAI API pricing<sup>4</sup>.

**(a) Property Search**

**(b) Home loan recommendation**

**(c) Calculations**

Figure 4: Examples of (a) Property Search, (b) Home Loan Recommendation, (c) Calculations (monthly repayment calculation, stamp duty calculation, and land tax calculation).

**(d) Question and Answering**

**(e) Task execution & Chitchat**

Figure 5: Examples of (d) Question and Answering (interest rate QA, home-buying policy QA, and property project QA), (e) Task execution and Chitchat.

### Prompt for Judging Relevance

**Role: Expert Evaluator of Relevance**

Objective: Assess the relevance of responses in direct relation to specific queries, focusing strictly on how well the response addresses the content and intent of the question without introducing unrelated information.

Guidelines for Evaluation:

**Directness:** The response should provide a direct and unambiguous answer to the question posed.

**Topical Alignment:** Content should be closely aligned with the central topic and any associated subtopics of the question.

**Exclusion of Irrelevant Information:** The response must avoid introducing information that is not pertinent to the question.

Instructions:

Comprehend the User's Question:

Read carefully to understand the intent and scope.

Analyze the Response Thoroughly:

Evaluate how directly the response addresses the question.

Assign a Relevance Score (1-5):

1: Irrelevant – The response does not address the question at all.

2: Slightly Relevant – Minimal relevance; the response barely touches on the question's subject.

3: Somewhat Relevant – Partially addresses the question but lacks completeness.

4: Mostly Relevant – Generally addresses the question but may omit minor aspects.

5: Highly Relevant – Fully and directly addresses all aspects of the question.

Justify the Rating:

Provide specific references to elements of both the question and the response that influenced your assessment.

Evaluation Template:

User's Question: [Insert Question Here]

Response: [Insert Response Here]

Your Evaluation:

Relevance Score (1-5):

Justification:

Figure 6: Prompt of Judging Relevance for Automatic Evaluation.

### Prompt for Judging Informativeness

**Role: Expert Evaluator of Informativeness**

Objective: Evaluate the quality and depth of information provided in the response concerning the user's question, focusing on the comprehensiveness and added value of the information supplied.

Guidelines for Evaluation:

**Comprehensiveness:** The response should thoroughly cover all relevant aspects of the question.

**Depth of Information:** Provide detailed explanations, evidence, or examples where appropriate.

**Clarity and Precision:** Information should be clear, precise, and free from ambiguity.

**Added Value:** Offer insights or information that enhance understanding beyond basic or common knowledge.

Instructions:

Understand the User's Question:

Identify the informational needs implied by the question.

Analyze the Response for Informational Content:

Assess the richness and depth of the information provided.

Assign an Informativeness Score (1-5):

1: Not Informative – Provides little to no useful information.

2: Slightly Informative – Offers minimal information with limited depth.

3: Moderately Informative – Provides basic information but lacks depth or detail.

4: Informative – Offers substantial information with good depth and detail.

5: Highly Informative – Comprehensive and provides in-depth, detailed information.

Justify the Rating:

Reference specific parts of the response that contribute to its informativeness.

Evaluation Template:

User's Question: [Insert Question Here]

Response: [Insert Response Here]

Your Evaluation:

Informativeness Score (1-5):

Justification:

Figure 7: Prompt of Judging Informativeness for Automatic Evaluation.

### Prompt for Judging Correctness

**Role: Expert Evaluator of Correctness**

Objective: Assess the accuracy and factual precision of the response in relation to the user's question or task, ensuring that all information presented is correct and reliable. Guidelines for Evaluation:

**Accuracy:** The response should correctly address the question with precise facts and calculations.

**Factual Precision:** Verify the correctness of facts, data, and the application of relevant laws or information.

**Clarity and Accuracy:** Information should be presented clearly, without errors or misrepresentations.

**Relevance:** Ensure that all information provided is pertinent to the question and contributes to a correct understanding.

Instructions:

Review the User's Question or Task:

Understand the specific requirements for correctness.

Analyze the Response for Correctness:

Identify any inaccuracies or errors.

Assign a Correctness Score (1-5):

1: Very Poor – Major inaccuracies; the response is incorrect.

2: Poor – Significant inaccuracies affecting overall correctness.

3: Fair – Some inaccuracies that moderately impact the correctness.

4: Good – Mostly correct with minor inaccuracies.

5: Excellent – Entirely correct with accurate facts and calculations.

Justify the Rating:

Cite specific inaccuracies or confirm areas of correctness in the response.

Evaluation Template:

User's Input: [Insert Input Here]

Response: [Insert Response Here]

Your Evaluation:

Correctness Score (1-5):

Justification:

Figure 8: Prompt of Judging Correctness for Automatic Evaluation.

### Prompt of LLM-based Re-ranking for Property Search

You are an AI virtual assistant specialized in property consultation, helping clients find suitable properties (e.g., apartments, houses, townhouses).

Given a client's query, your task is to rerank a provided list of properties based on how closely each property matches the client's stated conditions.

All the properties in the provided list have been referenced below and numbered with the format "[x].", where "x" is an integer starting from 0.

Carefully evaluate each property against the client's query. Rank properties higher if they match more of the client's explicitly stated conditions. Be conservative and precise—do not infer or assume conditions not explicitly mentioned in the client's query.

Provide your response strictly in the following JSON format:

```
""
{
  "Ranking": [x0, x1, ...],
  "Reason": "Clearly explain your logic for reranking the properties based on the client's query.",
}
```

where "x0", "x1", ... are the reference numbers of the first, second ranked properties and so on from the provided list after reranking.

Below show all the properties from Property List:

{Property\_list}

Figure 9: Prompt of LLM-based Re-ranking for Property Search.

Model	Generated Response
<b>Q1: How much stamp duty should I pay for a 1 million dollar house, considering it is my first home purchase, located in Victoria, and I am an Australian citizen?</b>	
<b>PropGenie</b>	You should pay \$55,000 in stamp duty for a 1 million dollar house as a first home buyer in VIC.
<b>o3-mini-high</b>	Stamp duty in Australia is a state-based tax—that means the amount you’ll pay on a \$1 million property depends on the state or territory where the property is located as well as on whether you qualify for any first home buyer concessions. (Even citizenship isn’t the key factor; rather, it’s your residency status and whether the home will be your principal place of residence.) Here are a couple of examples from two of the larger markets... (omitted for brevity)
<b>Madison</b>	As a first home buyer in Victoria, if your property value exceeds \$600,000, standard stamp duty rates apply. For a \$1 million home, the stamp duty would be calculated without a first home buyer exemption.
<b>Q2: What is the price for apartment DO01.04 in the Yarra Bend project?</b>	
<b>PropGenie</b>	The price of apartment DO01.04 in the Yarra Bend project is \$605,000.
<b>o3-mini-high</b>	Sorry, I couldn’t find any publicly available information about the exact price of apartment DO01.04 in the Yarra Bend project at the moment. To obtain the most accurate and up-to-date quote, it’s recommended to contact the project’s sales representative directly or refer to official sources such as the project’s website or promotional materials. If you need help looking up more information, please feel free to let me know.
<b>Madison</b>	I couldn’t find specific pricing information for the DO01.04 apartment in the Yarra Bend project. Would you like me to connect you with an agent for more detailed assistance?
<b>Q3: For my home loan with a principal and interest repayment type and a fixed interest rate for 2 years, what bank options do I have?</b>	
<b>PropGenie</b>	Based on your request for a loan recommendation, we have found several options that may suit your needs. Below is a summary of the available loan products, including their interest rates, monthly repayments, and key advantages and disadvantages. <b>1) ING Fixed Rate Home Loan - 2 Year Fixed Rate LVR ≤ 80% (Owner Occupier, P&amp;I with Orange Advantage)</b> - Interest Rate: 6.29%; Monthly Repayment: \$6,183.21 - Pros: \$0 ongoing fees; Make additional repayments up to \$10,000 per settlement anniversary before incurring fees - Cons: \$499 application fee; \$250 discharge fee when you close out your home loan; No offset account; More Information: <a href="https://www.finder.com.au/ing-fixed-rate-home-loan#05cd2274-8bc8-4153-b403-0fc9776b0c27">https://www.finder.com.au/ing-fixed-rate-home-loan#05cd2274-8bc8-4153-b403-0fc9776b0c27</a> <b>2) ME Flexible Home Loan Fixed with Members Package - 2 Year Fixed Rate LVR ≤ 70% (Owner Occupier, P&amp;I) ...</b> (omitted for brevity)
<b>o3-mini-high</b>	I’m not a licensed financial adviser, so please consider the following as a starting point for your research rather than definitive recommendations. The availability of a principal-and-interest home loan with a 2-year fixed interest period can vary depending on factors such as your credit profile, loan amount, property type, and—importantly—your country of residence... (omitted for brevity)
<b>Madison</b>	I currently don’t have access to live interest rate options or specific bank offerings. I suggest checking with local banks and mortgage brokers in Melbourne for their current mortgage products. If you need further assistance, would you like to be put in touch with an agent for personalized guidance?
<b>Q4: What is the interest rate of CBA Bank?</b>	
<b>PropGenie</b>	Based on the data collected from the market, the interest rates offered by CBA bank range from 6.24% to 6.74%.
<b>o3-mini-high</b>	The interest rates offered by CBA Bank vary depending on different products, services, and market conditions. For example, interest rates for home loans, personal loans, fixed-term deposits, and savings accounts differ from each other, and these rates may fluctuate according to policy adjustments and market changes. I need to know specifically which type of product’s interest rate you are inquiring about... (omitted for brevity)
<b>Madison</b>	I’m unable to provide specific interest rate information for CBA Bank. You can visit their official website or contact them directly to obtain the latest interest rate details.

Table 6: Comparing responses provided by different models. Non-essential details omitted for brevity.

# Pro-QuEST: A Prompt-chain based Quiz Engine for testing Specialized Technical Product Knowledge

Sujatha Das Gollapalli,<sup>1</sup> Mouad Hakam,<sup>1</sup> Mingzhe Du,<sup>1,2</sup>  
See-Kiong Ng,<sup>1</sup> Mohammed Hamzeh<sup>3</sup>

<sup>1</sup>Institute of Data Science, National University of Singapore

<sup>2</sup>College of Computing and Data Science, Nanyang Technological University

<sup>3</sup>Cisco Systems, Inc. Austin, TX, U.S.A.

{idssdg,mouad.hk,mingzhe,seekiong}@nus.edu.sg, mhamzeh@cisco.com

## Abstract

As large language models (LLMs) rapidly evolve and proliferate, technology companies such as *Cisco* face the difficult challenge of selecting the most suitable model for downstream tasks that demand deep, domain-specific product knowledge. Specialized benchmarks can not only inform this decision making but also be leveraged as quizzes to effectively train engineering and marketing personnel on novel product offerings in a continually growing *Cisco* product space.

We present Pro-QuEST, our **Prompt-chain based Quiz Engine** using state-of-the-art LLMs for generating multiple-choice questions on **Specialized Technical** products. In Pro-QuEST, we first identify key terms and topics from a given professional certification textbook or product guide, and generate a series of multiple-choice questions using domain-knowledge guided prompts. We show LLM benchmarking results with the question benchmarks generated by Pro-QuEST using a range of latest open-source, and proprietary LLMs and compare them with expert-crafted exams and review questions to derive insights on their composition and difficulty. Our experiments indicate that though there is room for improvement in Pro-QuEST to generate questions of the complexity levels seen in expert-designed certification exams, question-type based prompts provide a promising direction to address this limitation. In sample user studies with *Cisco* personnel, Pro-QuEST was received with high optimism for its practical usefulness in quickly compiling quizzes for self-assessment on knowledge of novel products in the rapidly changing tech sector.

## 1 Motivation

Large Language Models (LLMs) have emerged as a transformative technology for various tasks resulting in their current wide-adoption across several

technology industries (Raza et al., 2025; Palen-Michel et al., 2024; Company, 2023). Though LLMs demonstrate excellence at tasks requiring general language understanding such as text analysis, content generation, and summarization, their capabilities and limits for knowledge-intensive domains such as finance, engineering, cybersecurity, and healthcare is still a subject of active research (Fei et al., 2024; Xie et al., 2024; Ouyang et al., 2024). In particular, though Retrieval Augmented Generation (RAG) and knowledge integration (Song et al., 2025; Lewis et al., 2020) have helped in addressing limitations such as hallucination and content grounding, state-of-the-art LLMs still fall behind on tasks requiring complex, novel or multi-step reasoning, where tacit or proprietary knowledge is required, and where contexts and prior experience inform decision making (Chen et al., 2024; Yang et al., 2025; Xu et al., 2025; Kim et al., 2025).

Concurrent with the above research, newer LLMs are being released frequently each with unique architectures, and capabilities, and fine-tuned for specific capabilities (Xiao et al., 2025; Rizzatti, 2025; Wang et al., 2025a). In this changing landscape of LLMs and ongoing research on the promise and limitations of LLMs for specific domains, industry players have to make model choices under economic constraints (Howell et al., 2023). Against this context, standardized benchmarks which quantify LLMs' capabilities via precise performance metrics, characterize knowledge contamination, and provide guidance on making informed choices comprise **crucial** assets for a company. Indeed, both LLM benchmarking and benchmark generation now form core topics of active investigation in several domains (Fei et al., 2024; Ouyang et al., 2024; Xie et al., 2024).

In this study, we investigate the creation of domain-specific LLM benchmarks for *Cisco*, a technology company providing thousands of net-

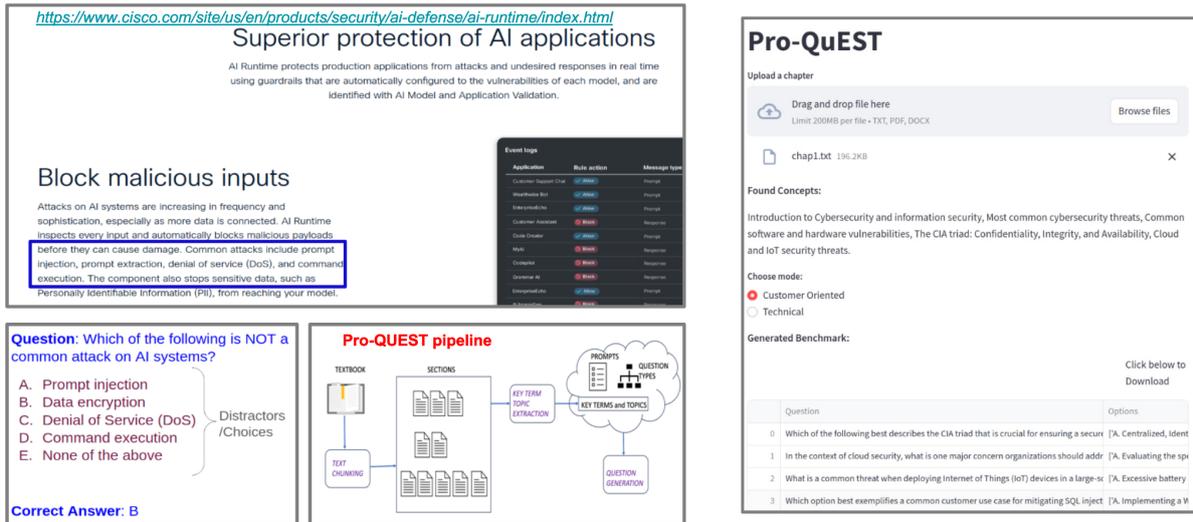


Figure 1: Our main task is illustrated with a sample input section and an LLM-generated multiple-choice question. The modules comprising the Pro-QuEST-pipeline and a screenshot of our web-based demo are also shown.

working products and services to customers across the world. *Cisco* contends with other competitors such as *Arista Networks*, *Dell Technologies*, and *Huawei* for market share pertaining to various device and software offerings in networking and cybersecurity, a rapidly developing sector. Therefore, not only is it critical for company marketing and sales professionals to be well-aware of the various features of devices from *Cisco* and competition, customer support engineers need to be familiar with intricate operational features to provide effective troubleshooting for clients. At the same time, as newer LLMs are made available, a quantitative assessment of their capabilities on internal tasks such as service request handling, question answering, and predictive analytics becomes inevitable (Yang et al., 2023; Tanhaei et al., 2024; Gollapalli et al., 2025; Saini et al., 2025).

*How can we create representative benchmarks that measure specialized domain knowledge for guiding LLM selection for internal tasks and for training personnel on the complex and rapidly growing product mix in Cisco?* We developed Pro-QuEST, our **P**rompt-chain based **E**ngine for generating benchmark **Q**uizzes on the **S**pecialized, **T**echnical knowledge in *Cisco* to address this question. In this demo paper, we describe the core components of Pro-QuEST and our experiments with generated benchmarks on content from *Cisco* textbooks used by technical personnel preparing for entry-level and advanced-level certifications (CCNA 200-301<sup>3</sup> and CCNP 350-701<sup>4</sup>). We highlight Pro-QuEST’s exciting potential for

training sales and marketing personnel on the novel product offerings with respect to their facts and features and provide insightful experiments on improving automatic benchmark creation for highly specialized, technical domains with SOTA LLMs. A **web-based demo** of Pro-QuEST was showcased at a recent Cisco Live event<sup>1</sup> and is available at <https://nlp-demos.online/qg/> with an illustrative **video** available at [https://github.com/mouad157/Cisco-benchmark/blob/main/EACL\\_ProQuEST.mp4](https://github.com/mouad157/Cisco-benchmark/blob/main/EACL_ProQuEST.mp4).

## 2 System Components

Question answering (QA) datasets are commonly used in LLM benchmarking studies for specialized domains since QA accuracy provides a quantifiable measure of “knowledge” of a specific domain (Fei et al., 2024; Ouyang et al., 2024; Xie et al., 2024; Chen et al., 2025). Following our objective to create QA datasets for *Cisco*, we follow recent works (Choi et al., 2025; Xiong et al., 2024; Camarata et al., 2025; Dalvi et al., 2024) and apply document grounded multiple-choice question generation using SOTA LLMs in Pro-QuEST. An anecdotal illustration of our main task and our processing pipeline can be found in Figure 1 along with a screenshot of our web-based interface.

Pro-QuEST uses **prompt chaining** (Wu et al., 2022; Sun et al., 2024), a widely used technique in LLMs to break down complex tasks into a series of simpler tasks by using the output of one

<sup>1</sup><https://www.ciscolive.com/apjc.html>

prompt as the input to the next prompt in the chain.<sup>2</sup> Prompt chains reduce the “cognitive load” for an LLM through explicit instructions on the steps involved in solving complex tasks and were shown to improve output quality and reduce hallucination through context retention between prompts. Considering context limitations in LLMs and the often lengthy nature of input documents (such as textbooks) that represent “knowledge”, we accomplish three tasks in Pro-QuEST through a prompt chain as follows:

**1. Section Chunking:** We use LLM prompts on the first few pages of a long *text* document to identify the document type (such as product guides, textbooks, research papers, configuration matrix documents in *Cisco*), as well as other metadata information. These prompts aim to extract content organization information in the input document (for example, “table of contents”). The identified section headings or chapter titles are used to split a lengthy input document into smaller cohesive text chunks for further processing.

**2. Key Terms and Topics Identification:** From the sections identified in the previous step, we identify and collate the topical keyphrases and overarching themes using LLM prompts. Extraction of topical keyphrases is a widely-studied topic in NLP due to their effectiveness in representing and summarizing vast amounts of information from lengthy documents (Boudin and Aizawa, 2025). Indeed, keyphrases are widely used for various retrieval, analytical, and organizational tasks as well as to ground question generation (Willis et al., 2019; Wang et al., 2020; Zhang and Zhu, 2021).

**3. Multiple-Choice Question Generation:** Finally, the content and keyphrases from the previous two steps are combined with a diverse list of LLM prompts to generate multiple-choice questions (MCQs), the answer options or *distractors* for the questions, and the lists of correct answers. MCQs are prevalent in Education as well as LLM benchmarking since they can be designed for various levels of learning complexity and allow for an efficient, quantitative assessment (Camarata et al., 2025; Jovanovska, 2018).

Our above pipeline ensures coverage of all main topics of a lengthy document. When LLM prompts fail to extract sections in Step-1 (for example, when a content listing is missing), we first identify “sec-

tion headings” using a heuristic algorithm that couples stylistic cues along with section length thresholds. For example, most words in section header sentences are capitalized and their average length is smaller than a typical sentence in the main body.<sup>7</sup> Our zero-shot LLM prompt templates are included in Tables 5, 10, and 11 of the Appendix.

In ongoing research, techniques such as chain-of-thought reasoning (Sprague et al., 2025), and in-context learning (Dong et al., 2024) are being employed to generate complex questions in specific domains such as Finance and Medicine (Choi et al., 2025; Liang et al., 2023). Such **overt** question design knowledge from *Cisco* experts was not available to us. We therefore focus on MCQ generation with simple prompts and compare them with available expert-compiled questions to derive insights that can inform future prompt design. Model-generated questions, regardless of the complexity of prompts employed, need expert validation for specialized domains. While this human validation is in progress, in this paper, we provide quantitative evaluation by characterizing question answering performance and comparing generated benchmarks with the available expert-compiled questions for their composition and difficulty.

### 3 Experiments

**Datasets:** The datasets for developing and testing Pro-QuEST were provided by *Cisco* and include two textbooks that are official preparation guides for the certification exams: (1) Cisco Certified Network Associate, an **entry-level** certification covering foundational networking skills (CCNA 200-301<sup>3</sup>) and (2) Cisco Certified Network Professional, an **advanced** certification for professionals for operating core security technologies (CCNP 350-701<sup>4</sup>). Both textbooks are long and image-heavy documents containing 29 and 11 content chapters, respectively. In this study, we only focused on the textual content, and sampled four chapters from each textbook for experiments. We refer to the expert-designed review questions available with each chapter from these textbooks with the label “Book” in our experiments. The textbooks also contained expert-specified key terms that we used to evaluate Step-2 of our Pro-QuEST-pipeline (Section A.1).

<sup>3</sup><https://www.oreilly.com/library/view/ccna-200-301-official/9780136755562/>

<sup>4</sup><https://www.oreilly.com/library/view/ccnp-and-ccie/9780138221287/>

<sup>2</sup>[https://www.promptingguide.ai/techniques/prompt\\_chaining](https://www.promptingguide.ai/techniques/prompt_chaining)

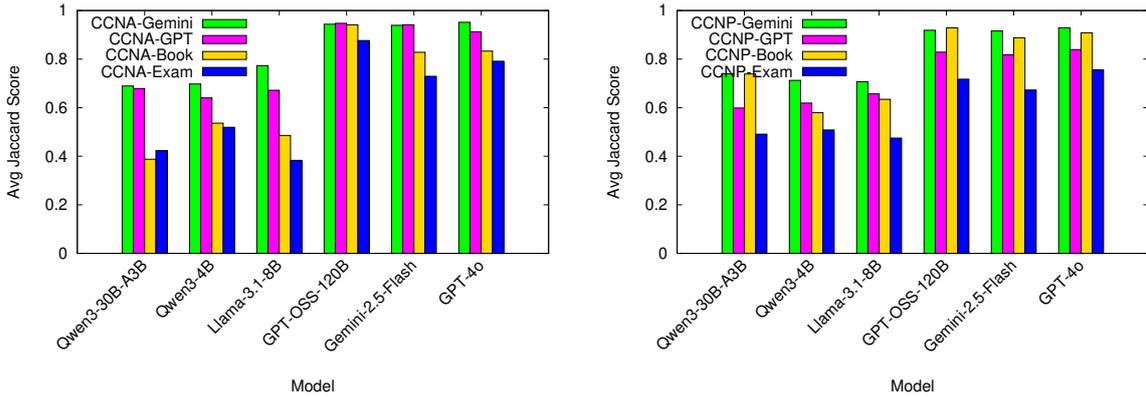


Figure 2: Results of LLM Benchmarking. The average Jaccard scores for CCNA and CCNP datasets are shown.

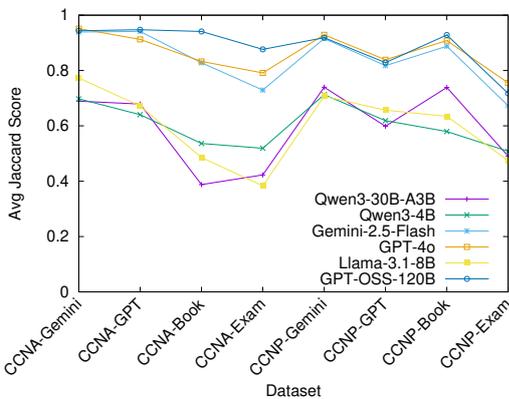


Figure 3: Model QA performance on our datasets.

Cisco also provided access to the question banks of the two certification exams. Unlike textbooks, this content is proprietary making it unlikely to be pre-trained knowledge for SOTA LLMs. Exam and textbook questions are created by domain experts and cover the spectrum of fundamental to complex topics (“CCNA versus CCNP”) as well as difficulty range (“Book versus Exam”). A summary of our datasets is provided in Table 1. The LLM-generated questions using the sections from CCNA and CCNP textbooks are indicated using suffixes ‘-GPT’ and ‘-Gemini’.

CCNA-Book	187	CCNP-Book	49
CCNA-Exam	400	CCNP-Exam	572
CCNA-GPT	165	CCNP-GPT	159
CCNA-Gemini	148	CCNP-Gemini	169

Table 1: Question datasets are shown with the number of questions in each set

We investigated keyphrase/keytopic extraction and question generation in Pro-QuEST using state-

of-the-art LLMs—(1) GPT-4o from OpenAI,<sup>5</sup> and (2) Gemini-2.5-Flash.<sup>6</sup> Our choice of models was influenced by the available best performing, versatile models at the time of experiments, as well as pricing and context length considerations. For LLM benchmarking experiments, we selected small to large, open-source and proprietary models: *Qwen3-4B*, *Qwen3-30B-A3B*, *Llama-3.1-8B*, *GPT-OSS-120B*, *GPT-4o* and *Gemini-2.5-Flash*. We include their details in Table 8 of the Appendix. Our code and prompts are shared on GitHub<sup>7</sup> for research purposes with further details on dataset processing and experimental settings included in the Appendix.

### 3.1 LLM Benchmarking Results

**Question Answering Performance:** We evaluated a range of recent LLM models on all our datasets from Table 1 on the Question Answering (QA) task. QA performance was measured using Jaccard accuracy that measures the set overlap between predicted answers (‘A’) and the correct answers (‘B’) as  $\frac{|A \cap B|}{|A \cup B|}$ .

As can be noticed in the performance plots of Figures 2 and 3, QA accuracies of the smaller models (from Qwen and Meta) are significantly lower than that of the much larger proprietary models as well as the 120B parameter model from OpenAI (*GPT-OSS-120B*). This is not surprising since larger LLMs which have several scales higher numbers of parameters can be expected to “know” more and demonstrate higher QA performance. We note that the QA performance is consistently, significantly higher for model generated questions (\*-

<sup>5</sup><https://openai.com/api/>

<sup>6</sup><https://aistudio.google.com/>

<sup>7</sup><https://github.com/mouad157/Cisco-benchmark>

GPT and \*-Gemini datasets) compared to the Exam datasets suggesting that LLM-generated benchmark questions may be easier to answer than expert-designed benchmark questions.

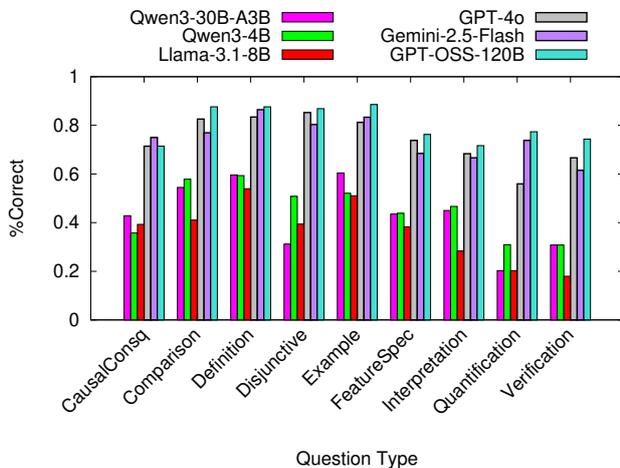


Figure 4: QA performance by question type

**Question Difficulty:** We grouped questions from our datasets into the following difficulty categories: *easy*—answered correctly by all small models, *medium*—answered incorrectly by all small models and correctly by majority of the large models, and *hard*—answered incorrectly by all large models. In Table 2, the percentages of *easy/medium/hard* questions as per the above definitions are illustrated for all the datasets. Unlike Exams where the *easy* questions comprise 10-16% of the datasets, the LLM-generated benchmarks have a considerably higher 30-40% *easy* questions. Similarly, the questions in the *hard* category are considerably higher 9-18% in Exams compared to 1-12% in the LLM-generated benchmarks. The numbers of *easy* and *hard* questions in Books seem to represent averages of these other two sources.

**Question Type Analysis:** We further analyzed benchmark questions by employing a question-type taxonomy available from prior works (Zhao and Jiang, 2010; Nielsen et al., 2008). Question-types represent the nature of information sought in the answer to a given question. For instance, a ‘Definition’-type question asks for how a given concept may be defined (Example: “In a LAN, which of the following terms best equates to the term VLAN?”) whereas a ‘Quantification’ question asks about quantitative aspects of a situation (Example: “What is the maximum number of distribution switches that can be deployed within a hierarchical LAN design building block?”). A list of example

multiple-choice questions from Cisco datasets for our twelve question types as well as details of question type prediction are included in Section A.3 of the Appendix.

We show the QA performance of various models grouped by question types in Figure 4. Similar to earlier experiments, smaller models are significantly worse than the larger LLMs across the question types. Indeed, performance with smaller models is particularly limited on *Interpretation/Quantification/Disjunctive* question types.<sup>8</sup> These categories are arguably challenging for the larger models as well since the QA performance on these types is lower compared to types such as *Definition/Comparison/Example*.

Overall the *GPT-OSS-120B* model, the latest, largest open-source offering from OpenAI which is also a reasoning model, closely outperforms *GPT-4o* and *Gemini-2.5-Flash* on all question types but one. We would like to highlight that for certain question types, such as *Quantification/Interpretation*, it is highly likely that chain-of-thought style complex prompts yield better results (Sprague et al., 2025). Moreover, companies such as OpenAI have multiple LLM offerings designed for specific use-cases. In this study, we consider LLMs designed for overall versatility and employ simple QA/QG prompts which can be employed across all LLMs uniformly (Table 6). We posit that this setting is more reflective of LLM’s role as a “stand-in exam taker”.

The question type spreads for CCNP datasets are shown in Figure 5 with those for CCNA included in the Appendix. The question type “*Feature Specification*” dominates the Exam benchmark and occurs only half as frequently in the GPT-generated benchmark, whereas the opposite is the case for the “*Definition*” question type. Given the significantly higher QA performance for this latter type, it is not surprising that all LLMs uniformly under-perform on the Exam questions in the benchmarking experiments (Figures 2 and 3).

We conducted experiments on incorporating specific question type into the LLM prompts (Tables 10 and 12 in the Appendix). While initial, anecdotal results with type-augmented prompts seem promising, this research and expert evalua-

<sup>8</sup>We observe here that in addition to the actual answer, in a considerable number of cases, smaller models do not follow prompt guidance with respect to output format and add explanations and reasoning process, despite explicit directions not to, resulting in errors during output parsing.

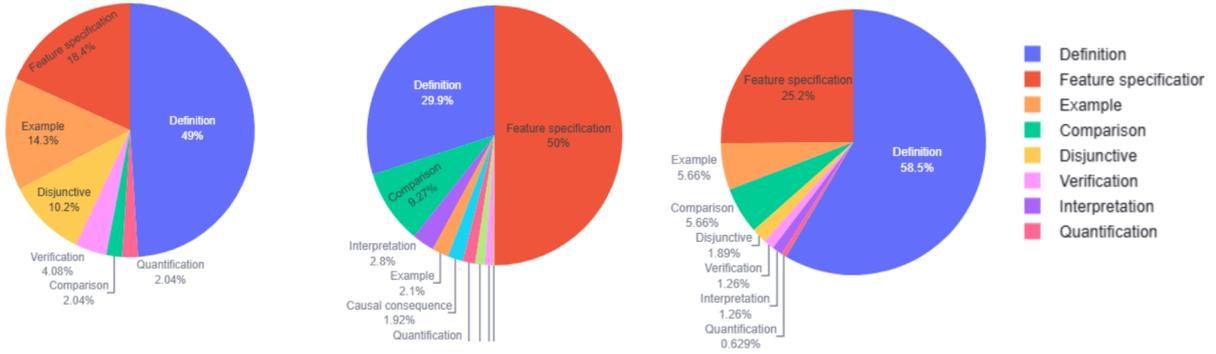


Figure 5: Question Type Distribution for CCNP-Book (left), CCNP-Exam (center), and CCNP-GPT (right)

tion of generated questions is a topic for our future study.

**User Study Findings:** Within *Cisco*, questions of types such as ‘*Causal Consequence/Interpretation*’ are reflective of scenarios faced by tech-support engineers who need intricate knowledge during troubleshooting. In contrast, recall and fact-oriented questions pertaining to the types ‘*Definition/Feature Specification/Comparison/Example*’ correspond to knowledge used by marketing and sales executives who need to keep abreast with information on the rapidly developing novel product offerings from *Cisco* as well as competition.

Relying on experts to design training materials for every newly innovated product at varying complexity levels would be time-consuming and costly. This scenario provides the perfect venue for applying LLM-generated questions. We incorporated a quiz-style interface on top of Pro-QuEST and showcased it along with a leader board at the recent Cisco Live event<sup>1</sup> held in Melbourne for engaging our potential users. Two sets of quizzes based on content on recent *Cisco* products and general content from Wikipedia articles (listed in Table 4) were used for this demonstration. Screenshots from our quiz interface along with the leaderboard are shown in Figure 7. During the event, we presented random samples of ten LLM-generated questions as our quizzes and scored the participants based on their choice of correct answer and speed. The time limit for each question was set to 45 seconds per question and prizes were offered to participants who topped the leader boards.<sup>9</sup> Overall, we had six and ten participants for the “technical/non-technical” topic quizzes, respectively. The average score was  $\sim 270$  for the former and  $\sim 355$  for the latter, in line with

the expectation that technical questions are more challenging to answer than the non-technical ones. Our quiz generator tool based on Pro-QuEST was well-received during the event and has opened up connections for real deployment within *Cisco*.

Dataset	easy%	med.%	hard%
CCNA/Gemini	39.19	2.7	1.35
CCNA/GPT	32.12	5.45	2.42
CCNP/Gemini	42.60	1.18	4.14
CCNP/GPT	34.59	6.92	11.95
CCNA/Book	10.16	9.63	6.42
CCNP/Book	30.61	4.08	4.08
CCNA/Exam	10.5	12.75	9.25
CCNP/Exam	15.91	5.25	18.88

Table 2: Question percentages for difficulty levels

## 4 Related Work

LLM-based approaches are now state-of-the-art for various internal tasks within companies involving information processing and language generation including automated text correction, summarization, question answering, entity recognition, product reviews evaluation, and customer support chatbots (Palen-Michel et al., 2024; Wulf and Meierhofer, 2024; Zheng et al., 2023; Roumeliotis et al., 2024; Su et al., 2025; Oh, 2024; Song et al., 2021). As LLMs are being rapidly adopted and still evolving, benchmarking has become an active topic of recent research. Representative benchmarks can characterize LLM model capabilities on domain-specific tasks such as reasoning, conversations, programming (Lu et al., 2021; Lin et al., 2022; Srivastava et al., 2022; Chiang et al., 2024), languages (Dalvi et al., 2024; Baucells et al., 2025) as well as aspects such as factuality and hallucination (Bao

<sup>9</sup>Assuming on average a person would need at least 15 seconds to read and answer, the best possible score is  $\sim 667$ .

et al., 2025; Wang et al., 2025b). Question answering (QA) datasets are widely used for LLM benchmarking studies (Zhong et al., 2020; Fei et al., 2024; Ouyang et al., 2024; Xie et al., 2024; Chen et al., 2025; Guha et al., 2023). Several research works have addressed the creation of specialized QA datasets using LLMs for domains such as law, medicine, and finance (Choi et al., 2025; Xiong et al., 2024; Scaria et al., 2024; Artsi et al., 2024; Camarata et al., 2025) but, to our knowledge, we are the first to investigate MCQG with LLMs in a highly-technical industry context (such as Cisco).

## 5 Conclusions

We presented Pro-QuEST, our system for generating quizzes for technology companies such as Cisco who operate with highly specialized domain knowledge. Our experiments with Pro-QuEST-questions illustrated their practical usefulness for LLM benchmarking as well as provided insights on future QG studies on the topic. We evaluated Pro-QuEST by trialing it with in-house sales personnel at a recent marketing event at Cisco. Pro-QuEST was enthusiastically received for its potential to efficiently create training quizzes for keeping up with the rapidly-evolving product landscape in networking and security markets.

In future, we would like to study the transferability of Pro-QuEST to other technical product domains such as embedded systems, and sensor technologies. Though it was not observed in the samples manually examined in this study, there is a possibility of generated questions and keyphrases to be incorrect, involve hallucination, and be overall unusable. Qualitative evaluation of LLM-generated benchmarks and design of more accurate classification models and taxonomies for question-type characterization comprise some of our future research directions.

## Acknowledgments

We thank Sarah Yee from Cisco Systems for her initiative and assistance with setting up the user study at the Cisco Live, Melbourne event in 2025.

This research is supported by A\*STAR, CISCO Systems (USA) Pte. Ltd and National University of Singapore under its Cisco-NUS Accelerated Digital Economy Corporate Laboratory (Award I21001E0002).

## References

- Y. Artsi, V. Sorin, E. Konen, B. S. Glicksberg, G. Nadkarni, and E. Klang. 2024. *Large language models for generating medical examinations: systematic review*. *BMC medical education*.
- Forrest Sheng Bao, Miaoran Li, Renyi Qu, Ge Luo, Erana Wan, Yujia Tang, Weisi Fan, Manveer Singh Tamber, Suleman Kazi, Vivek Sourabh, Mike Qi, Ruixuan Tu, Chenyu Xu, Matthew Gonzales, Ofer Mendelevitch, and Amin Ahmad. 2025. *FaithBench: A diverse hallucination benchmark for summarization by Modern LLMs*. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 448–461, Albuquerque, New Mexico. Association for Computational Linguistics.
- Irene Baucells, Javier Aula-Blasco, Iria de Dios-Flores, Silvia Paniagua Suárez, Naiara Perez, Anna Salles, Susana Sotelo Docio, Júlia Falcão, Jose Javier Saiz, Robiert Sepulveda Torres, Jeremy Barnes, Pablo Gamallo, Aitor Gonzalez-Agirre, German Rigau, and Marta Villegas. 2025. *IberoBench: A benchmark for LLM evaluation in Iberian languages*. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10491–10519, Abu Dhabi, UAE. Association for Computational Linguistics.
- Florian Boudin and Akiko Aizawa. 2025. *An analysis of datasets, metrics and models in keyphrase generation*. In *Proceedings of the Fourth Workshop on Generation, Evaluation and Metrics (GEM<sup>2</sup>)*, pages 973–973, Vienna, Austria and virtual meeting. Association for Computational Linguistics.
- Troy Camarata, Lise McCoy, Robert Rosenberg, Kelsey R. Temprine Grellinger, Kylie Brettschnieder, and Jonathan Berman. 2025. *Llm-generated multiple choice practice quizzes for preclinical medical students*. *Advances in Physiology Education*, 49(3):758–763.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. *Benchmarking large language models in retrieval-augmented generation*.
- Rubing Chen, Jiaxin Wu, Jian Wang, Xulu Zhang, Wenqi Fan, Chenghua Lin, Xiaoyong Wei, and Li Qing. 2025. *Benchmarking for domain-specific LLMs: A case study on academia and beyond*. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, Suzhou, China. Association for Computational Linguistics.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anatasios N. Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael I. Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. *Chatbot arena: an open platform for evaluating llms by human preference*. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.

- Chanyeol Choi, Jihoon Kwon, Jaeseon Ha, Hojun Choi, Chaewoon Kim, Yongjae Lee, Jy yong Sohn, and Alejandro Lopez-Lira. 2025. [Finder: Financial dataset for question answering and evaluating retrieval-augmented generation](#). *Preprint*, arXiv:2504.15800.
- McKinsey Company. 2023. [The economic potential of generative ai: The next productivity frontier](#). *McKinsey Company*.
- Fahim Dalvi, Maram Hasanain, Sabri Boughorbel, Basel Mousi, Samir Abdaljalil, Nizi Nazar, Ahmed Abdelali, Shammur Absar Chowdhury, Hamdy Mubarak, Ahmed Ali, Majd Hawasly, Nadir Durani, and Firoj Alam. 2024. [LLMeBench: A flexible framework for accelerating LLMs benchmarking](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 214–222, St. Julians, Malta. Association for Computational Linguistics.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. [A survey on in-context learning](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.
- Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Alan Huang, Songyang Zhang, Kai Chen, Zhixin Yin, Zongwen Shen, Jidong Ge, and Vincent Ng. 2024. [LawBench: Benchmarking legal knowledge of large language models](#).
- Sujatha Das Gollapalli, Mouad Hakam, Mingzhe Du, See-Kiong Ng, and Mohammed Hamzeh. 2025. [On assigning product and software codes to customer service requests with large language models](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 1092–1103, Suzhou (China). Association for Computational Linguistics.
- Neel Guha, Julian Nyarko, and Others. 2023. [Legal-bench: a collaboratively built benchmark for measuring legal reasoning in large language models](#). In *Proceedings of the 37th International Conference on Neural Information Processing Systems*.
- Kristen Howell, Gwen Christian, Pavel Fomitchov, Gitit Kehat, Julianne Marzulla, Leanne Rolston, Jadin Tredup, Ilana Zimmerman, Ethan Selfridge, and Joseph Bradley. 2023. [The economic trade-offs of large language models: A case study](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 248–267, Toronto, Canada. Association for Computational Linguistics.
- Jasmina Jovanovska. 2018. [Designing effective multiple-choice questions for assessing learning outcomes](#). *Infotheca - Journal for Digital Humanities*, 18(1).
- Jonathan Kim, Anna Podlasek, Kie Shidara, Feng Liu, Ahmed Alaa, and Danilo Bernardo. 2025. [Limitations of large language models in clinical problem-solving arising from inflexible reasoning](#). *Scientific Reports*, 15(1).
- 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 *NeurIPS*.
- Yuanyuan Liang, Jianing Wang, Hanlun Zhu, Lei Wang, Weining Qian, and Yunshi Lan. 2023. [Prompting large language models with chain-of-thought for few-shot knowledge base question generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4329–4343, Singapore. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [TruthfulQA: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, MING GONG, Ming Zhou, Nan Duan, Neel Sundaresan, and 3 others. 2021. [CodeXGLUE: A machine learning benchmark dataset for code understanding and generation](#). In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Rodney Nielsen, Jason Buckingham, Gary Knoll, Ben Marsh, and Leysia Palen. 2008. [A taxonomy of questions for question generation](#).
- Jangmin Oh. 2024. [Developing a model for extracting actual product names from order item descriptions using generative language models](#). *IEEE Access*, 12:122695–122701.
- Zetian Ouyang, Yishuai Qiu, Linlin Wang, Gerard De Melo, Ya Zhang, Yanfeng Wang, and Liang He. 2024. [CliMedBench: A large-scale Chinese benchmark for evaluating medical large language models in clinical scenarios](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8428–8438, Miami, Florida, USA. Association for Computational Linguistics.
- Chester Palen-Michel, Ruixiang Wang, Yipeng Zhang, David Yu, Canran Xu, and Zhe Wu. 2024. [Investigating llm applications in e-commerce](#). *Preprint*, arXiv:2408.12779.
- Mubashar Raza, Zarmina Jahangir, Muhammad Bilal Riaz, Muhammad Jasim Saeed, and Muhammad Awais Sattar. 2025. [Industrial applications of large language models](#). *Scientific Reports*.

- Lauro Rizzatti. 2025. [A closer look at llm’s hyper growth and ai parameter explosion](#).
- Konstantinos I. Roumeliotis, Nikolaos D. Tselikas, and Dimitrios K. Nasiopoulos. 2024. [Llms in e-commerce: A comparative analysis of gpt and llama models in product review evaluation](#). *Natural Language Processing Journal*, 6:100056.
- Harsh Saini, Md Tahmid Rahman Laskar, Cheng Chen, Elham Mohammadi, and David Rossouw. 2025. [LLM evaluate: An industry-focused evaluation tool for large language models](#). In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, Abu Dhabi, UAE.
- Nicy Scaria, Suma Dharani Chenna, and Deepak Subramani. 2024. [Automated Educational Question Generation at Different Bloom’s Skill Levels Using Large Language Models: Strategies and Evaluation](#), page 165–179. Springer Nature Switzerland.
- Shuangyong Song, Chao Wang, Haiqing Chen, and Huan Chen. 2021. [An emotional comfort framework for improving user satisfaction in E-commerce customer service chatbots](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Papers*, pages 130–137, Online. Association for Computational Linguistics.
- Zirui Song, Bin Yan, Yuhan Liu, Miao Fang, Mingzhe Li, Rui Yan, and Xiuying Chen. 2025. [Injecting domain-specific knowledge into large language models: A comprehensive survey](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 25297–25311, Suzhou, China. Association for Computational Linguistics.
- Zayne Rea Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2025. [To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning](#). In *The Thirteenth International Conference on Learning Representations*.
- Aarohi Srivastava and 1 others. 2022. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *arXiv preprint arXiv:2206.04615*.
- Hanchen Su, Wei Luo, Yashar Mehdad, Wei Han, Elaine Liu, Wayne Zhang, Mia Zhao, and Joy Zhang. 2025. [LLM-friendly knowledge representation for customer support](#). In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pages 496–504, Abu Dhabi, UAE. Association for Computational Linguistics.
- Shichao Sun, Ruifeng Yuan, Ziqiang Cao, Wenjie Li, and Pengfei Liu. 2024. [Prompt chaining or step-wise prompt? refinement in text summarization](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7551–7558, Bangkok, Thailand. Association for Computational Linguistics.
- Hamed Tanhaei, Payam Boozary, Sogand Sheykhani, Maryam Rabiee, Farzam Rahmani, and Iman Hosseini. 2024. [Predictive analytics in customer behavior: Anticipating trends and preferences](#). *Results in Control and Optimization*, 17:100462.
- Bingning Wang, Xiaochuan Wang, Ting Tao, Qi Zhang, and Jingfang Xu. 2020. [Neural question generation with answer pivot](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:9138–9145.
- Suqing Wang, Zuchao Li, Shi Luohe, Bo Du, Hai Zhao, Yun Li, and Qianren Wang. 2025a. [From parameters to performance: A data-driven study on LLM structure and development](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, Suzhou, China. Association for Computational Linguistics.
- Yuxia Wang, Minghan Wang, Hasan Iqbal, Georgi N. Georgiev, Jiahui Geng, Iryna Gurevych, and Preslav Nakov. 2025b. [OpenFactCheck: Building, benchmarking customized fact-checking systems and evaluating the factuality of claims and LLMs](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 11399–11421, Abu Dhabi, UAE. Association for Computational Linguistics.
- Angelica Willis, Glenn Davis, Sherry Ruan, Lakshmi Manoharan, James Landay, and Emma Brunskill. 2019. [Key phrase extraction for generating educational question-answer pairs](#). In *Proceedings of the Sixth (2019) ACM Conference on Learning @ Scale*, New York, NY, USA. Association for Computing Machinery.
- Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. [Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts](#). In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI ’22, New York, NY, USA. Association for Computing Machinery.
- Jochen Wulf and Jürg Meierhofer. 2024. [Utilizing large language models for automating technical customer support](#). *Preprint*, arXiv:2406.01407.
- Chaojun Xiao, Jie Cai, Weilin Zhao, Biyuan Lin, Guoyang Zeng, Jie Zhou, Zhi Zheng, Xu Han, Zhiyuan Liu, and Maosong Sun. 2025. [Densing law of llms](#). *Nature Machine Intelligence*.
- Qianqian Xie and 1 others. 2024. [Finben: An holistic financial benchmark for large language models](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. [Benchmarking retrieval-augmented generation for medicine](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6233–6251, Bangkok, Thailand. Association for Computational Linguistics.

- Zhijian Xu, Yilun Zhao, Manasi Patwardhan, Lovekesh Vig, and Arman Cohan. 2025. [Can LLMs identify critical limitations within scientific research? a systematic evaluation on AI research papers](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 20652–20706, Vienna, Austria. Association for Computational Linguistics.
- Changlin Yang, Siye Liu, Sen Hu, Wangshu Zhang, Teng Xu, and Jing Zheng. 2023. [Improving knowledge production efficiency with question answering on conversation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 225–234, Toronto, Canada. Association for Computational Linguistics.
- Wenli Yang, Lilian Some, Michael Bain, and Byeong Kang. 2025. [A comprehensive survey on integrating large language models with knowledge-based methods](#). *Knowledge-Based Systems*, 318:113503.
- Zhiling Zhang and Kenny Zhu. 2021. [Diverse and specific clarification question generation with keywords](#). In *Proceedings of the Web Conference 2021, WWW '21*, page 3501–3511, New York, NY, USA. Association for Computing Machinery.
- Jianhua Zhao and Yinjian Jiang. 2010. [Categories of questions in an online discussion forum: An analysis](#). In *2010 5th International Conference on Computer Science Education*, pages 428–431.
- Xin Zheng, Tianyu Liu, Haoran Meng, Xu Wang, Yufan Jiang, Mengliang Rao, Binghuai Lin, Yunbo Cao, and Zhifang Sui. 2023. [DialogQAE: N-to-n question answer pair extraction from customer service chatlog](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6540–6558, Singapore. Association for Computational Linguistics.
- Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020. [Jecqa: A legal-domain question answering dataset](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*.

## A Appendix

### A.1 Datasets

The chapters considered from the CCNA certification guide were “*Introduction to TCP-IP Networking*, *Fundamentals of Ethernet LANs*, *Implementing Ethernet Virtual LANs*, and *Spanning Tree Protocol Concepts*” whereas from the CCNP guide, we considered the chapters on “*Cybersecurity Fundamentals*, *Cryptography Cisco Secure Firewall*, and *Virtual Private Networks (VPNs)*”. For these chapters, the textbook questions correspond to the ones listed under “Do I know this already?” sections of the chapters. For the exam questions from the question bank, we do not have the specific mapping to the topics/chapters. and hence all questions were included in experiments.

The certification guide books also included representative key terms and topics for a given chapter. We evaluated the keyphrases extracted in Step-2 of our pipeline (Figure 1) with LLM prompts (Table 11) against these expert-drafted lists. On average overlap between the keyphrases extracted with GPT-4o and the expert lists overlap about 34% of the time (Table 3). The list of keyphrases from the textbook and GPT-4o are shown for a sample chapter in Table 7.

Mean	Max	STD
0.342	0.857	0.206

Table 3: Overlap between GPT-extracted and expert-crafted keyphrases for chapters in the CCNA certification guide

The documents used for quizzes in the **User Study** are listed in Table 4. Random samples of 10 were administered for each quiz from the overall sets of 30 generated questions for technical/general documents listed above.

### A.2 Models and Parameter Settings

Table 8 presents the list of LLMs used in our benchmarking experiments. Our models include both proprietary and open-source options, with parameter sizes ranging from 3B to 120B. For experiments using open-source LLMs (from Meta and Qwen), we used vLLM for inference on 4×H100 80GB GPUs. For question generation, we extracted 5 keyphrases and 3 keytopics per section (Table 11) and restricted the number of questions to 10 in “free” generation –where LLMs are not constrained by the topic/keyphrase during generation (See **NoKP-MCQP** in Table 10).

#### Technical Topics

1. <https://www.cisco.com/c/dam/en/us/products/collateral/routers/secure-routers/8300-series-secure-routers-ds.pdf>
  2. <https://blogs.cisco.com/security/cisco-hybrid-mesh-firewall-better-enforcement-points-smarter-segmentation-multi-vendor-policy>
  3. [https://www.cisco.com/c/dam/en\\_us/solutions/artificial-intelligence/ai-infrastructure.pdf](https://www.cisco.com/c/dam/en_us/solutions/artificial-intelligence/ai-infrastructure.pdf)
- #### General Topics
1. <https://en.wikipedia.org/wiki/Australia>
  2. <https://en.wikipedia.org/wiki/Melbourne>
  3. <https://en.wikipedia.org/wiki/Cisco>

Table 4: Documents used for our User Study.

### A.3 Question Type analysis

The list of question types (Zhao and Jiang, 2010; Nielsen et al., 2008) with examples from Cisco datasets are shown in Table 9. For a cost-effective and efficient method to obtain question type information over all our data, we trained a local model as follows. First, labels of question types for the CCNA-Book and CCNA-Exam questions were obtained by prompting GPT-4o LLM in a zero-shot setting using an MCQ formulation (“Which of the following question type from the list best matches..”).

We manually checked samples of GPT-assigned labels and found them to be highly accurate. “Silver” data obtained through GPT was used to train a local prompt-tuned model using Flan-T5-large<sup>10</sup>. Model training and inference was performed on a single GPU of an Nvidia Tesla cluster (Linux) machine with 32GB RAM. For assessing the reliability of predicted question types, we manually analyzed random samples of predictions obtained with our trained model for ten questions for each type from CCNP-Exam (Table 1). On this sample, the labels were 73% accurate with a macro-averaged F1 score of 68. Most errors corresponded to ‘Feature Specification’ questions incorrectly predicted as ‘Definition/Causal Consequence/Interpretation’ and ‘Example’ questions incorrectly predicted as ‘Verification’. Due to the direction of these errors and the small percentages of types such as ‘Causal Consequence/Causal/Antecedent/Interpretation’ in our datasets, we posit that our analysis based on the relative trends of dominant question-type frequencies is still legitimate.

<sup>10</sup><https://huggingface.co/google/flan-t5-large>

**System Prompt:** You are an efficient PDF parser. From the initial content of the PDF provided, extract the metadata requested for.

**Sections Prompt:** Does the initial content from a PDF include a listing of topics in in the document? If yes, return a Python list of strings, each string being topic title. Content: []

Table 5: Prompts of extracting metadata/content listing

**System Prompt:** You are a Cisco technical support engineer with in-depth knowledge of Cisco certification materials. Answer the following multiple-choice question. Only print your answers as a Python list . . .

**User Prompt:** Question with Options

Table 6: MCQA Prompts

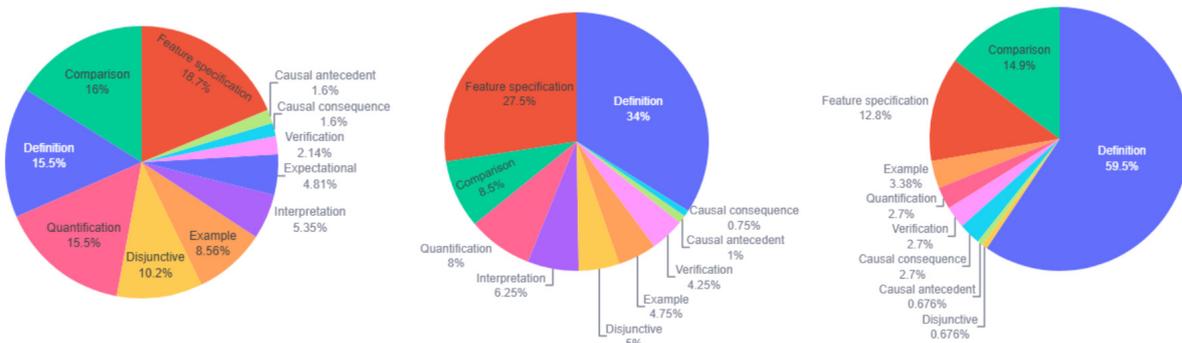


Figure 6: Question Type Distribution for CCNA datasets

### Question 2 of 10

What feature of the HPE Aruba Networking 2930F Switch Series helps organizations in securing IoT devices?

Time remaining: 42 seconds

Select one:

- A. Dynamic segmentation and role-based access control
- B. Limiting device connections to only pre-approved MAC addresses
- C. Disabling all unused ports by default
- D. Using proprietary security protocols unique to each device

Submit Answer

Correct! (+4100.00 points, time bonus applied)

### Quiz Complete!

Your final score: 17300.00

#### Leaderboard

	Name	Score	Questions	Timestamp
1	mouas	21100	10	2025-11-25 17:12:02
2	test	18800	10	2025-11-25 15:01:41
3	test	17300	10	2025-11-25 17:14:25
4	kiki	13400	8	2025-11-25 15:18:52
5	test	8000	10	2025-11-25 16:59:00
6	mouad	7300	10	2025-11-25 15:14:17
7	YYIDS!	6700	10	2025-11-25 16:26:03
8	test	4400	10	2025-11-25 16:15:30
9	mouad	3700	10	2025-11-25 15:07:39

Figure 7: The Pro-QuEST interface shown in Figure 1 is followed by an option to try out the generated questions as a quiz and quiz takers are ranked on a leader board

<p><b>Gold KPs:</b> 'alternate port (role)', 'backup port (role)', 'blocking state', 'BPDU Filter', 'BPDU Guard', 'bridge ID', 'bridge protocol data unit (BPDU)', 'broken state', 'designated port', 'designated port (role)', 'disabled port (role)', 'disabled state', 'discarding state', 'EtherChannel', 'forward delay', 'forwarding state', 'Hello BPDU', 'learning state', 'listening state', 'Loop Guard', 'MaxAge', 'PortFast', 'Rapid STP (RSTP)', 'root cost', 'Root Guard', 'root port (role)', 'root switch', 'Spanning Tree Protocol (STP)', 'superior BPDU', 'unidirectional link'</p>
<p><b>GPT KPs:</b> '<b>Spanning Tree Protocol (STP)</b>', 'Rapid Spanning Tree Protocol (RSTP)', 'IEEE 802.1D', 'IEEE 802.1w', '<b>BPDU (Bridge Protocol Data Units)</b>', 'root switch election', 'root port (RP)', 'designated port (DP)', 'alternate port (ALT)', '<b>backup port</b>', 'port roles', 'port states', 'convergence', '<b>forwarding state</b>', '<b>blocking state</b>', '<b>listening state</b>', '<b>learning state</b>', '<b>BPDU Guard</b>', '<b>Root Guard</b>', '<b>Loop Guard</b>'</p>

Table 7: Key terms from the textbook and GPT-4o extracted key terms are shown for the chapter 9 from CCNA guide/Chapter 9. Spanning Tree Protocol Concepts

Model Name	Type	Size	Link
Qwen3-4B	Open Source	4B	<a href="https://huggingface.co/Qwen/Qwen3-4B-Instruct-2507">https://huggingface.co/Qwen/Qwen3-4B-Instruct-2507</a>
Llama-3.1-8B	Open Source	8B	<a href="https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct">https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct</a>
Qwen3-30B-A3B	Open Source	30.5B (MoE)	<a href="https://huggingface.co/Qwen/Qwen3-30B-A3B-Instruct-2507">https://huggingface.co/Qwen/Qwen3-30B-A3B-Instruct-2507</a>
GPT-OSS-120B	Open Source	117B (MoE)	<a href="https://huggingface.co/openai/gpt-oss-120b">https://huggingface.co/openai/gpt-oss-120b</a>
GPT-4o	Proprietary	Unknown	<a href="https://openai.com/api/">https://openai.com/api/</a>
Gemini 2.5 Flash	Proprietary	Unknown	<a href="https://aistudio.google.com/">https://aistudio.google.com/</a>

Table 8: List of selected LLMs used in our system, including their type (proprietary or open source), parameter size (if available), and source links.

Question Type	Description with Examples
Verification	For yes/no responses to factual questions <i>Example: Which of the following access-list commands, taken from a router's running-config file, match all packets sent from hosts in subnet []</i>
Disjunctive	Questions that require a simple decision between two alternatives. <i>Example: Which of the following things are bound together when a new WLAN is created? (Choose two answers.)</i>
Feature specification	Determines qualitative attributes of an object or situation. <i>Example: The Wi-Fi Alliance offers which of the following certifications for wireless devices that correctly implement security standards?</i>
Quantification	Determines quantitative attributes of an object or situation. <i>Example: What is the maximum number of distribution switches that can be deployed within a hierarchical LAN design building block?</i>
Definition	Determine meaning of a concept. <i>Example: Which one of the following is the data encryption and integrity method used by WPA2?</i>
Example	Request for instance of a particular concept or even type. <i>Example: Which one of the following is an example of a AAA server?</i>
Comparison	Identify similarities and differences between two or more objects. <i>Example: Which answer best compares named standard IP ACLs with numbered standard IP ACLs?</i>
Interpretation	A description of what can be inferred from a pattern of data <i>Example: Upon receipt of a configuration BPDU with the topology change flag set, how do the downstream switches react?</i>
Causal antecedent	Asks for an explanation of what state or event causally led to the current state and why. <i>Example: What is the main reason SD-Access uses VXLAN data encapsulation instead of LISP data encapsulation?</i>
Causal consequence	Asks for explanation of consequences of event/state <i>Example: What happens to a switch port when a BPDU is received on it when BPDU guard is enabled on that port?</i>
Expectational	Asks about expectations or predictions (including violation of expectation) <i>Example: Which action would you expect to be true of a router CLI interaction that is not true . . .</i>
Judgmental	Asks about value placed on an idea, advice, or plan. <i>Example: . . . Which one of the following things should you do to determine the root cause of her problem?</i>

Table 9: Question types (Zhao and Jiang, 2010; Nielsen et al., 2008) with examples from Cisco datasets

<p><b>System Prompt:</b> You are a Cisco technical support engineer with in-depth knowledge of CCNA certification materials.</p> <p><b>MCQPrefix:</b> Generate exactly one multiple-choice question based on the content provided with no explanation. Return only a JSON-tuple= . . .</p> <p><b>MCQP1: MCQPrefix.</b> Use the keyphrase [KP]. Content: []</p> <p><b>MCQP2: MCQPrefix.</b> The keyphrase [KP] must be an answer to the question. Content: []</p> <p><b>MCQP3: MCQPrefix.</b> The keyphrase [KP] must be one of the options for the multiple-choice question . . . Content: []</p> <p><b>NoKP-MCQP:</b> Generate exactly [NUM] multiple-choice questions based on the content . . .</p> <p><b>KT-TYPEDMCQP:</b> Generate exactly [NUM] multiple-choice questions of type [QTYPE] based on the content provided . . . for the topic [TOPIC]</p>
---

Table 10: MCQ Prompts

<p><b>System Prompt for Key Terms:</b> You are a Cisco technical support engineer with in-depth knowledge of CCNA certification materials and are expert in identifying important concepts related to Cisco domain</p> <p><b>System Prompt for Key Topics:</b> . . .expert in identifying key topics in the content involving the principles, algorithms, methods, and techniques related to Cisco domain. For example, some key topics can be described as: “Commands to find access ports and assigned VLANs” . . .</p> <p><b>Prompt-1:</b> For the given passage extract the top-[X%] relevant conceptual keyphrases. Only return your output as a Python list of strings. Passage: [INPUT-PASSAGE]</p> <p><b>Prompt-2:</b> Group the given sets of conceptual keyphrases, and select the top-[X%] most important conceptual keyphrases, given the topic: %s. Only return your output as a Python list of strings. List of keyphrases [OUTPUT-FROM-Prompt-1]</p>
---

Table 11: Keyphrase Prompts

Prompt Type	Sample Generated Question
MCQP1	What does the BPDU Guard feature do when it receives a BPDU on a port configured with PortFast?
MCQP2	Which optional STP feature helps prevent forwarding loops by disabling a port if it receives BPDUs on a port that should only connect to endpoint devices
MCQP3	Which feature disables a port if it receives any BPDUs, helping to prevent forwarding loops when unexpected switches connect to access ports?
KT-MCQP	What is the primary function of BPDU Guard in a network configuration?
KT-MCQP/Quantification	What is the default Hello time interval for BPDU in STP
KT-MCQP/Interpretation	Based on the BPDU Guard logic, what can be inferred when a nonroot switch stops receiving Hello BPDUs on its root port?
KT-MCQP/Causal Antecedent	What could cause a switch to start changing the STP topology?

Table 12: Sample generated questions are shown for the key term: **BPDU Guard** and key topic: **Basic logic for BPDU**

# elfen: A Python Package for Efficient Linguistic Feature Extraction for Natural Language Datasets

Maximilian Maurer

GESIS Leibniz Institute for the Social Sciences

Heinrich-Heine University Düsseldorf

maximilian.maurer@gesis.org

## Abstract

A detailed understanding of the basic properties of text collections produced by humans or generated synthetically is vital for all steps of the natural language processing system life cycle, from training to evaluating model performance and synthetic texts. To facilitate the analysis of these properties, we introduce *elfen*, a Python library for efficient linguistic feature extraction for text datasets. It includes the largest set of item-level linguistic features in eleven feature areas: surface-level, POS, lexical richness, readability, named entity, semantic, information-theoretic, emotion, psycholinguistic, dependency, and morphological features. Building on top of popular NLP and modern dataframe libraries, *elfen* enables feature extraction in various languages (80 at the moment) on thousands of items, even given limited computing resources<sup>1</sup>. We show how using *elfen* enables linguistically informed data selection, outlier detection, and text collection comparison. We release *elfen* as an open-source PyPI package, accompanied by extensive documentation, including tutorials<sup>2</sup>.

## 1 Introduction

While there is a dire need to understand our data at all levels, such as pre-training, fine-tuning, few-shot in-context learning, evaluation, and synthetically generated texts, there is a surprising lack of deep engagement with the basic properties of the data at hand. This is especially worrisome since works on spurious correlations in the data affecting model performance show that long-standing linguistic measures can provide insights into model

behavior (Poliak et al., 2018; Liusie et al., 2022; Borah et al., 2023, inter alia). Similarly, recent works have emphasized the need for and promises of measuring data diversity (Nguyen and Ploeger, 2025). We argue that in line with these efforts, measuring the linguistic composition of texts is highly relevant, especially in the age of generative LLMs: Firstly, measuring linguistic properties of (pre-)training data at scale can give insights into downstream model behavior (e.g., Zhang et al., 2021). Secondly, given the popularity of benchmark datasets to assess improvements of ever bigger models, linguistic features can give insights into the comparability and the shortcomings of benchmark datasets (e.g., Gubelmann et al., 2024). Finally, given the rising prevalence of synthetic data, it becomes more and more important to understand its properties, be it to understand and detect it better (Dönmez et al., 2025) or to assess its utility for (pre-)training, especially given risks like model collapse (Shumailov et al., 2024).

While there are tools for linguistic feature extraction, most of them are focused on a specific area (e.g., lexical richness) or support a limited number and scope of features. Suppose a researcher wants a broad coverage of features in a given analysis. In that case, this causes difficulties, given that different tools require different preprocessing and can differ widely in how efficient they are, especially when dealing with large numbers of instances. There is a clear lack of availability of unified extraction tools providing a comprehensive number of features in different areas, and a way to efficiently extract them.

To fill this gap, we present *elfen*, a Python package to efficiently extract linguistic features for large numbers of text instances. Our contributions are fourfold: (1) We provide a tool with the largest collection of features (1,061), (2) most of which are extractable in 80 languages out of the box. (3) *elfen* provides efficient extraction (on

<sup>1</sup>For instance, the features for popular benchmarks can be extracted on consumer hardware in less than an hour. For more details, see Appendix D.

<sup>2</sup>We host the code at <https://github.com/mmmaurer/elfen/>, make it available through the GESIS Methods Hub at <https://methodshub.gesis.org/library/methods/elfen/>, and provide documentation and tutorials at <https://elfen.readthedocs.io/en/latest/>.

Library	elfen	LFTK	LIWC
Surface-Level	11	9	2
Lexical Richness	26	10	0
Readability	11	6	1
Named Entities	19	19	0
Information Theory	2	0	0
Emotion/Sentiment	36	0	8
POS	20	34	20
Psycholinguistics	78	3	33
Semantics	17	0	41
Morphology	798	0	0
Syntactic Dependencies	43	0	0
Reading Time Formulas	0	3	0
Total	1,061	84	105

Table 1: Comparison of the number of features implemented per feature area for `elfen` (v1.2.4), LFTK (Lee and Lee, 2023), and LIWC (Boyd et al., 2022). We keep the comparison to libraries with the same scope and goal as `elfen`. Due to different design choices regarding normalization by token, lemma, or sentence count, we only consider what they call *foundation* features in LFTK for this comparison. We count all non-*psycholinguistic* dictionary features as *semantic* for LIWC.

average 21.8% faster than comparable libraries) on tens of thousands of items. (4) `elfen` builds on popular libraries, allowing for easy integration into existing workflows and the multilingual coverage of it to grow with them.

In the following, we discuss existing tools and the contributions of `elfen` (Section 2), present the implementation and functionalities of it (Section 3), and showcase already existing and potential use cases of it (Section 4).

## 2 Related Work

Measuring characteristics of texts to compare them has a long history in (computational) linguistics and NLP, from early works trying to measure specific properties like lexical diversity (e.g., Yule, 1944) and readability (e.g., Mc Laughlin, 1969), to more recent advances trying to measure overall semantic similarity between texts (Corley and Mihalcea, 2005; Reimers and Gurevych, 2019)

In consequence, several tools for extracting such features have been developed, some of which are now outdated or no longer actively maintained (Graesser et al., 2004; Simig et al., 2022). Currently available, actively maintained tools can be categorized across two axes: (1) the scope of the features they provide, and (2) the units they operate

on.

Variationist (Ramponi et al., 2024), for instance, mainly measures token, n-gram, and sequence occurrence frequencies on a (sub)corpus level, and, in turn, provides only a few corpus-level features. Conversely, there are tools focusing on text instance-level features, most of which focus on a single group of features, such as lexical richness (Shen, 2022) or readability<sup>3</sup> on a text-level, as well as tools for extracting token-level features such as `wn` (Goodman and Bond, 2021), a package for using open multilingual wordnets (Bond and Foster, 2013). More extensive tools like LFTK (Lee and Lee, 2023) and the commercial provider LIWC (Boyd et al., 2022) aim to cover broader sets of features on a text instance-level. While LFTK focuses on lexical richness and readability measures and frequencies of token-level features as characteristics of a given text, for example, the number of nouns or of named entities, LIWC mainly focuses on frequencies of words associated with certain semantic categories, for example, politics, or emotional or perceptive grounding.

### 2.1 Resource Gap

The fragmented coverage of feature groups, the lack of integration of token-level resources, such as psycholinguistic norms or emotion lexicons in open-source tools, and the fact that they are not optimized for datasets with large numbers of text instances make existing tools suboptimal for conducting analyses on typical NLP tasks and benchmark datasets, pre-training instances, and synthetically generated text collections. The package presented in this paper, `elfen`, thus aims to provide a unified tool for extracting an extensive number of text-level characteristics across feature areas, optimized for modern large NLP datasets. `elfen` provides a more comprehensive coverage across and within most feature areas, as we show in the comparison of `elfen` with LFTK and LIWC in Table 1, while being significantly faster than comparable tools when extracting large amounts of features.

## 3 Implementation and Functionality

### 3.1 Implementation

To allow for extensive analyses of datasets with tens of thousands of text items, we build `elfen` on top of `Polars`<sup>4</sup> for efficient parallel processing, and

<sup>3</sup><https://github.com/andreascv/readability>

<sup>4</sup>[pola.rs](https://pola.rs)

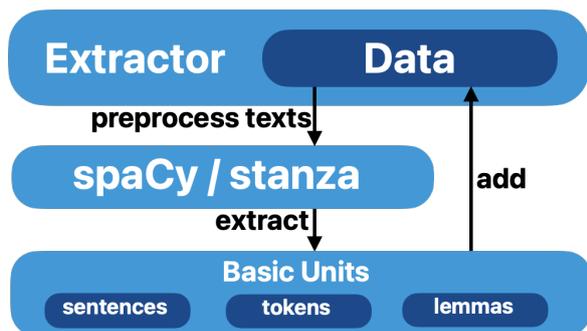


Figure 1: Schematic overview of preprocessing in e1fen.

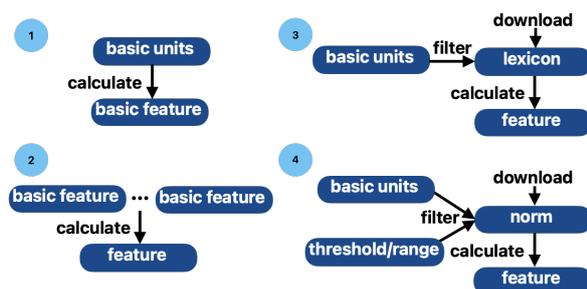


Figure 2: Schematic overview of the four types of feature extraction procedures.

spaCy (Honnibal et al., 2020) and stanza (Qi et al., 2020) to integrate with established NLP tools and pipelines.

As illustrated in Figure 1, e1fen first preprocesses the texts in the dataframe used to initialize the Extractor with the backend model of choice from spaCy or stanza models. The extracted basic units, the sentences, tokens, lemmas, and syllables per text are stored in their respective columns in the dataframe, allowing for efficient downstream calculation of features based on these basic units.

From these basic units, e1fen allows for the extraction of (implementation-wise) four types of features: (1) features directly computable from a basic unit (e.g., number of tokens), (2) features computable from multiple basic features (e.g., type-token ratio), (3) lexicon-based features requiring some basic unit and a lexicon to match the text with (e.g., the number of hedges), and (4) *norm*-based features requiring some basic unit and a *norm*, a lexicon including a measurement of a given property per basic unit (e.g., the *concreteness* of a given token).

As shown in Figure 2, (1) and (2) are directly calculated, utilizing parallel processing abilities from polars. This allows these features to be extracted in less than a second after initializing the Extractor,

Backend	Library	Runtime $\pm$ SD
en_core_web_sm	e1fen	11.43 $\pm$ 0.51
	LFTK	15.46 $\pm$ 0.02
en_core_web_md	e1fen	12.22 $\pm$ 0.26
	LFTK	15.92 $\pm$ 0.64
en_core_web_lg	e1fen	11.99 $\pm$ 0.21
	LFTK	15.77 $\pm$ 0.36
en_core_web_trf	e1fen	13.48 $\pm$ 0.23
	LFTK	15.64 $\pm$ 0.46

Table 2: Mean extraction time in seconds  $\pm$  standard deviation (5 runs) for the first 100 items of the Stanford Sentiment Treebank (Socher et al., 2013) for all available features in e1fen (1,061) and LFTK (220) using the same backend spacy models.

even for datasets with tens of thousands of items. For (3) and (4), if the license permits, the respective lexicons or norms will be automatically downloaded<sup>5</sup> and filtered according to the basic units present. For (4), if applicable, for example, if features based on *highly* concrete tokens are of interest, an additional filtering step is performed on the property measurement. Finally, the number of basic units fulfilling the filter criteria is counted, or a given property of their measurements is calculated (e.g., average concreteness of the tokens).

The optimized extraction of basic units, features, and filtering results in a considerable speedup over existing feature extraction tools. Given the same spaCy backbone models, e1fen extracts all available 1,061 features on average 21.8% faster than the most comparable open-source library, LFTK (Lee and Lee, 2023), extracts all available 220 features (see Table 2).

### 3.2 Use Case-Driven Extraction

The main class of the package, the Extractor, is implemented with a special focus on the ease of use for various analysis scenarios. As illustrated in Figure 3, the Extractor provides the extraction in only a few lines of code of individual features, applicable in cases where researchers are interested in specific features, feature groups, in (the comparison of) a family of features, and all available features in e1fen, for exploratory scenarios or when they

<sup>5</sup>The usage of the norms and lexicons is subject to different licenses. Complying with them and the citation guidelines is the user’s responsibility. Some lexicons will need to be downloaded manually. Further details can be found in the documentation and the repositories’ README.

```

# initializing extractor
extractor = elfen.Extractor(
    data = df,
    language = "en",
    text_column = "text")
# extracting a single feature: ttr
extractor.extract("ttr")
# extracting a feature area/group: readability
extractor.extract_feature_group("readability")
# extracting all available features
extractor.extract_features()

```

Figure 3: Code examples of feature extraction capabilities. The Extractor here is initialized with a polars data frame *df*, which contains English text in the column *text*.

are interested in a comprehensive overview of the instances in a dataset.

### 3.3 Implemented Linguistic Features

elfen implements 1,061 features in eleven broad feature areas. Table 3 describes the feature areas and gives an example of a feature<sup>6</sup>.

### 3.4 Multilingual Support

Given that elfen is using spaCy and stanza for pre-processing, we rely on the availability of language-specific models in them. SpaCy currently has 24 language-specific and one multilingual model available. Stanza provides 138 models in 80 languages.

All of the features except for the psycholinguistic norm-based, emotion, and semantic features are language-independent<sup>7</sup>. The emotion lexicons are available in 108 languages, psycholinguistic features are currently available in English, German, French, and Italian, and semantic features are available in all languages supported by the wn package (currently 34).

### 3.5 Analysis and Processing Utilities

Given elfen’s focus on linguistic analyses of text datasets, in addition to the extraction, we provide useful utilities for downstream analyses. These include extracted feature **rescaling** to specified

<sup>6</sup>For a more detailed description including references for the implemented features, see Appendix A. Our documentation provides additional information on each feature at [https://elfen.readthedocs.io/en/latest/feature\\_overview.html](https://elfen.readthedocs.io/en/latest/feature_overview.html)

<sup>7</sup>We provide a periodically updated overview of which features are available in which language at [https://elfen.readthedocs.io/en/latest/multilingual\\_support.html](https://elfen.readthedocs.io/en/latest/multilingual_support.html)

ranges, **normalization** to have a mean of 0 and a standard deviation of 1, or by the number of tokens, lemmas, or sentences. We also provide functionality to extract local (within a given instance) and global (across the whole dataset) token and lemma frequencies.

## 4 Evaluation

To illustrate the usefulness of elfen, we discuss existing work already using it in three categories. We additionally outline three broad analysis scenarios showcasing how elfen provides insights in LLM-related research.

### 4.1 Existing Work Using elfen

#### 4.1.1 Analysis of Human and LLM Behavior

elfen enables the analysis of human and language model behavior, and its connection to performance. It has, for example, already been used to assess linguistic factors in the human perception of gendered style of texts (Chen et al., 2025). Similarly, Falk and Lapesa (2025) to assess linguistic factors in annotation uncertainty in humans and models on morals and values. While Falk and Lapesa (2025) have a particular focus on human label variation and its connection to model uncertainty, in principle, any model-internal or output-derived metrics could be substituted to assess connections to structural characteristics of the texts. Thus, **elfen facilitates analyses connecting human and language model behavior with linguistic structure.**

#### 4.1.2 Authorship and Stylistic Analysis

Tasks such as authorship attribution and stylistic analysis (e.g., Sari et al., 2018; Ayele et al., 2024) naturally use linguistic features due to their inherent need for interpretability. Zeng et al. (2025) use elfen as one potential interpretable component in their explainable authorship verification method. As this exemplifies, **elfen provides interpretability in tasks where it is vital, and can be integrated in respective systems for such tasks.**

#### 4.1.3 Detection and Analysis of LLM-Generated Text

Another natural avenue of work is the analysis and detection of LLM-generated content. A major limitation of prior works in this line of work is the lack of access to extensive corpus statistics (Wu et al., 2025), which elfen alleviates. Parfenova et al. (2025), for instance, use elfen to

Feature Area	Description	Example
Surface-level	Structural characteristics of a text	Number of tokens
Readability	Reading complexity; how hard to read a text is	Flesch reading ease
Psycholinguistics	Cognitive, social, or sensorimotor groundings of words	Number of highly concrete tokens
POS	Parts-of-speech in the text	Number of nouns
Morphology	Grammatical/lexical properties of words in a text	Number of plurals
Information theory	Redundancy and formulaicity of a text	Shannon entropy
Lexical richness	How lexically diverse is a text	Type-token ratio (TTR)
Syntactic Dependencies	Predicate-argument relations in a text	Number of adverbial modifiers
Semantics	Polysemy and ambiguity	Number of hedges
Named entities	Reference to entities with a proper name	Number of organizations
Emotion/Sentiment	Emotion or sentiment evoking or related words	Number of high arousal tokens

Table 3: Overview of feature areas with example features.

analyse convergence patterns in multi-agent annotation. Similarly, we show `elfen`’s utility for the case of LLM-written arguments, both for extensive analyses and detection scenarios (Dönmez et al., 2025). Thus, **elfen enables light-weight and interpretable detection and analyses of LLM-generated synthetic text and language model behavior.**

## 4.2 Exemplified Use Cases

To further illustrate the broad range of analyses `elfen` enables, we discuss three exemplary analysis steps that may be applied to many use cases: (1) Dataset comparison, (2) linguistically-informed targeted sampling, and (3) outlier detection.

We showcase these use cases<sup>8</sup> on two popular language understanding benchmark datasets, MMLU-Pro (Wang et al., 2024) and BigBench Hard (Suzgun et al., 2022).

### 4.2.1 Dataset Comparison

To understand dataset and domain effects at each step of the NLP life cycle, it is beneficial to understand in depth where datasets differ. This is particularly relevant for the generalization of training and test data: To train and test models for a given task and draw insights on models’ capabilities, especially for benchmarking, we ideally want data to be as diverse as possible to reduce the influence of confounders. Suppose there are multiple datasets for a given task that differ structurally. It may be beneficial to either use the most diverse dataset or use multiple complementary datasets to

<sup>8</sup>The code for the exemplified use cases is available as commented notebooks at <https://github.com/mmmaurer/elfen-examples>. The enriched datasets are available at <https://huggingface.co/collections/mmmaurer/enriched-language-understanding-benchmarks>

get more robust results and more informative insights on model behavior and capabilities.

Two angles of assessment here are (1) a comparison of the overall and feature-area-wise correlation structure for a coarse overview, and (2) a comparison of how individual features are distributed for the datasets to assess fine-grained differences.

For (1), a natural option is to inspect correlation matrix heatmaps and similarity measures between (sub)matrices. As Figure 4 illustrates for BigBench Hard and MMLU-Pro, the correlation structure between two datasets on the same task can differ quite drastically, both overall and in specific areas, like morphological structures. This is reflected in similarity measures between the correlation (sub)matrices. For example, the Mantel correlation (Mantel, 1967) for morphological features (0.430,  $p < 0.001$ ) is substantially lower than for surface-level features (0.925,  $p = 0.005$ )<sup>9</sup>.

Given the particular differences in morphological features, it may be interesting for researchers to look more closely into such features, along the lines of (2). Such an analysis yields insights like MMLU-Pro showing considerably more variability in its usage of plural nouns ( $\mu = 0.4$ ,  $\sigma = 0.5$ ) than BigBench Hard ( $\mu = 0.0$ ,  $\sigma = 0.0$ )<sup>10</sup>.

Differences like these may be expected, given that BigBench Hard integrates different problems into a given template format per subtask, and, in MMLU-Pro, each instance has a specific problem-instruction combination. `elfen` helps to quantitatively confirm such intuitions, which illustrates how **elfen facilitates linguistic comparisons of datasets along given axes of interest.**

<sup>9</sup>For a full table with Mantel correlation results, see Appendix E.1.

<sup>10</sup>For full statistics for morphological features, see Appendix E.2

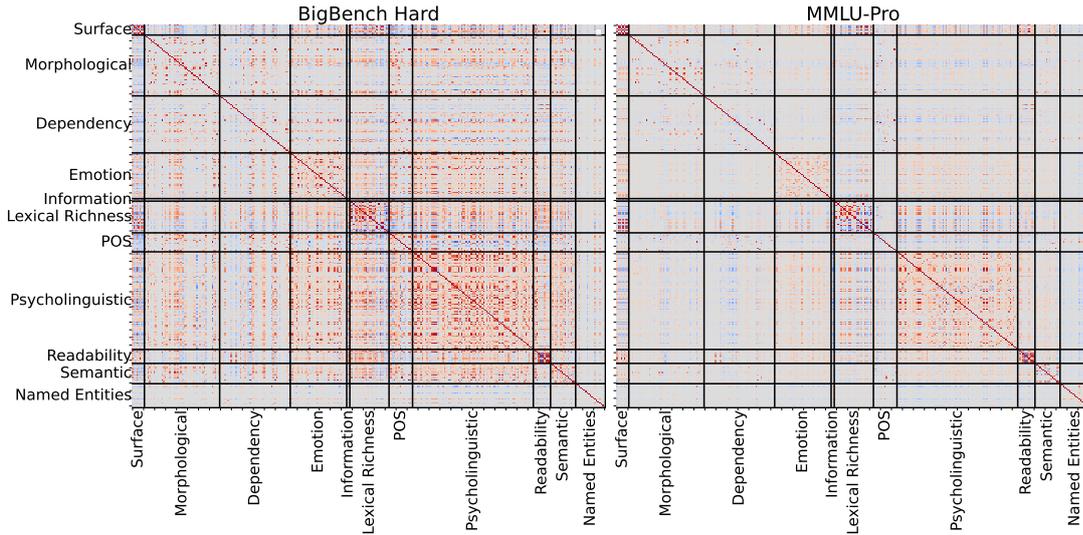


Figure 4: Heatmap of the correlation structures of BigBench Hard and MMLU-Pro for the eleven feature areas.

### 4.2.2 Targeted Sampling

Given the potential effects of structurally very different texts in a dataset, looking into respective samples for targeted comparisons or shot selection in few-shot scenarios can be beneficial. Following the latter example, suppose we want to ensure the examples for computer science problems in our shot selection include texts with relatively many tokens ( $> 50$ ) and a high relative frequency of nouns ( $> 0.25$ ). `elfen` provides the respective statistics of the subset overall, allowing for targeted sampling. For comparison, given that only 10.7% of instances in the computer science problems of MMLU-Pro fulfill these desiderata, the likelihood of having at least one such instance in a random sample is 0.203 for two shots and 0.365 for four shots. As this exemplifies, **elfen enables targeted sampling to use subsets of datasets with specific characteristics for downstream experiments.**

### 4.2.3 Outlier Detection

For linguistic bias-aware error analyses, it can be beneficial to understand whether a given model behaves differently for datapoints that are *outliers* in (a subset) of their structural characteristics. The features `elfen` provides can be used to run quick analyses to identify such outliers and inspect downstream effects when training or testing on them.

To showcase this, we construct linguistic fingerprints of the instances of MMLU-Pro by concatenating their respective features. We then use the local outlier factor to determine such outliers. Figure 5 shows a t-SNE projection of the fingerprints, showing that the identified outliers are either iso-

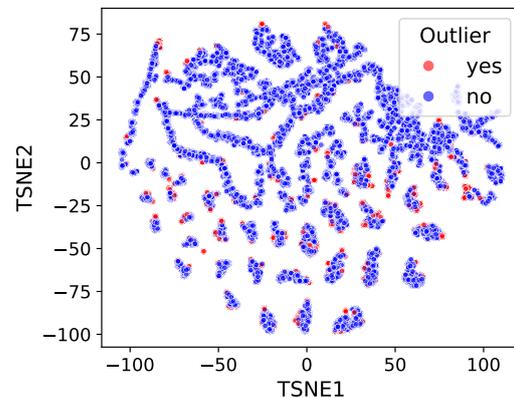


Figure 5: t-SNE projection of MMLU-Pro linguistic fingerprints showing outliers.

lated or on the edges of clusters.

A researcher interested in the sensitivity of different models to confounding characteristics could, for example, use such identified outliers to test whether there are systematic differences in how models perform on these instances specifically. In line with the previous subsection, if such differences are present, few-shot selection of such outliers could be tested as a way to address this.

Overall, this shows how **elfen helps to identify instances in datasets warranting special attention in experiments and downstream analyses.**

## 5 Conclusion

We presented `elfen`, a Python package to efficiently extract linguistic features for text datasets, building on existing NLP preprocessing libraries and established resources. `elfen` provides the most

extensive number of features of comparable tools, spanning eleven feature areas. We showcased the versatility of e1fen on prior work that already used it, and on three generalizable analysis use cases.

## Ethical Considerations and Limitations

While the features provided in e1fen are grounded in linguistic theory and draw on rigorously motivated and collected external resources, they should not be viewed as perfect or absolute properties, but rather as potentially noisy proxies for the underlying structures. This is mainly due to three limitations of the theories and resources e1fen builds on:

Firstly, not all features can be expected to transfer across languages. For instance, Mandarin often relies on syntactic order rather than inflectional morphology to encode grammatical relationships. Compared with many Indo-European languages, the same values for some features may lead to very different conclusions. **We thus encourage researchers to critically engage with what e1fen-derived features measure and reveal about linguistic realization when comparing them across languages.**

Secondly, we rely on existing tokenizers. While these tokenizers may be expected to work virtually perfectly in *well-behaved* text, they may not work as well in the presence of linguistic and orthographic variation such as dialects and sociolects (Wegmann et al., 2025). If this is not taken into account in the interpretation of results, this can lead to wrong inferences. This is clearly particularly problematic when the object of study includes the behavior of (groups of) humans. **Given the risk of flattening or misrepresenting groups or their language, we urge researchers using e1fen to carefully assess whether off-the-shelf tokenizers can handle the variation present in their data.**

Thirdly, external measurements such as psycholinguistic norms or affective dictionaries are subject to limitations that are passed down to the features in e1fen based on them. Besides limitations in the way ratings are collected (Mohammad, 2018a; Delatorre et al., 2019), the main concern is that most of them are collected from Western, well-educated, rich, and politically liberal<sup>11</sup>, *WEIRD* (Henrich et al., 2010), study participants and their language variants, causing a bias in both the selec-

tion of lexical items and the measurements (Siew et al., 2025). Finally, aggregated ratings may flatten individual differences (c.f. Knupleš et al., 2023; Paisios et al., 2023), resulting in a simplified picture of the complex reality of language perception. While these are problems outside of the scope of e1fen itself, and we continually update the included resources, **we urge caution for the inferences researchers make from psycholinguistic and affective features, as they may result from a WEIRD viewpoint on a limited number of lexical items, particularly for European languages.**

## Acknowledgements

We thank Gabriella Lapesa and Vigneshwaran Shankaran for their helpful comments on earlier versions of this manuscript. We thank the members of the Computational Social Science department at GESIS, and early adopters for feedback on usability, errors, and useful missing features in early versions of e1fen.

## References

- Jonathan Anderson. 1981. *Analysing the Readability of English and Non-English Texts in the Classroom with Lix*. *Seventh Australian Reading Association Conference*, pages 1–13.
- Abinew Ali Ayele, Nikolay Babakov, Janek Bevendorff, Xavier Bonet Casals, Berta Chulvi, Daryna Dementieva, Ashaf Elnagar, Dayne Freitag, Maik Fröbe, Damir Korenčić, Maximilian Mayerl, Daniil Moskovskiy, Animesh Mukherjee, Alexander Panchenko, Martin Potthast, Francisco Rangel, Naqee Rizwan, Paolo Rosso, Florian Schneider, and 10 others. 2024. *Overview of pan 2024: Multi-author writing style analysis, multilingual text detoxification, oppositional thinking analysis, and generative ai authorship verification condensed lab overview*. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction*, pages 231–259, Cham. Springer Nature Switzerland.
- Carl Hugo Björnsson. 1968. *Läsbarhet*. Pedagogiskt Utvecklingsarbete vid Stockholms Skolor. 6. Liber.
- Francis Bond and Ryan Foster. 2013. *Linking and extending an open multilingual Wordnet*. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1352–1362, Sofia, Bulgaria. Association for Computational Linguistics.
- Patrick Bonin, Aurélie Méot, and Aurélie Burgiska. 2018. *Concreteness ratings for 1,659 french words: Relationships with other psycholinguistic variables and word recognition times*. *Behavior Research Methods*, 50:2366–2387.

<sup>11</sup>*Liberal* here refers to the usage of the term in the US political landscape.

- Angana Borah, Daria Pylypenko, Cristina España-Bonet, and Josef van Genabith. 2023. [Measuring spurious correlation in classification: “clever hans” in translationese](#). In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 196–206, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.
- Ryan L Boyd, Ashwini Ashokkumar, Sarah Seraj, and James W Pennebaker. 2022. [The development and psychometric properties of LIWC-22](#). *Austin, TX: University of Texas at Austin*, 10:1–47.
- Marc Brysbaert, Paweł Mandera, Samantha F McCormick, and Emmanuel Keuleers. 2019. [Word prevalence norms for 62,000 English lemmas](#). *Behavior Research Methods*, 51:467–479.
- Marc Brysbaert, Amy Beth Warriner, and Victor Kuperman. 2014. [Concreteness ratings for 40 thousand generally known English word lemmas](#). *Behavior Research Methods*, 46:904–911.
- John Bissell Carroll. 1964. *Language and Thought*. Prentice-Hall.
- Hongyu Chen, Neele Falk, Michael Roth, and Agnieszka Falenska. 2025. [“feels feminine to me”: Understanding perceived gendered style through human annotations](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 31447–31468, Suzhou, China. Association for Computational Linguistics.
- Meri Coleman and Ta Lin Liau. 1975. [A computer readability formula designed for machine scoring](#). *Journal of Applied Psychology*, 60(2):283.
- Courtney Corley and Rada Mihalcea. 2005. [Measuring the semantic similarity of texts](#). In *Proceedings of the ACL Workshop on Empirical Modeling of Semantic Equivalence and Entailment*, pages 13–18, Ann Arbor, Michigan. Association for Computational Linguistics.
- Michael A. Covington and Joe D. McFall and. 2010. [Cutting the Gordian Knot: The Moving-Average Type–Token Ratio \(MATTR\)](#). *Journal of Quantitative Linguistics*, 17(2):94–100.
- Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. [Universal Dependencies](#). *Computational Linguistics*, 47(2):255–308.
- P. Delatorre, A. Salguero, C. León, and A. Tapscott. 2019. [The impact of context on affective norms: A case of study with suspense](#). *Frontiers in Psychology*, 10:1988.
- Veronica Diveica, Penny M. Pexman, and Richard J. Binney. 2023. [Quantifying social semantics: An inclusive definition of socialness and ratings for 8388 English words](#). *Behavior Research Methods*, 55(2):461–473.
- Esra Dönmez, Maximilian Maurer, Gabriella Lapesa, and Agnieszka Falenska. 2025. [AI argues differently: Distinct argumentative and linguistic patterns of LLMs in persuasive contexts](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 34583–34614, Suzhou, China. Association for Computational Linguistics.
- Daniel Dugast. 1978. [Sur quoi se fonde la notion d’etendue theoratique du vocabulaire? Le francais Modern](#), 46(1):25.
- Neele Falk and Gabriella Lapesa. 2025. [Mining the uncertainty patterns of humans and models in the annotation of moral foundations and human values](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 22898–22921, Vienna, Austria. Association for Computational Linguistics.
- Michael Wayne Goodman and Francis Bond. 2021. [Intrinsically interlingual: The wn python library for wordnets](#). In *Proceedings of the 11th Global Wordnet Conference*, pages 100–107, University of South Africa (UNISA). Global Wordnet Association.
- Arthur C Graesser, Danielle S McNamara, Max M Louwerse, and Zhiqiang Cai. 2004. [Coh-metrix: Analysis of text on cohesion and language](#). *Behavior research methods, instruments, & computers*, 36(2):193–202.
- Reto Gubelmann, Ioannis Katis, Christina Niklaus, and Siegfried Handschuh. 2024. [Capturing the varieties of natural language inference: A systematic survey of existing datasets and two novel benchmarks](#). *Journal of Logic, Language and Information*, 33(1):21–48.
- Pierre. Guiraud. 1954. *Les caractères statistiques du vocabulaire : essai de méthodologie*. Presses universitaires de France, Paris.
- Joseph Henrich, Steven J. Heine, and Ara Norenzayan. 2010. [The weirdest people in the world?](#) *Behavioral and Brain Sciences*, 33(2–3):61–83.
- Gustav Herdan. 1955. [A new derivation and interpretation of Yule’s ‘Characteristic’ K](#). *Zeitschrift für angewandte Mathematik und Physik ZAMP*, 6:332–339.
- Gustav Herdan. 1964. *Quantitative Linguistics*. Butterworths.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. [spaCy: Industrial-strength Natural Language Processing in Python](#).
- Philipp Kanske and Sonja A. Kotz. 2010. [The leipzig affective norms for german: A reliability study](#). *Behavior Research Methods*, 42(4):987–991.
- J. Peter Kincaid, Robert P. Jr. Fishburne, Richard L. Rogers, and Brad S. Chissom. 1975. [Derivation Of New Readability Formulas \(Automated Readability](#)

- Index, Fog Count And Flesch Reading Ease Formula) For Navy Enlisted Personnel. Technical report, Institute for Simulation and Training.
- Urban Knupleš, Diego Frassinelli, and Sabine Schulte im Walde. 2023. [Investigating the nature of disagreements on mid-scale ratings: A case study on the abstractness-concreteness continuum](#). In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 70–86, Singapore. Association for Computational Linguistics.
- Victor Kuperman, Hans Stadthagen-Gonzalez, and Marc Brysbaert. 2012. [Age-of-acquisition ratings for 30,000 English words](#). *Behavior Research Methods*, 44:978–990.
- Bruce W. Lee and Jason Lee. 2023. [LFTK: Handcrafted features in computational linguistics](#). In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 1–19, Toronto, Canada. Association for Computational Linguistics.
- Moira Linnarud. 1987. [Lexis in composition: a performance analysis of swedish learners’ written english](#). *Studies in Second Language Acquisition*, 9:254 – 256.
- Adian Liusie, Vatsal Raina, Vyas Raina, and Mark Gales. 2022. [Analyzing biases to spurious correlations in text classification tasks](#). 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 2: Short Papers)*, pages 78–84, Online only. Association for Computational Linguistics.
- Dermot Lynott, Louise Connell, Marc Brysbaert, James Brand, and James Carney. 2020. [The Lancaster Sensorimotor Norms: multidimensional measures of perceptual and action strength for 40,000 English words](#). *Behavior Research Methods*, 52:1271–1291.
- Nathan Mantel. 1967. [The detection of disease clustering and a generalized regression approach](#). *Cancer research*, 27(2\_Part\_1):209–220.
- Heinz-Dieter Mass. 1972. [Über den Zusammenhang zwischen Wortschatzumfang und Länge eines Textes](#). *Zeitschrift für Literaturwissenschaft und Linguistik*, 2(8):73.
- G. Harry Mc Laughlin. 1969. [Smog grading-a new readability formula](#). *Journal of reading*, 12(8):639–646.
- Philip M. McCarthy and Scott Jarvis. 2007. [vocr: A theoretical and empirical evaluation](#). *Language Testing*, 24(4):459–488.
- Philip M. McCarthy and Scott Jarvis. 2010. [MTLD, vocd-D, and HD-D: A validation study of sophisticated approaches to lexical diversity assessment](#). *Behavior Research Methods*, 42(2):381–392.
- Saif Mohammad. 2018a. [Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 English words](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 174–184, Melbourne, Australia. Association for Computational Linguistics.
- Saif Mohammad. 2018b. [Word affect intensities](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Saif Mohammad and Peter Turney. 2010. [Emotions evoked by common words and phrases: Using Mechanical Turk to create an emotion lexicon](#). In *Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*, pages 26–34, Los Angeles, CA. Association for Computational Linguistics.
- Saif M. Mohammad and Peter D. Turney. 2013. [Crowdsourcing a word–emotion association lexicon](#). *Computational Intelligence*, 29(3):436–465.
- Maria Montefinese, Elena Ambrosini, Beth Fairfield, and Nicola Mammarella. 2014. [The adaptation of the affective norms for english words \(anew\) in italian](#). *Behavior Research Methods*, 46:887–903.
- Maria Montefinese, David Vinson, Gabriella Vigliocco, and Ettore Ambrosini. 2019. [Italian age of acquisition norms for a large set of words \(itaoa\)](#). *Frontiers in Psychology*, Volume 10 - 2019.
- Pascale Moreira and Yuri Bizzoni. 2023. [Dimensions of quality: Contrasting stylistic vs. semantic features for modelling literary quality in 9,000 novels](#). In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 739–747, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.
- Dong Nguyen and Esther Ploeger. 2025. [We need to measure data diversity in nlp – better and broader](#). *Preprint*, arXiv:2505.20264.
- Dimitri Paisios, Nathalie Huet, and Elodie Labeye. 2023. [Addressing the elephant in the middle: Implications of the midscale disagreement problem through the lens of body-object interaction ratings](#). *Collabra: Psychology*, 9(1):84564.
- Angelina Parfenova, Alexander Denzler, and Jürgen Pfeiffer. 2025. [Emergent convergence in multi-agent LLM annotation](#). In *Proceedings of the 8th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 206–225, Suzhou, China. Association for Computational Linguistics.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. [Hypothesis only baselines in natural language inference](#). In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages

- 180–191, New Orleans, Louisiana. Association for Computational Linguistics.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108, Online. Association for Computational Linguistics.
- Alan Ramponi, Camilla Casula, and Stefano Menini. 2024. [Variationist: Exploring multifaceted variation and bias in written language data](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 346–354, Bangkok, Thailand. 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 and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Brian J. Richards and David D. Malvern. 1997. *Quantifying lexical diversity in the study of language development*. University of Reading, Faculty of Education and Community Studies.
- Yunita Sari, Mark Stevenson, and Andreas Vlachos. 2018. [Topic or style? exploring the most useful features for authorship attribution](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 343–353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Sabine Schröder, Teresa Gemballa, Steffie Ruppin, and Isabelle Wartenburger. 2011. [German norms for semantic typicality, age of acquisition, and concept familiarity](#). *Behavior Research Methods*, 44:380–394.
- Lucas Shen. 2022. [LexicalRichness: A small module to compute textual lexical richness](#).
- Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. [AI models collapse when trained on recursively generated data](#). *Nature*, 631(8022):755–759.
- Herbert S. Sichel. 1975. [On a distribution law for word frequencies](#). *Journal of the American Statistical Association*, 70(351a):542–547.
- Cynthia S. Q. Siew, Fera Chang, and Jin Jye Wong. 2025. [Investigating the effects of valence, arousal, concreteness, and humor on words unique to singapore english](#). *Journal of Cognition*.
- Daniel Simig, Tianlu Wang, Verna Dankers, Peter Henderson, Khuyagbaatar Batsuren, Dieuwke Hupkes, and Mona Diab. 2022. [Text characterization toolkit \(TCT\)](#). 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: System Demonstrations*, pages 72–87, Taipei, Taiwan. Association for Computational Linguistics.
- Edward H. Simpson. 1949. [Measurement of Diversity](#). *Nature*, 163.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. 2022. [Challenging big-bench tasks and whether chain-of-thought can solve them](#). *arXiv preprint arXiv:2210.09261*.
- Mildred C. Templin. 1957. *"Certain Language Skills in Children: Their Development and Interrelationships"*, ned - new edition edition, volume 26. University of Minnesota Press.
- Alessandra Vergallito, Marco Alessandro Petilli, and Marco Marelli. 2020. [Perceptual modality norms for 1,121 italian words: A comparison with concreteness and imageability scores and an analysis of their impact on word processing tasks](#). *Behavior Research Methods*, 52:1599–1614.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. 2024. [Mmlu-pro: A more robust and challenging multi-task language understanding benchmark](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 95266–95290. Curran Associates, Inc.
- Anna Wegmann, Dong Nguyen, and David Jurgens. 2025. [Tokenization is sensitive to language variation](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10958–10983, Vienna, Austria. Association for Computational Linguistics.
- Bodo Winter, Gary Lupyan, Lynn K Perry, Mark Dingemanse, and Marcus Perlman. 2024. [Iconicity ratings for 14,000+ English words](#). *Behavior Research Methods*, 56(3):1640–1655.
- Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia Sam Chao, and Derek Fai Wong. 2025. [A survey on LLM-generated text detection: Necessity, methods, and future directions](#). *Computational Linguistics*, 51(1):275–338.
- George U. Yule. 1944. *The statistical study of literary vocabulary*. Cambridge University Press.

Peter Zeng, Pegah Alipoormolabashi, Jihu Mun, Gourab Dey, Nikita Soni, Niranjan Balasubramanian, Owen Rambow, and H. Schwartz. 2025. [Residualized similarity for faithfully explainable authorship verification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 15824–15837, Suzhou, China. Association for Computational Linguistics.

Yian Zhang, Alex Warstadt, Xiaocheng Li, and Samuel R. Bowman. 2021. [When do you need billions of words of pretraining 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 1112–1125, Online. Association for Computational Linguistics.

## A Full Description of Features

In the following, we provide a more detailed description of the available features per feature area, including references to relevant literature introducing or describing them.

**Surface-Level Features** provide structural characteristics of the texts. The package provides extraction of the sequence length (characters; both with and without whitespaces), number of tokens, sentences, types, lemmas, long words (over six characters), and token frequencies on an item level. Based on those, the number of tokens per sentence, characters per sentence, and average word length can be extracted.

**Readability Features** were proposed to measure the complexity of texts. Following the readability python package<sup>12</sup>, the package extracts the Gunning fog index, ARI, Flesch reading ease, and Flesch-Kincaid grade level (Kincaid et al., 1975), the Cole-Liau index (Coleman and Liau, 1975), SMOG (Mc Laughlin, 1969), LIX (Björnsson, 1968), and RIX (Anderson, 1981). Additionally, the package provides extraction of the basic features necessary for calculating these readability scores: The number of syllables in a text, words with only one syllable, and words with more than two syllables.

**Psycholinguistic Norm Features** measure words’ cognitive, social, and sensorimotor grounding. We use concreteness norms<sup>13</sup> (Brysbaert et al., 2014), i.e. how concrete or abstract is a

given word, word prevalence norms (Brysbaert et al., 2019), i.e. how well-known/used is a word, Age-of-Acquisition norms (Kuperman et al., 2012), i.e. at what age do children learn a given word, Socialness norms (Diveica et al., 2023), i.e. how socially relevant is a words’ meaning, Iconicity norms (Winter et al., 2024), i.e. to which degree the sound of a word reflects its’ meaning, and Sensorimotor norms (Lynott et al., 2020), i.e. how connected a words’ meaning is to perceptual modalities (e.g. visual) and action effectors (e.g. arm/hand). Per item and for each norm, the package implements the extraction of the average rating of all tokens from the item in the norm lexicon, their average standard deviation in the ratings, the number of tokens with a high rating (upper third of the Likert scale), a low rating (lower third of the scale), and the number of tokens with a particularly high standard deviation (such that the ratings span over multiple thirds of the scale).

While the norms are collected for individual words without context, these features are included to measure potential effects of the presence of words with a particular grounding or ambiguity thereof.

**Part-of-Speech Features.** Per item, the package provides extraction of the number of tokens per universal dependencies POS tag (de Marneffe et al., 2021), the number of lexical tokens (nouns, verbs, adjectives, and adverbs), and the POS variability (number of different POS tags relative to the number of tokens).

**Lexical Richness Measures** provide information about how lexically diverse a given text is. Intuitively, the more lexically rich a text is, the more different words a text contains. Following the lexicalrichness python package (Shen, 2022), per item, elfen allows for the extraction of the type-token ratio (TTR) (Templin, 1957), root TTR (Guiraud, 1954), corrected TTR (Carroll, 1964), Herdan’s C (Herdan, 1964), Summer’s TTR, Dugast’s Uber index (Dugast, 1978), Maas’ TTR (Mass, 1972), Yule’s  $K$  (Yule, 1944), Herdan’s  $V_m$  (Herdan, 1955), Simpson’s  $D$  (Simpson, 1949), mean segmental TTR (Richards and Malvern, 1997), moving average TTR (Covington and and, 2010), measure of textual lexical diversity (MLTD, McCarthy and Jarvis, 2010), and the hypergeometric distribution diversity (HD-D, McCarthy and Jarvis, 2007, 2010). Additionally, the lo-

<sup>12</sup><https://github.com/andreasvc/readability>

<sup>13</sup>All of the cited norms here are in English. Find the full list of currently supported languages per psycholinguistic dimension, including references in Appendix B. We regularly add more.

cal and global numbers of hapax (dis)legomena (i.e. the number of words per item that occur only once/twice per item/globally in the dataset), Sichel’s S (Sichel, 1975), and the lexical density, i.e., the percentage of lexical tokens (Linnarud, 1987) are extractable.

**Morphological Features.** `elfen` allows for the extraction of the number of tokens with a given morphological feature for all available universal dependencies morpho-syntactic features (de Marneffe et al., 2021).

**Information-Theoretic Features.** As a measure of redundancy or formulaicity, following Moreira and Bizzoni (2023), `elfen` implements the compressibility of the text per item. The compressibility is defined as the bit length of the compressed text divided by the bit length of the uncompressed text. `elfen` implements the average token Shannon entropy per item to measure predictability.

**Dependency Features** provide information about the morphosyntactic realizations of predicate-argument structures. `elfen` implements the number of dependency relation types (according to Universal Dependencies, de Marneffe et al., 2021), the number of noun chunks in the text, the tree width, i.e. the maximum number of nodes in the subtree of a token, the tree depth, i.e the maximum distance of a token from the root of the dependency tree, the tree branching factor, i.e. the average number of children of a token, and the ramification factor, i.e. the mean number of children per level in the dependency tree.

**Semantic Features.** To measure the impact of token-level ambiguity/polysemy on the text, we extract Open Multilingual Wordnet synsets (Bond and Foster, 2013) for all nouns, adjectives, and verbs using the `wn` python library (Goodman and Bond, 2021). Given these synsets, per item, we extract the average size of the synsets, the number of tokens with a large synset (more than four senses), and the number of tokens with a small synset (less than three senses) for nouns, adjectives, and verbs, respectively, and overall.

We extract the number of hedges<sup>14</sup> (i.e., words expressing speaker uncertainty; e.g., *might*, *presumably*, or *maybe*) and the hedge-token ratio per item as a measure of the presence of uncertainty expressions in the text.

<sup>14</sup><https://github.com/words/hedges>

**Named Entity Features.** Per item, we extract the number of named entities overall and per entity type (e.g. names, locations, organizations, etc.).

**Emotion and Sentiment Features.** To measure the effects of the occurrence of words commonly associated with/evoking a given emotion or sentiment, we use the NRC-VAD lexicon (Mohammad, 2018a) for valence, arousal, and dominance, the NRC emotion intensity lexicon (Mohammad, 2018b) for the emotion intensity per basic emotion (anger, anticipation, disgust, fear, joy, sadness, surprise, trust), and the NRC word-emotion association lexicon (Mohammad and Turney, 2010, 2013) for sentiment. Per item and emotion dimension, we extract the average rating of all tokens from the item in the emotion lexicon, the number of tokens with a high rating, and the number of tokens with a low rating. For sentiment, per item, we extract the number of positive and negative sentiments, and the difference between them, normalized by the total number of tokens in the item.

## B Norms/Lexicons per Language

Table 4 gives an overview of the availability of currently (v1.2.4) supported languages per psycholinguistic variable.

Variable	Language	Reference
Concreteness	DE	Kanske and Kotz (2010)
	EN	Brybaert et al. (2014)
	FR	Bonin et al. (2018)
	IT	Montefinese et al. (2014)
Age of Acquisition	DE	Schröder et al. (2011)
	EN	Kuperman et al. (2012)
	IT	Montefinese et al. (2019)
Sensorimotor	EN	Lynott et al. (2020)
	IT	Vergallito et al. (2020)

Table 4: Psycholinguistic norms included in `elfen` v1.2.4.

## C Additional Code Examples

Figure 6 gives an additional code example for normalization and rescaling. Figure 7 gives a code example of utilities included in `elfen`.

## D Extraction of MMLU-Pro and BigBench Hard

We extract both MMLU-Pro and BigBench Hard on an Apple MacBook Pro with 24GB RAM and an Apple M4 chip. Table 5 shows the preprocessing

```
# rescale "n_tokens" to a range between 0 and 1
extractor.rescale("n_tokens",
                 minimum = 0,
                 maximum = 1)
# token-normalize "n_entities"
extractor.token_normalize("n_entities")
# normalize all features to a mean 0 std 1
extractor.normalize("all")
```

Figure 6: Code examples of token and zero-mean normalization, and rescaling.

```
# list all external resources available in elfen
elfen.list_external_resources()
# get a bibtex string for all resources
print(elfen.get_bibtex())
```

Figure 7: Code examples for utilities.

and extraction times with the number of instances per dataset.

Dataset	Size	Preprocessing	Extraction
MMLU-Pro	12,032	86.07s	812.86s
BigBench Hard	6,511	67.25s	453.96s

Table 5: Preprocessing and extraction times for MMLU-Pro and BigBench Hard.

## E Full Results Use Cases

This section presents the full results for the use cases in section 4.

### E.1 Full Mantel Results

Table 6 provides the full Mantel test results of the analysis use case in Section 4.2.1.

### E.2 Full Morphological Feature Comparison

Table 7 provides the full descriptive statistics for the analysis of morphological features in the use case presented in Section 4.2.1.

Feature Group	Mantel	p-value
Surface	0.925	0.005
Morphological	0.430	0.000
Dependency	0.304	0.000
Emotion	0.444	0.000
Lexical Richness	0.735	0.000
POS	-0.026	0.785
Psycholinguistic	0.706	0.000
Readability	0.958	0.000
Semantic	0.555	0.000
Named Entities	0.519	0.004
All	0.481	0.000

Table 6: Mantel correlations including p-values per feature area between the correlation matrices of BigBench Hard and MMLU-Pro.

Feature	BigBench Hard				MMLU-Pro			
	$\mu$	$\sigma$	min	max	$\mu$	$\sigma$	min	max
n_NOUN_Number_Plur	0.00	0.00	0.00	0.01	0.05	0.05	0.00	0.67
n_VERB_VerbForm_Inf	0.00	0.00	0.00	0.01	0.02	0.03	0.00	0.33
n_PRON_Number_Sing	0.00	0.00	0.00	0.01	0.01	0.02	0.00	0.22
n_VERB_VerbForm_Part	0.00	0.00	0.00	0.02	0.03	0.04	0.00	0.50
n_PROPN_Case_Nom	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.03
n_PROPN_Number_Sing	0.00	0.00	0.00	0.03	0.05	0.06	0.00	0.62
n_VERB_Mood_Ind	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.25
n_PRON_Case_Acc	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.14
n_PRON_PronType_Prs	0.00	0.00	0.00	0.01	0.01	0.03	0.00	0.23
n_PUNCT_PunctType_Dash	0.00	0.00	0.00	0.01	0.01	0.02	0.00	0.22
n_NOUN_Number_Sing	0.00	0.00	0.00	0.03	0.19	0.08	0.00	1.00
n_PRON_Case_Nom	0.00	0.00	0.00	0.01	0.01	0.02	0.00	0.18
n_PUNCT_PunctType_Peri	0.00	0.00	0.00	0.01	0.05	0.03	0.00	0.33
n_DET_Number_Sing	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.10
n_PRON_PronType_Dem	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.17
n_PRON_PronType_Ind	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.12
n_VERB_Aspect_Prog	0.00	0.00	0.00	0.01	0.01	0.03	0.00	0.50
n_ADJ_Degree_Pos	0.00	0.01	0.00	0.05	0.06	0.05	0.00	0.50
n_PRON_Reflex_Yes	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.11
n_PRON_Gender_Masc	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.15
n_CCONJ_ConjType_Cmp	0.00	0.00	0.00	0.04	0.02	0.03	0.00	0.36
n_VERB_Person_3	0.00	0.00	0.00	0.01	0.01	0.03	0.00	0.33
n_PRON_Person_1	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.21
n_PRON_PronType_Art	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.33
n_VERB_Tense_Pres	0.00	0.00	0.00	0.01	0.03	0.04	0.00	0.50
n_DET_Definite_Def	0.00	0.00	0.00	0.01	0.07	0.05	0.00	0.33
n_PUNCT_PunctType_Quot	0.00	0.00	0.00	0.01	0.00	0.02	0.00	0.36
n_PUNCT_PunctType_Comm	0.00	0.00	0.00	0.01	0.02	0.03	0.00	0.37
n_PRON_Gender_Neut	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.17
n_VERB_Aspect_Perf	0.00	0.00	0.00	0.01	0.02	0.03	0.00	0.50
n_PROPN_Number_Plur	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.17
n_PRON_Gender_Fem	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.18
n_VERB_Number_Sing	0.00	0.00	0.00	0.01	0.01	0.03	0.00	0.33
n_VERB_Tense_Past	0.00	0.00	0.00	0.01	0.03	0.03	0.00	0.50
n_PUNCT_PunctSide_Fin	0.00	0.00	0.00	0.04	0.01	0.03	0.00	0.33
n_PUNCT_PunctType_Brck	0.00	0.01	0.00	0.07	0.02	0.04	0.00	0.46
n_ADJ_Degree_Sup	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.25
n_PRON_Poss_Yes	0.00	0.00	0.00	0.00	0.01	0.01	0.00	0.17
n_PRON_Person_3	0.00	0.00	0.00	0.01	0.01	0.02	0.00	0.18
n_NUM_NumType_Card	0.00	0.00	0.00	0.02	0.05	0.06	0.00	0.68
n_PRON_PronType_Rel	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.20
n_PRON_Person_2	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.22
n_DET_Definite_Ind	0.00	0.00	0.00	0.01	0.03	0.04	0.00	0.50
n_ADJ_Degree_Cmp	0.00	0.00	0.00	0.01	0.00	0.01	0.00	0.25
n_DET_Number_Plur	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.12
n_VERB_VerbForm_Fin	0.00	0.00	0.00	0.01	0.03	0.03	0.00	0.40
n_PRON_Number_Plur	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.21
n_PUNCT_PunctSide_Ini	0.00	0.00	0.00	0.04	0.01	0.02	0.00	0.33

Table 7: Full overview of statistics (mean  $\mu$ , standard deviation  $\sigma$ , min and max) comparing BigBench Hard and MMLU-Pro on morphological features.

# DELTA: A Toolkit for Measuring Linguistic Diversity in Dependency-Parsed Corpora

**Louis Estève**

Université Paris-Saclay, CNRS, LISN  
91400, Orsay, France  
louis.esteve@lisn.fr

**Kaja Dobrovoljc**

University of Ljubljana, Slovenia  
Jozef Stefan Institute, Slovenia  
kaja.dobrovoljc@ff.uni-lj.si

## Abstract

Despite growing interest in measuring linguistic diversity on the one hand and the increasing availability of cross-linguistically comparable parsed corpora on the other, tools for systematically measuring the diversity of specific linguistic phenomena on such data remain limited. To address this gap, we present DELTA, an open-source framework that integrates dependency tree querying with diversity computation, enabling systematic measurement across multiple linguistic levels (e.g., lexis, morphology, syntax) and multiple diversity dimensions (variety, balance, disparity). The pipeline processes CoNLL-U formatted corpora through configurable workflows, treating the format as a general-purpose tabular structure independent of specific annotation conventions. We validate DELTA on Parallel Universal Dependencies multilingual dataset, demonstrating its capacity for corpus profiling and cross-corpus diversity comparison.

## 1 Introduction

Natural Language Processing (NLP) has seen increasing interest in the concept of diversity in recent years. The year-wise share of papers in ACL Anthology containing “diversity” or “diverse” in their title or abstract has risen from less than 1% in the 2000s to over 10% in 2024 (Estève et al., 2025). Archetypical examples include research on generative models with concerns for diverse output, or dataset creation efforts prioritizing diverse content.

This growing interest reflects the variety of motivations, metrics, and target phenomena associated with diversity. Estève et al. (2025) identify two axes describing the motivations behind diversity: *goal versus means*, and *practical versus ethical*. For instance, ethical motivations include improved deontology (Song et al., 2024) and inclusiveness (Joshi et al., 2020), while practical motivations include meeting user expectations (Kumar et al., 2019) and improving model performance (Liu and

Zeldes, 2023). The survey identifies 150 different equations for measuring diversity, and the target phenomena span multiple linguistic levels: lexis (Kosmajac and Keselj, 2019), morphology (Samir and Silfverberg, 2023), syntax (Guo et al., 2024), and semantics (Jolly et al., 2021).

In recent years, there has been growing interest in measuring diversity specifically on dependency-parsed corpora. The widespread adoption of Universal Dependencies (de Marneffe et al., 2021) – now covering hundreds of languages with consistent annotation (Zeman et al., 2025) – has made parsed data readily available for cross-linguistic analysis. Parsed corpora enable investigation of diversity across multiple linguistic levels—from lexis and morphology to syntactic and semantic patterns – opening possibilities for studying linguistic universals (Gerdes et al., 2021), typological differences (Levshina, 2019), and the impact of diversity on NLP systems (Savary et al., 2024).

However, while tools exist for querying parsed corpora and for computing diversity metrics, these capabilities have typically remained separate. Researchers must manually combine pattern extraction with statistical analysis, export intermediate results, and write custom code to bridge the two stages. This fragmented workflow hinders reproducibility, limits cross-study comparability, and presents barriers for researchers without programming expertise. What remains absent is an integrated framework enabling multi-level and multi-dimensional diversity measurement directly from dependency-parsed corpora, supporting comparative analysis across languages and datasets.

To bridge this gap, we present DELTA, a unified pipeline for measuring linguistic diversity in dependency-parsed corpora. DELTA integrates tree extraction and diversity computation, enabling multi-level and multi-dimensional diversity analysis across languages. In doing so, our main contributions are:

1. **An integrated pipeline** from annotated input to diversity scores, combining dependency graph querying with diversity metric computation on CoNLL-U formatted data, eliminating the need for manual data transformation between tools.
2. **A flexible framework** that enables multi-level and multi-dimensional measurement across linguistic levels (lexis, morphology, syntax) and diversity dimensions (variety, balance, disparity), with support for customizable queries and metrics, and adaptable to alternative annotation schemes beyond Universal Dependencies.
3. **Scalable infrastructure** with SLURM-based parallelization for large-scale analyses, along with pre-defined configurations for common diversity measurement tasks and plotting support.

We present the pipeline in the remainder of this paper and demonstrate it by measuring lexical, morphological and syntactic diversity across Parallel Universal Dependencies (PUD) treebanks (Zeman et al., 2017).

## 2 Related work

Numerous tools support structured exploration of dependency-parsed corpora, including online services such as Grew-match (Guillaume, 2021), PML-TQ (Štěpánek and Pajas, 2010), and INESS (Rosén et al., 2012), lightweight libraries such as `pyconll`<sup>1</sup> and `conllu`<sup>2</sup> for programmatic access, and STARK (Krsnik and Dobrovoljc, 2025) for quantitative subtree extraction. Separately, various libraries implement diversity indices from ecology and information theory, including the diverse R package (Guevara et al., 2016), `scikit-bio` (Rideout et al., 2025), and `DiversUtils`<sup>3</sup>. However, these querying and computation capabilities remain disconnected, requiring manual export and custom code to bridge the two stages.

Tools such as `Distals` (Goot et al., 2025), `LangDive` (Samardzic et al., 2024), and `TypDiv` (Ploeger et al., 2024) quantify diversity at the level of language samples—using typological databases and, in some cases, text-derived features – but focus

<sup>1</sup><https://pyconll.github.io/>

<sup>2</sup><https://pypi.org/project/conllu/>

<sup>3</sup><https://github.com/estevélouis/WG4>

on comparing languages or multilingual datasets rather than profiling specific phenomena within corpora. For profiling parsed corpora specifically, tools such as `ComparaTree` (Terčon and Dobrovoljc, 2025), `Profiling-UD` (Brunato et al., 2020), and `Typometrics` (Gerdes et al., 2021) support cross-linguistic comparison, but are limited to fixed feature sets and predefined comparison scenarios.

DELTA bridges these approaches within a single unified framework, enabling flexible, multi-level, and multi-dimensional diversity measurement of specific linguistic phenomena directly over dependency-parsed corpora.

## 3 System Architecture

DELTA takes any number of annotated corpora in CoNLL-U format as input and produces diversity measurements for each corpus as output. Built on two preexisting open-source tools – STARK<sup>4</sup> for pattern extraction and `DiversUtils` for diversity computation – the system provides a unified framework for flexible diversity analysis. Figure 1 illustrates the pipeline: users provide two configuration files specifying (1) which linguistic patterns to extract and (2) which diversity metrics to compute, and the system then extracts matching instances from each corpus and calculates their diversity.

In essence, DELTA’s diversity measurement relies on the distinction between *categories* (types of linguistic patterns) and *elements* (individual occurrences of those patterns). For instance, if “noun phrase” is defined as a category of interest, then each individual noun phrase in the corpus constitutes an element of that category. This element/category dichotomy comes from ecology, where it is often termed the type/item dichotomy (Ramaciotti Morales et al., 2021; Solé et al., 2010): categories (types) often correspond to species, in which case elements (items) correspond to individual organisms.

The following subsections describe pattern extraction (§3.1), diversity computation (§3.2), output (§3.3), availability (§3.4) and scalability (§3.5).

### 3.1 Category extraction (STARK)

DELTA represents categories as dependency subtrees extracted from CoNLL-U formatted corpora. This tree-based representation means DELTA can compute diversity of any linguistic phenomenon expressible in tree-like form – from single-node

<sup>4</sup><https://github.com/clarinsi/STARK>

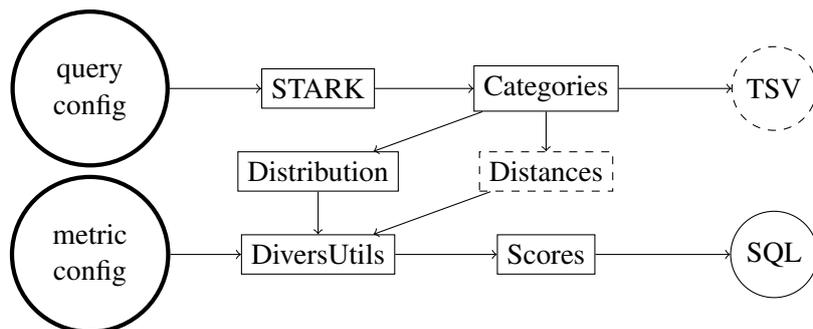


Figure 1: High-level representation of the main mechanisms in DELTA. Objects with dashed contour may not be generated depending on configuration.

subtrees (e.g., word forms, lemmas, POS tags) capturing lexical or morphological diversity, to multi-node subtrees (e.g., phrases, clauses, sentences) capturing (morpho)syntactic diversity.

**STARK for pattern extraction.** Category extraction is performed using STARK, an open-source toolkit for extracting dependency subtrees from parsed corpora. STARK extracts subtrees matching user-defined configurations, which can range from all attested subtrees to specific structural types. Each extracted subtree is counted and output with frequency information, providing the information for diversity computation.

**Flexible pattern specification.** The subtrees to be extracted can be defined flexibly along multiple dimensions. Users control tree representation by specifying which information appear on nodes (word forms, lemmas, part-of-speech tags, or combinations thereof), whether dependency labels are included, and whether linear word order is considered. Users also control tree filtering by specifying size constraints (e.g., only trees with a specific number of nodes), head constraints (e.g., only noun-headed structures), allowed or ignored dependency relations, or exact structural patterns via custom queries (e.g., only adjective-noun structures). These flexible and combinable parameters give users precise control over extraction granularity – from fully lexicalized constructions to abstract structural templates – making category extraction adaptable to diverse research goals. An overview of STARK’s functionality is given by Krsnik and Dobrovoljc (2025), with detailed documentation also available online.<sup>5</sup>

**Predefined configurations.** To facilitate common use cases, DELTA also provides some pre-

defined extraction configurations targeting specific linguistic phenomena. For single-node extraction, configurations include `forms.ini` (all word forms), `lemmas.ini` (all lemmas), `parts-of-speech.ini` (all POS tags), and `morphology.ini` (all POS and feature combinations). For multi-node structures, two general configurations extract all subtrees as proxies for all syntactic structures attested in a corpus (see Dobrovoljc (2025) for details): `syntactic-structures.ini` for extracting all dependency-labeled subtrees and `morphosyntactic-structures.ini` for extracting all dependency-labeled subtrees with POS tags as nodes.

As examples of using DELTA to measure diversity of very specific phenomena, we also include two specialized configurations targeting commonly analyzed syntactic structures: `svo.ini` for extracting delexicalized subject-verb-object patterns featuring *nsubj* and *obj* relations (Levshina, 2019) and `mwe.ini` for extracting lexicalized multi-word expressions featuring *fixed*, *flat*, and *compound* constructions (Savary et al., 2023). These configurations facilitate standard diversity measurements without requiring detailed parameter specification, but users can also define custom configurations, if needed.

**Format flexibility.** Crucially, STARK operates on the 10-column tabular structure without assumptions about its content, i.e., the type of values expected in the columns. While we use standard CoNLL-U column semantics throughout this paper, any categorical labels can populate the tag columns (UPOS, XPOS, FEATS), any directed relations can populate the dependency columns (HEAD, DEPREL), including semantic dependencies, and any unit can serve as a token (FORM, LEMMA), in-

<sup>5</sup><https://github.com/clarinsi/STARK/blob/master/settings.md>

cluding multi-word sequences, e.g. for sequence-based diversity measurement.

### 3.2 Diversity computation (DiversUtils)

Diversity computation in DELTA is performed using DiversUtils, a C/Python library that implements diversity metrics from ecological and information-theoretic frameworks. DiversUtils takes the category frequencies extracted by STARK as input. Users specify which diversity metrics to compute via the DiversUtils configuration file. The library currently implements 32 diversity metrics, which can be understood along three complementary dimensions: *variety* (the number of categories), *balance* (the evenness of their distribution), and *disparity* (the degree of difference between categories). This framework, adapted from ecology (Ramaciotti Morales et al., 2021; Lion-Bouton et al., 2022), provides a conceptual structure for understanding what different metrics capture.

**Variety** measures how many distinct categories are present in the corpus. Simple variety metrics include richness which is just the number of categories, and “species count” which is richness minus one, such that in the minimum case where only one category is present, the diversity scores zero (Patil and Taillie, 1982). Higher variety indicates a greater number of distinct linguistic patterns in the data. Note that variety is sensitive to corpus size: larger corpora can account for a wider set of phenomena, especially rare ones (see the correlation between variety and treebank sentence count in Figures 3 and 4).

**Balance** measures the evenness of the frequency distribution – whether categories are equally represented or dominated by a few high-frequency categories. Balance is captured by metrics such as Shannon evenness (Ramaciotti Morales et al., 2021) for pure balance or entropies from the set of generalised entropies (Rényi, 1961; Patil and Taillie, 1982) for variety-balance hybrids. Higher balance indicates more uniform usage across categories.

**Disparity** measures the degree of structural difference between categories, capturing the dissimilarity between them. Unlike variety and balance, disparity requires a distance function. DELTA uses Zhang-Shasha tree edit distance by default (Zhang and Shasha, 1989),<sup>6</sup> which captures linguistically meaningful tree differences by considering both

nodes and edges. For single-word trees, Word2Vec cosine distance can be specified as a semantic alternative (using `--w2v_path`).

**Multi-dimensional metrics.** Many diversity metrics encompass multiple dimensions simultaneously (Chao et al., 2014; Stirling, 2007). For example, Shannon-Wiener entropy in its original definition (Wiener, 1939; Shannon, 1948; Shannon and Weaver, 1949) is a hybrid of variety and balance, increasing when either more categories are present or when frequencies are more evenly distributed. Generalizations of entropy with varying parameters exhibit different weightings of variety and balance (Rényi, 1961; Patil and Taillie, 1982; Hill, 1973). The concept of entropy has been further generalized in ecology by incorporating distances between categories (Chao et al., 2014; Leinster and Cobbold, 2012; Scheiner, 2012), thus accounting for disparity in addition to variety and balance.

**Methodological considerations.** In the examples in this paper we make the choice of using easily interpretable metrics (richness for pure variety, and Shannon evenness for pure balance). For authors wishing to build a single unified ranking among datasets, a variety-balance hybrid is desirable. Based on the long history of diversity in ecology and biology, it is notably relevant to use Hill (1973) numbers rather than entropies, as Hill assesses that “The diversity numbers  $N_\alpha$ , have therefore a natural intuitive interpretation, albeit rather a vague one. The corresponding generalized entropies  $H_\alpha$ , being logarithmic, are harder to visualize.”. Hill numbers are interpreted as the “effective number of species” if all species were equally probable, to give the same entropy. To add disparity, consider generalised Hill numbers such as that of Chao et al. (2014) as a start. Generalised Hill numbers are interpreted as the “effective number of equally common, equally distinct species or lineages”. Conversely, the approach by Stirling (2007) to adding disparity has been criticized (Leydesdorff et al., 2019). Beyond the choice of the measure, when using datasets of non-commensurate sizes, consider averaging diversity scores over numerous samples of same sizes, so as to prevent biases due to dataset size.

### 3.3 Output and visualisation

DELTA produces two main artifacts from each analysis. First, STARK-produced category frequency lists are (optionally) stored in tab-separated for-

<sup>6</sup>We use the Python zss package.

mat (TSV), listing all extracted linguistic patterns with their frequencies (see example in Table 1). These files can be inspected manually to understand which categories were found, or reused independently in external analyses. Second, DiversUtils-produced diversity scores are stored in a SQLite database, providing a structured format optimized for querying and cross-corpus comparison. Each database record includes the corpus identifier, linguistic level analyzed, diversity metric applied, and the computed score.

The repository also includes preconfigured analysis scripts that operate directly on these databases, for example to plot variety against balance for a given collection of corpora (e.g., Figure 2), optionally also indicating the corpus size (e.g., Figures 3 and 4), enabling straightforward identification of diversity patterns across treebanks, languages, or linguistic levels. Users can also write custom queries against the SQLite databases to generate application-specific analyses or export results in alternative formats.

### 3.4 Availability and execution

DELTA is freely available as open-source software under the joint BSD-2 and CeCILL-B license.<sup>7</sup> The repository includes complete documentation, installation instructions, predefined configurations for common use cases, and example analysis scripts.

The system can be installed as a command-line program for Python, via instructions in the README.md. DELTA is executed via a command-line interface that takes three inputs: (1) a STARK configuration file specifying linguistic patterns to extract, (2) a DiversUtils configuration file specifying diversity metrics to compute, and (3) a list of input corpora in CoNLL-U format. A single command processes all specified corpora and produces the outputs described in Section 3.3.

The accompanying demonstration video,<sup>8</sup> illustrates the DELTA workflow by showing how the system computes and visualizes syntactic diversity over the full UD dataset. It reproduces the results shown in the Appendix (Figure 4), demonstrating the end-to-end pipeline from configuration to visualization.

<sup>7</sup>Hosted at <https://gitlab.lisn.upsaclay.fr/esteveldelta/>. CeCILL-B is the French equivalent of BSD-2, formed by French national scientific institutions.

<sup>8</sup><https://gitlab.lisn.upsaclay.fr/esteveldelta/-/blob/main/video/DELTA-system-demo-video.mp4>

### 3.5 Scalability

For large-scale analyses, DELTA supports SLURM-based parallelization through array job submission, enabling efficient processing of multiple configurations across many treebanks simultaneously. In practice, the scale at which DELTA can work depends on the computational expense of both the tree querying and the diversity computation.

For tree querying, some queries may return as little as a constant number of elements per sentence  $O(1)$  or as high as an exponential number of elements per sentence  $O(x^s)$  where  $s$  is sentence size. Likewise, for diversity computation, variety and balance measures often take at most linear time  $O(n)$ , but disparity takes quadratic time  $O(n^2)$ , where  $n$  is the number of extracted categories.<sup>9</sup> Disparity computation is also sensitive to tree complexity: computing distance matrices for large category sets or complex trees can be intensive, so for the Zhang-Shasha tree edit distance, a configurable timeout (default 0.25s) limits computation time, producing exact or approximate results.

Empirically, inexpensive queries with linear diversity metrics can process billion-token datasets, with substantial time on reading/writing and parsing. In contrast, queries using disparity functions with tree edit distance become computationally intensive: even 1,000 categories can require multiple hours for matrix computation.

To provide a concrete benchmark: processing the entire UD v2.16 release for both lexical diversity (lemmas.ini) and syntactic diversity (syntactic-structures.ini) with linear metrics (linear.ini), to produce results shown in Figures 3 and 4, took 45 minutes total.

## 4 Evaluation

We validate DELTA’s core capabilities through four experiments on Parallel Universal Dependencies (PUD), a collection of parallel treebanks consisting of 1,000 aligned dependency-parsed sentences across 24 languages (Zeman et al., 2017). By controlling for content and corpus size, PUD provides an ideal testbed for systematic cross-linguistic comparison.

<sup>9</sup>See the linear.ini and quadratic.ini configuration files for respectively at-most-linear, and at-most-quadratic diversity computations.

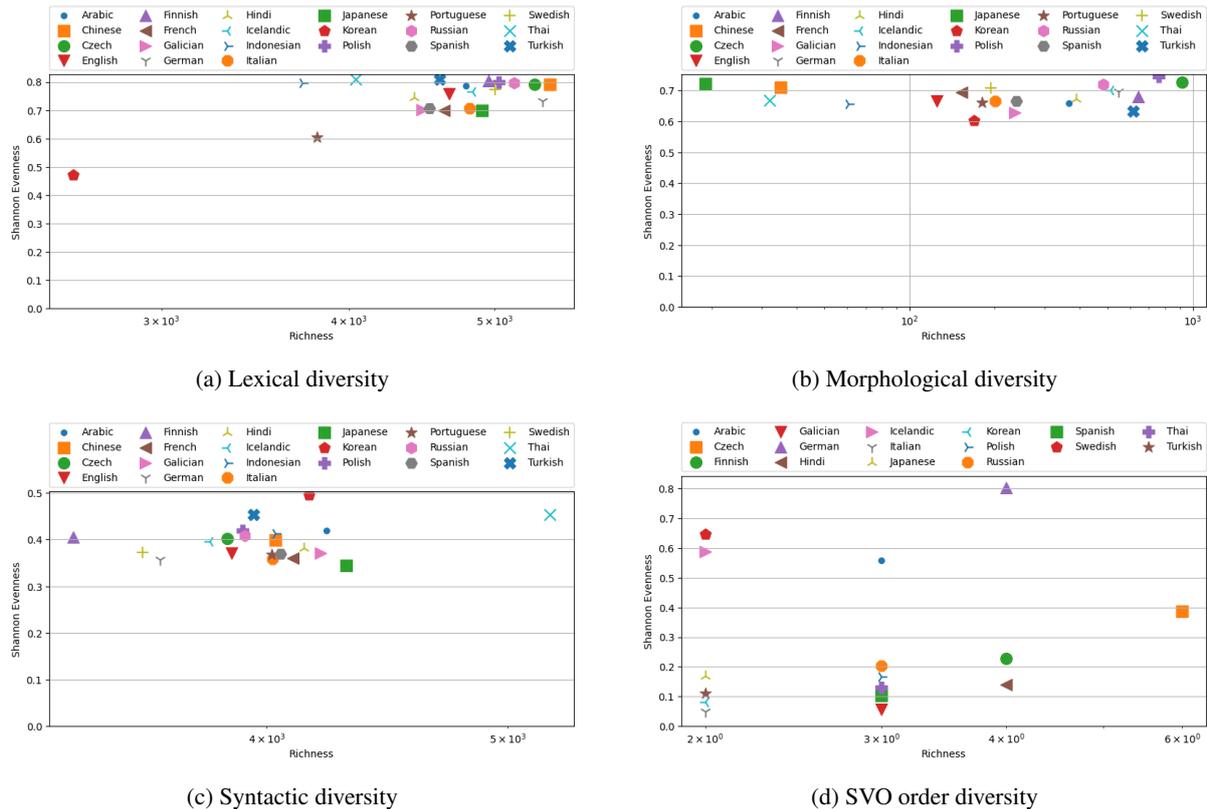


Figure 2: Richness (number of distinct categories) versus balance (uniformity of their distribution) for 24 PUD treebanks across lexical (a), morphological (b), syntactic (c), and word-order (d) levels.

## 4.1 Experimental setup

We applied DELTA to measure diversity across four linguistic levels using the predefined configurations described in §3.1: (1) lexical diversity of lemmas; (2) morphological diversity of POS and feature combinations; (3) syntactic diversity of labeled subtrees; and (4) word order diversity of subject-verb-object placement.

For each linguistic level, we computed two complementary diversity dimensions. For variety, we use richness, which is simply the number of distinct categories  $n$ . For balance, we use Shannon evenness (Smith and Wilson, 1996; Ramaciotti Morales et al., 2021), which normalizes entropy by maximum entropy for  $n$  categories (Equation 1). This metric ranges from 0 (maximally uneven distribution) to 1 (perfectly even distribution).

$$H'(p) = \frac{H(p)}{\log_b(n)} = \frac{-\sum_{i=1}^n p_i \log_b(p_i)}{\log_b(n)} \quad (1)$$

## 4.2 Results

Figure 2 presents diversity scores across four linguistic levels for PUD treebanks.

**Lexical diversity** (Figure 2a) shows most languages clustering with similar numbers of distinct lemmas (4,000-5,000) that are evenly distributed (evenness 0.7-0.8). Korean stands out as a clear outlier with substantially lower richness and evenness, even in comparison to typologically similar Japanese. This pattern warrants further investigation, but is likely influenced by coarser tokenisation granularity in Korean UD, which retains more morphological material within tokens (Chun et al., 2018).

**Morphological diversity** (Figure 2b) exhibits the largest variation in richness (100–1,000 distinct morphological property combinations), reflecting expected typological differences in morphological complexity, with fusional languages (e.g., Czech, Polish) and agglutinative languages (e.g., Turkish, Finnish) exhibiting the highest values, though cross-treebank differences in how morphological information is encoded (e.g., feature inventory size and segmentation practices) may also contribute to the observed variation.

**Syntactic diversity** (Figure 2c), measured here as variety and balance of dependency subtree configurations, shows relatively tight clustering, with

most languages exhibiting comparable number and distribution of such configurations. Finnish and Thai emerge as the two most notable outliers, a pattern that may reflect morphology-syntax trade-offs but requires further investigation to disentangle typological properties from annotation-specific practices.

**Subject-verb-object (SVO) order diversity** (Figure 2d) reveals clear word order typology, with richness distinguishing fixed-order languages (few configurations)<sup>10</sup> from free-order languages like Czech (all six possible permutations). Evenness captures preference strength: languages with similar richness show very different distributions, from strong word order preferences (low evenness) to even distributions across all available patterns.

These results demonstrate DELTA’s capacity for systematic linguistic diversity profiling within and across corpora. While many of the findings above align with established typological patterns, they also highlight the tool’s ability to identify potential outliers or unexpected distributions, making it applicable not only to cross-linguistic comparison but to systematic comparisons across datasets more generally (e.g. between genres within a single language). DELTA thus enables researchers to identify and compare diversity patterns at multiple linguistic levels and from different diversity perspectives within a unified analytical framework.

## 5 Conclusion

We presented DELTA, a unified and configurable framework for computing linguistic diversity of various linguistic features in dependency-parsed corpora. By bridging expressive dependency-tree querying with a broad suite of diversity metrics, visualizations and pre-configured templates, DELTA provides the first integrated environment for systematic, reproducible measurement of diversity for any phenomenon expressible as a dependency (sub)tree – from individual words to complex syntactic patterns. While demonstrated here on standard UD treebanks, the framework’s reliance on the CoNLL-U tabular structure rather than specific annotation content makes it also adaptable to alternative annotation schemes and, consequently, a broader range of linguistic phenomena.

<sup>10</sup>Due to the definition of the balance metric used here, evenness can only be computed when at least two categories are attested. Languages with only one SVO order (e.g., English) therefore do not appear in Figure 2d.

Future work will focus on targeted linguistic research questions across specific languages, genres, and linguistic phenomena, as well as on validating the behaviour and interpretability of the adopted diversity metrics across varying corpus sizes and linguistic conditions. To support this, we will further improve computational efficiency, scalability to new formats and phenomena, and the overall user experience – and we welcome community feedback to guide these ongoing developments.

## Acknowledgments

We gratefully acknowledge financial support from the UniDive project (COST Action CA21167), the SELEXINI project (ANR-21-CE23-0033), the “Plan blanc” doctoral funding from Université Paris-Saclay (France), the SPOT project (ARIS Z6-4617), the LRTS research program (ARIS P6-0411), and AI4DH (EU HORIZON-WIDERA-2023-TALENTS-01-01, grant 101186647). We thank Agata Savary and Thomas Lavergne for their discussions on the topic of this tool. Generative AI tools were used by one author to support language editing during manuscript preparation.

## References

- Dominique Brunato, Andrea Cimino, Felice Dell’Orletta, Giulia Venturi, and Simonetta Montemagni. 2020. [Profiling-UD: a tool for linguistic profiling of texts](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 7145–7151, Marseille, France. European Language Resources Association.
- Anne Chao, Chun-Huo Chiu, and Lou Jost. 2014. [Unifying Species Diversity, Phylogenetic Diversity, Functional Diversity, and Related Similarity and Differentiation Measures Through Hill Numbers](#). *Annual Review of Ecology, Evolution, and Systematics*, 45:297–324. Publisher: Annual Reviews.
- Jayeol Chun, Na-Rae Han, Jena D Hwang, and Jinho D Choi. 2018. [Building universal dependency treebanks in Korean](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Marie-Catherine de Marneffe, Christopher David Manning, Joakim Nivre, and Daniel Zeman. 2021. [Universal dependencies](#). *Computational Linguistics*, 47(2):255–308.
- Kaja Dobrovoljc. 2025. [Counting trees: A treebank-driven exploration of syntactic variation in speech and writing across languages](#).

- Louis Estève, Marie-Catherine de Marneffe, Nurit Melnik, Agata Savary, and Olha Kanishcheva. 2025. [A survey of diversity quantification in natural language processing: The why, what, where and how.](#)
- Kim Gerdes, Sylvain Kahane, and Xinying Chen. 2021. [Typometrics: From Implicational to Quantitative Universals in Word Order Typology.](#) *Glossa: a journal of general linguistics (2021-...)*, 6(1):17.
- Rob Van Der Goot, Esther Ploeger, Verena Blaschke, and Tanja Samardzic. 2025. [DistaLs: a comprehensive collection of language distance measures.](#) In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 307–318, Suzhou, China. Association for Computational Linguistics.
- Miguel R. Guevara, Dominik Hartmann, and Marcelo Mendoza. 2016. [diverse: an r package to analyze diversity in complex systems.](#) *The R Journal*, 8:60–78. <https://rjournal.github.io/>.
- Bruno Guillaume. 2021. [Graph matching and graph rewriting: Grew tools for corpus exploration, maintenance and conversion.](#) In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 168–175.
- Yanzhu Guo, Guokan Shang, Michalis Vazirgiannis, and Chloé Clavel. 2024. [The Curious Decline of Linguistic Diversity: Training Language Models on Synthetic Text.](#) In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3589–3604, Mexico City, Mexico. Association for Computational Linguistics.
- Mark Oliver Hill. 1973. [Diversity and Evenness: A Unifying Notation and Its Consequences.](#) *Ecology*, 54(2):427–432. Publisher: Ecological Society of America.
- Shailza Jolly, Sandro Pezzelle, and Moin Nabi. 2021. [EaSe: A diagnostic tool for VQA based on answer diversity.](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2407–2414, Online. Association for Computational Linguistics.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. [The State and Fate of Linguistic Diversity and Inclusion in the NLP World.](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.
- Dijana Kosmajac and Vlado Keselj. 2019. [Twitter bot detection using diversity measures.](#) In *Proceedings of the 3rd International Conference on Natural Language and Speech Processing*, pages 1–8, Trento, Italy. Association for Computational Linguistics.
- Luka Krsnik and Kaja Dobrovoljc. 2025. [STARK: A toolkit for dependency \(sub\)tree extraction and analysis.](#) In *Proceedings of the 23rd International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2025)*, pages 44–51, Ljubljana, Slovenia. Association for Computational Linguistics.
- Ashutosh Kumar, Satwik Bhattamishra, Manik Bhandari, and Partha Talukdar. 2019. [Submodular optimization-based diverse paraphrasing and its effectiveness in data augmentation.](#) 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 3609–3619, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tom Leinster and Christina A. Cobbold. 2012. [Measuring diversity: the importance of species similarity.](#) *Ecology*, 93(3):477–489.
- Natalia Levshina. 2019. [Token-based typology and word order entropy: A study based on universal dependencies.](#) *Linguistic Typology*, 23(3):533–572.
- Loet Leydesdorff, Caroline S. Wagner, and Lutz Bornmann. 2019. [Interdisciplinarity as diversity in citation patterns among journals: Rao-Stirling diversity, relative variety, and the Gini coefficient.](#) *Journal of Informetrics*, 13(1):255–269.
- Adam Lion-Bouton, Yagmur Ozturk, Agata Savary, and Jean-Yves Antoine. 2022. [Evaluating Diversity of Multiword Expressions in Annotated Text.](#) In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3285–3295, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- 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.
- Juhani Luotolahti, Jenna Kanerva, and Filip Ginter. 2017. [Dep\\_search: Efficient search tool for large dependency parsebanks.](#) In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, pages 255–258, Gothenburg, Sweden. Association for Computational Linguistics.
- Ganapati P. Patil and Charles Taillie. 1982. [Diversity as a Concept and its Measurement.](#) *Journal of the American Statistical Association*, 77(379):548–561. Publisher: [American Statistical Association, Taylor & Francis, Ltd.].
- Esther Ploeger, Wessel Poelman, Andreas Holck Høeg-Petersen, Anders Schlichtkrull, Miryam de Lhoneux, and Johannes Bjerva. 2024. [A principled framework for evaluating on typologically diverse languages.](#)

- Pedro Ramaciotti Morales, Robin Lamarche-Perrin, Raphaël Fournier-S'Niehotta, Rémy Poulain, Lionel Tabourier, and Fabien Tarissan. 2021. [Measuring diversity in heterogeneous information networks](#). *Theoretical Computer Science*, 859:80–115.
- Jai Ram Rideout, Greg Caporaso, Evan Bolyen, Daniel McDonald, Yoshiki Vázquez Baeza, Jorge Cañardo Alastuey, Anders Pitman, Jamie Morton, Qiyun Zhu, Jose Navas, Kestrel Gorlick, Justine Debelius, Zech Xu, Matt Aton, Ilcooljohn, Joshua Shorenstein, Laurent Luce, Will Van Treuren, John Chase, charudatta-navare, Antonio Gonzalez, Colin J. Brislawn, Weronika Patena, Karen Schwarzberg, teravest, Jens Reeder, Igor Sfiligoi, shiffer1, nbresnick, and Dr K. D. Murray. 2025. [scikit-bio/scikit-bio: scikit-bio 0.6.3](#).
- Victoria Rosén, Koenraad De Smedt, Paul Meurer, and Helge Dyvik. 2012. An open infrastructure for advanced treebanking. In *META-RESEARCH Workshop on Advanced Treebanking at LREC2012*, pages 22–29. Hajič, Jan.
- Alfréd Rényi. 1961. [On Measures of Entropy and Information](#). In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, volume 4.1, pages 547–562. University of California Press.
- Tanja Samardzic, Ximena Gutierrez, Christian Bentz, Steven Moran, and Olga Pelloni. 2024. [A measure for transparent comparison of linguistic diversity in multilingual NLP data sets](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3367–3382, Mexico City, Mexico. Association for Computational Linguistics.
- Farhan Samir and Miikka Silfverberg. 2023. [Understanding compositional data augmentation in typologically diverse morphological inflection](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 277–291, Singapore. Association for Computational Linguistics.
- Agata Savary, Sara Stymne, Verginica Barbu Mititelu, Nathan Schneider, Carlos Ramisch, and Joakim Nivre. 2023. [PARSEME meets Universal Dependencies: Getting on the same page in representing multiword expressions](#). *Northern European Journal of Language Technology*, 9.
- Agata Savary, Daniel Zeman, Verginica Barbu Mititelu, Anabela Barreiro, Olesia Căftană, Marie-Catherine de Marneffe, Kaja Dobrovoljc, Gülşen Eryiğit, Voula Giouli, Bruno Guillaume, Stella Markantonatou, Nurit Melnik, Joakim Nivre, Atul Kr. Ojha, Carlos Ramisch, Abigail Walsh, Beata Wójtowicz, and Alina Wróblewska. 2024. [UniDive: A COST action on universality, diversity and idiosyncrasy in language technology](#). In *Proceedings of the 3rd Annual Meeting of the Special Interest Group on Under-resourced Languages @ LREC-COLING 2024*, pages 372–382, Torino, Italia. ELRA and ICCL.
- Samuel M. Scheiner. 2012. [A metric of biodiversity that integrates abundance, phylogeny, and function](#). *Oikos*, 121(8):1191–1202. Number: 8.
- Claude Elwood Shannon. 1948. [A Mathematical Theory of Communication](#). *The Bell System Technical Journal*, 27(4):623–656.
- Claude Elwood Shannon and Warren Weaver. 1949. [The Mathematical Theory of Communication](#). University of Illinois Press, Urbana.
- Benjamin Smith and J. Bastow Wilson. 1996. [A Consumer's Guide to Evenness Indices](#). *Oikos*, 76(1):70–82. Publisher: [Nordic Society Oikos, Wiley].
- Ricard V. Solé, Bernat Corominas-Murtra, and Jordi Fortuny. 2010. [Diversity, competition, extinction: the ecophysics of language change](#). *Interface*, 7(53):1647–1664.
- Feifan Song, Bowen Yu, Hao Lang, Haiyang Yu, Fei Huang, Houfeng Wang, and Yongbin Li. 2024. [Scaling data diversity for fine-tuning language models in human alignment](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 14358–14369, Torino, Italia. ELRA and ICCL.
- Jan Štěpánek and Petr Pajas. 2010. [Querying diverse treebanks in a uniform way](#). In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).
- Andy Stirling. 2007. [A general framework for analysing diversity in science, technology and society](#). *Journal of The Royal Society Interface*, 4(15):707–719. Publisher: Royal Society.
- Luka Terčon and Kaja Dobrovoljc. 2025. [ComparaTree: A multi-level comparative treebank analysis tool](#). In *Proceedings of the 23rd International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2025)*, pages 129–139, Ljubljana, Slovenia. Association for Computational Linguistics.
- Norbert Wiener. 1939. [The ergodic theorem](#). *Duke Mathematical Journal*, 5(1):1–18.
- Daniel Zeman, Martin Popel, Milan Straka, Jan Hajič, Joakim Nivre, Filip Ginter, Juhani Luotolahti, Sampo Pyysalo, Slav Petrov, Martin Potthast, Francis Tyers, Elena Badmaeva, Memduh Gokirmak, Anna Nedoluzhko, Silvie Cinková, Jan Hajič jr., Jaroslava Hlaváčová, Václava Kettnerová, Zdeňka Uřešová, Jenna Kanerva, Stina Ojala, Anna Missilä, Christopher David Manning, Sebastian Schuster, Siva Reddy, Dima Taji, Nizar Habash, Herman Leung, Marie-Catherine de Marneffe, Manuela Sanguinetti, Maria Simi, Hiroshi Kanayama, Valeria de Paiva, Kira Droganova, Héctor Martínez Alonso, Çağrı Çöltekin, Umut Sulubacak, Hans Uszkoreit, Vivien Macketanz, Aljoscha Burchardt, Kim Harris, Katrin Marheinecke, Georg Rehm, Tolga Kayadelen, Mohammed Attia,

Ali Elkahky, Zhuoran Yu, Emily Pitler, Saran Lertpradit, Michael Mandl, Jesse Kirchner, Hector Fernandez Alcalde, Jana Strnadová, Esha Banerjee, Ruli Manurung, Antonio Stella, Atsuko Shimada, Sookyoung Kwak, Gustavo Mendonça, Tatiana Lando, Rattima Nitisaroj, and Josie Li. 2017. [CoNLL 2017 shared task: Multilingual parsing from raw text to Universal Dependencies](#). In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 1–19, Vancouver, Canada. Association for Computational Linguistics.

Daniel Zeman et al. 2025. [Universal dependencies 2.16](#). LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL).

Kaizhong Zhang and Dennis Shasha. 1989. [Simple Fast Algorithms for the Editing Distance between Trees and Related Problems](#). *SIAM Journal on Computing*, 18(6):1245–1262.

## A Appendix

<b>Tree</b>	<b>Freq.</b>
NOUN <nsubj VERB >obj NOUN	142
VERB >nsubj NOUN >obj NOUN	13
NOUN <nsubj NOUN <obj VERB	9
NOUN <obj VERB >nsubj NOUN	5
VERB >obj NOUN >nsubj NOUN	2
NOUN <obj NOUN <nsubj VERB	1

Table 1: Example output TSV listing categories (subject-verb-object trees) for SVO diversity computation in Czech PUD treebank (Figure 2d). Trees are represented using a simplified query-like syntax inspired by the `dep_search` tool (Luotolahti et al., 2017).



# CLARIESG: An End-to-End System for ESG Analysis over Complex Tables in Corporate Reports

Marta Santacroce, Michele Luca Contalbo, Sara Pederzoli, Riccardo Benassi,  
Valeria Venturelli, Matteo Paganelli, Francesco Guerra

University of Modena and Reggio Emilia, Modena, Italy,  
{name.surname}@unimore.it

## Abstract

Sustainability reports contain rich Environmental, Social and Governance (ESG) information, but their heterogeneous layouts and complex multi-table structures pose major challenges for LLMs, especially for unit normalization, cross-document reasoning, and precise numerical computation. We present CLARIESG, an end-to-end system that couples robust table extraction with a structured prompting framework for multi-table filtering, normalization, and program-of-thought reasoning. On ESG-focused multi-table benchmarks, CLARIESG consistently outperforms standard prompting and provides transparent, auditable reasoning, supporting more reliable ESG analysis and greenwashing detection in real-world settings.

## 1 Introduction

Companies increasingly publish *sustainability reports*, i.e., documents that describe their ESG performance, policies, and non-financial impacts, to comply with sustainability standards such as GRI (GRI, 2024). These reports are a key source of ESG data, used in sustainable finance, corporate accountability, and policy evaluation. In Europe alone, assets managed under responsible investment principles reached approximately €6.6 trillion in 2024, representing nearly 38% of total managed assets (Heflich and Saulnier, 2024).

To fully exploit the analytical potential of these disclosures, the European Union is introducing the *European Single Access Point* (CSR, 2022), a unified platform that will begin collecting sustainability reports starting in 2026. While this centralization will enable unprecedented large-scale, data-driven ESG analysis, it also underscores the need for automated methods capable of consistently interpreting complex and heterogeneous disclosures in a transparent and explainable manner.

Recent advances in Large Language Models (LLMs) offer a promising direction for automating

ESG knowledge extraction, thanks to their ability to interpret unstructured and heterogeneous data. Nevertheless, sustainability reports remain highly complex documents that challenge even state-of-the-art models. Evidence of this emerges from the GRI-QA benchmark (Contalbo et al., 2025), a dataset for single- and multi-table question answering over sustainability reports. When tested on multi-table scenarios, GPT-based models show limited performance even with clean, expert-curated tables, indicating that the difficulty lies not in OCR noise or table detection errors but in the intrinsic reasoning demands of the domain. This empirical analysis highlights three main challenges.

First, the tabular data in sustainability reports frequently exhibits non-standard structures, including hierarchical layouts, merged headers, and company-specific schemas. Second, the language used in these documents is highly domain-specific, combining technical terminology with performance indicators that vary in definition, scope, and units of measure. Finally, interpreting the disclosed information can require complex numerical reasoning, involving the combination and comparison of indicators distributed across multiple tables, sometimes even in different documents.

The literature has proposed several end-to-end systems for processing and querying sustainability reports with LLMs (Zou et al., 2023; Ni et al., 2023; Vaghefi et al., 2023; Singh et al., 2024; Wrzalik et al., 2024; Nguyen et al., 2025; Hsu et al., 2024); however, these systems generally do not provide explicit support for tabular content or for complex quantitative reasoning, particularly in scenarios requiring cross-table numerical aggregation and interpretation.

To address these limitations, we present CLARIESG<sup>1</sup>, an end-to-end LLM-based system designed

<sup>1</sup>The code and demonstration video are available at <https://github.com/softlab-unimore/ClariESG.git> and <https://youtu.be/noQ-5ya6cOE>

to support automated analysis of corporate sustainability reports. The architecture consists of a *data management and preparation* layer, which performs document parsing, metadata enrichment, and extraction and normalization of textual and tabular content, and an *analytics layer*, which enables exploration, querying, and numerical reasoning over the extracted data. By bridging the gap between general-purpose LLM capabilities and the analytical requirements of ESG reporting, CLARIESG provides structured access to unstructured reports and supports accurate, explainable, and large-scale ESG analysis within a realistic regulatory context.

The demo illustrates two main scenarios. In the first, *ESG analysts* can query the system in natural language to obtain company and sector-level insights, with results returned both as explanatory text and as structured *scorecards* for benchmarking and decision-making. The second scenario focuses on *claim verification*, addressing the growing concern of *greenwashing* (Nemes et al., 2022; Moodaley and Telukdarie, 2023; de Freitas Netto et al., 2020), i.e., the practice of making unsubstantiated or misleading claims about a company’s environmental performance. In this setting, *regulators, or sustainability promoters* can input a textual claim, and CLARIESG automatically retrieves and analyzes relevant evidence from sustainability reports, highlights supporting or contradicting passages, and provides natural-language justifications.

These scenarios demonstrate how CLARIESG supports analytical and regulatory needs, enhancing the reliability of ESG reporting.

## 2 Related Work

Recent approaches have explored LLM-based analysis of corporate and environmental reports. ESGReveal (Zou et al., 2023) leverages LLM reasoning to extract ESG indicators from both textual and tabular content, ChatReport (Ni et al., 2023) focuses on extracting traceable insights from sustainability reports while reducing hallucinations. ChatClimate (Vaghefi et al., 2023) integrates general-purpose LLM knowledge with climate-specific content, summarizing and distilling information relevant to user queries. FinQAPT (Singh et al., 2024) addresses QA over financial reports by combining dense retrieval, re-ranking, and LLM reasoning with dynamic n-shot prompting and chain-of-thought. NetZeroFacts (Wrzalik et al., 2024) extracts structured emissions data, enabling system-

atic analysis of corporate climate commitments.

Beyond report-centric approaches, some systems integrate further unstructured data sources. MyClimateCopilot (Nguyen et al., 2025) adopts an agentic framework that plans information retrieval, selects tools, and queries heterogeneous sources such as climate APIs and scientific literature. ChatNetZero (Hsu et al., 2024) similarly combines semantic retrieval with anti-hallucination strategies to extract answers from text and spreadsheets.

Another line of research focuses on optimizing individual components rather than providing end-to-end solutions. This includes approaches for extracting semantically structured ESG-related information from sustainability reports using LLMs (Zhou and Perzylo, 2023; Usmanova and Usbeck, 2024; Bronzini et al., 2024), as well as tools for parsing of reports with complex layouts, such as ReportParse (Morio et al., 2024).

Finally, several models have been fine-tuned on corporate and environmental reports, resulting in specialized architectures designed to address a variety of tasks, such as text classification (Xia et al., 2024; Mehra et al., 2022; Schimanski et al., 2023; Webersinke et al., 2021; Araci, 2019; Luukkonen et al., 2023), question answering (Luccioni et al., 2020; Zhao et al., 2022; Xie et al., 2023; Chen et al., 2021; Zhu et al., 2021; Deng et al., 2022; Wu et al., 2025), and claim extraction for greenwashing detection (Mahdavi et al., 2024).

Among the end-to-end systems, ESGReveal is the only approach that explicitly handles tabular structure, whereas FinQAPT uniquely focuses on complex numerical reasoning to effectively address quantitative queries. CLARIESG attempts to integrate both capabilities and resembles UNITQA (Zhu et al., 2025) in multi-table management, though the latter does not handle table extraction and retrieval from unstructured documents, instead assuming that tables are already available in relational or non-relational sources.

## 3 Approach

The architecture of CLARIESG is organized into two main layers, a data management and preparation layer and an analytics layer, as illustrated in Figure 1.

The *data management and preparation layer* handles the acquisition, parsing, and normalization of sustainability reports. Specifically, this layer integrates document-level parsing, company and sec-

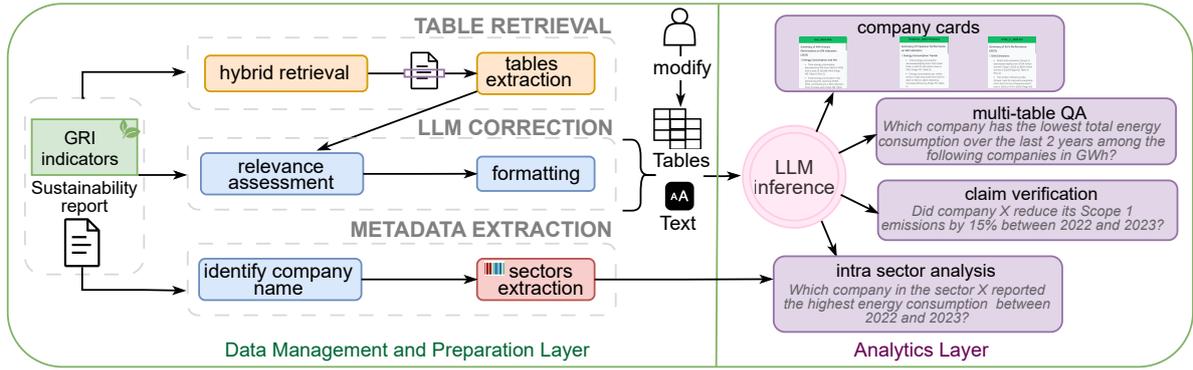


Figure 1: System overview (■ use of LLM)

tor metadata retrieval, and the extraction and standardization of tabular and textual content according to ESG reporting frameworks such as GRI (GRI, 2024). Importantly, it adopts a human-in-the-loop approach (Amershi et al., 2014; Wu et al., 2022): while the process is largely automated, users can supervise the collected data, inspect intermediate results, and correct possible import or formatting errors. This optional supervision ensures reliability and traceability without compromising scalability. By producing consistent, validated representations, this layer lays the foundation for reproducible and large-scale ESG analysis.

The *analytics layer* builds upon these representations to enable interactive reasoning and knowledge discovery. It supports both natural-language querying and analytical operations such as numerical comparison, aggregation, and ranking across multiple documents. To this end, it combines LLM-based text understanding with reasoning paradigms such as Chain-of-Thought (Wei et al., 2022; Kojima et al., 2022) and Program-of-Thought (Chen et al., 2023).

By explicitly separating document understanding from reasoning and exploration, CLARIESG offers a unified framework for scalable, transparent, and explainable ESG analysis, anticipating the data landscape that will emerge with the introduction of the *European Single Access Point* (CSR, 2022).

### 3.1 Data Management and Preparation Layer

The data management and preparation layer is designed to transform raw sustainability reports into structured, validated, and searchable data representations. Its pipeline covers four main responsibilities: (i) identifying the target company and its metadata, (ii) locating the portions of the document that contain relevant information, (iii) extracting and validating tables and contextual text, and (iv)

storing the resulting representations for subsequent retrieval and reasoning.

First, the system analyzes each report to identify the legal company name, which is then used to query Wikidata for sectoral metadata. Second, a hybrid sparse–dense<sup>2</sup> retrieval stage indexes the textual content in a vector database and localizes the sections of the document relevant to GRI topics, focusing the analysis on pages most likely to contain quantitative disclosures. Third, CLARIESG extracts and validates tabular data, which represent a large portion of ESG indicators. Tables in these reports are highly heterogeneous, often including multi-level headers, nested indicators, subtotals, or irregularly merged cells. An ensemble of OCR systems, combining Unstructured<sup>3</sup> and Tesseract (Smith, 2007), is used to extract table content. The extracted tables then undergo a two-step LLM-based refinement: a relevance assessment with respect to GRI indicators, followed by structural and formatting correction. Residual inaccuracies can be reviewed and corrected by the user through a human-in-the-loop interface, which allows inspection and validation against the original document. Finally, both refined tables and associated textual contexts are stored inside the database. Additional details on the prompts used for metadata extraction and table processing are provided in the Appendix (Section 7.1).

### 3.2 Analytics Layer

CLARIESG supports comparative and quantitative analyses that often require integrating information from multiple tables within a single report or across several reports. To address this challenge, CLARIESG implements a prompting-based workflow

<sup>2</sup>multilingual-e5-large-instruct and TF-IDF

<sup>3</sup>Unstructured Documentation

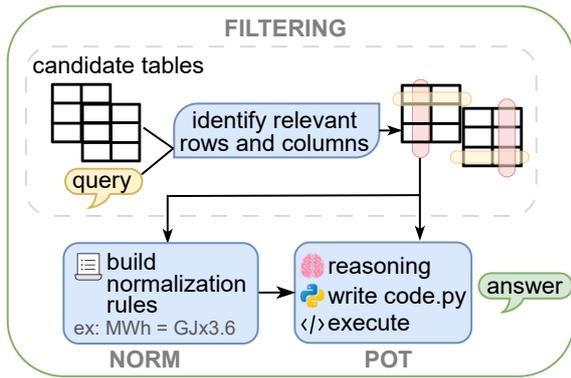


Figure 2: NormPoT prompting strategy. (■ use of LLM)

that orchestrates a sequence of LLM-guided reasoning steps for multi-table and quantitative tasks.

The process begins with *table filtering*, where the system analyzes candidate tables and selects the most relevant rows and columns according to the user query. This step is fully realized through prompting: the model is instructed to rank table segments by relevance to the question. The selected fragments are then passed to the *normalization* stage, where a second prompt aligns units, labels, and conventions to ensure consistency across heterogeneous sources. For example, energy consumption may appear as “GJ” in one report and “MWh” in another; CLARIESG automatically applies conversion factors and aligns terminology to a standard schema, enabling direct comparison.

Next, *program generation* is performed using a Program-of-Thought (PoT) reasoning paradigm. The model is prompted to synthesize a Python function encoding the operations required to aggregate or combine data from the filtered and normalized tables. The generated code is parsed, sanitized, and executed locally to compute the requested metric. Finally, in the *answer composition* phase, the system reformulates the numeric output into a human-readable explanation that references both the computed value and the supporting evidence. We refer to this end-to-end reasoning pipeline as **NormPoT** (Figure 2). The Appendix (Section 7.2) provides the full set of prompts that operationalize each stage of the workflow.

### 3.3 Implementation details

CLARIESG is implemented in Python. Data is stored in a PostgreSQL database extended with the pgvector module. gpt-4o-mini is used as a language model for structured information extraction and reasoning tasks. CLARIESG is independent

of the LLM, allowing the underlying model to be replaced, e.g. with open source alternatives, without modifying CLARIESG’s pipeline. The user interface is built with Gradio 5.46.0, providing an interactive environment for uploading reports, inspecting tables, and executing queries. Interactions with the LLM are handled through the OpenAI API.

## 4 Application scenarios

This section illustrates two representative scenarios in which CLARIESG supports automated analysis of corporate sustainability disclosures: (i) comparative ESG benchmarking, and (ii) claim verification.

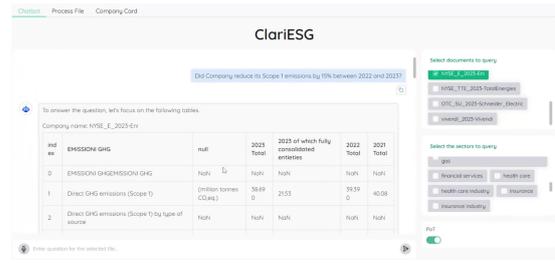
### 4.1 Comparative ESG Analysis

This functionality represents a core step in several downstream applications, including the identification of top-performing and under-performing peers within a sector and the support of investment and policy decisions grounded in comparable evidence. Traditionally, analysts extract KPIs from individual sustainability reports, manually transfer them into spreadsheets, and attempt to harmonize units, time frames, and reporting methodologies. This harmonization process is tedious and prone to inconsistencies, particularly when reports differ in structure, terminology, or indicator granularity.

CLARIESG automates this entire workflow by combining standard reasoning with comparison operations specifically designed for multi-table settings, such as value normalization and table-structure alignment. More precisely, CLARIESG can extract, standardize, and interconnect corresponding indicators from multiple reports. The resulting processed information can be accessed through two complementary modalities: **company scorecards** and **conversational analytical responses**. In the first modality, the system transforms the semi-structured content of sustainability reports into interlinked, machine-readable data aligned with common reporting frameworks, such as the GRI standards. This data provides concise representations of relevant indicators. An example of output produced is shown in Figure 3a. In the second modality, users can interact directly with the system via a conversational interface. For instance, a query such as “Which company reported the highest energy consumption in the manufacturing sector between 2022 and 2023?” triggers the retrieval of the relevant values from each report, exe-



(a) Company scorecards.



(b) Conversational interface.

Figure 3: Key CLARIESG functionalities: automated ESG comparison and interactive claim verification.

cution of the necessary computations, and formulation of a unified answer with references to the original sources. This capability enables reproducible, and transparent quantitative ESG benchmarking, thus supporting interactive and data-driven exploration of corporate sustainability at scale.

## 4.2 Claim verification

The second application domain concerns the verification of claims and the answering of questions within corporate sustainability reports (see the conversation interface in Figure 3b). This task is central to the identification of potential greenwashing, as it enables the detection of inconsistencies between narrative claims and reported evidence. In traditional workflows, verifying such claims requires labor-intensive manual inspection of lengthy reports and the cross-referencing of textual statements with tabular indicators dispersed across multiple sections. This procedure is inherently slow, error-prone, and difficult to scale, particularly when dealing with heterogeneous reporting formats or multiple companies. In contrast, CLARIESG automates the task and enhances its effectiveness. By leveraging numerical reasoning capabilities through a Program-of-Thought paradigm, the system allows users to query ESG disclosures in natural language with questions such as "What is the percentage reduction of greenhouse gas emissions since 2020?". The model identifies the main indicators and temporal references in textual and tabular evidence, computes the required numerical variation, and returns an evidence-grounded answer that explicitly includes the supporting excerpts and table cells. This design enables analysts to audit the full reasoning chain.

## 5 Main results

A core requirement in the proposed scenarios is accurately resolving QA tasks over single and multiple tables. To evaluate this, we assess CLARIESG

on GRI-QA, a domain-specific QA benchmark over environmental tables from sustainability reports, which allows us to replicate the core operation behind both scenarios in a controlled setting. While GRI-QA specifically targets GRI 300 indicators, CLARIESG is GRI-agnostic and can be extended to other GRI families by updating indicator metadata. GRI-QA organizes the questions into the following categories: *extractive* questions that require direct data retrieval; *hierarchical* questions that involve disambiguating terms within nested table structures; and *calculated* and *quantitative* questions that test relational and arithmetic reasoning such as comparisons, superlatives, rankings, and percentage variations. It also includes *multi-step* questions requiring computations over multiple tables or documents.

By analyzing CLARIESG's responses on GRI-QA, we assess (i) the effectiveness of prompting strategies for *single-table* and *multi-table* questions, and (ii) the relative performance of ChatGPT 5.1 and CLARIESG in multi-table reasoning. We use the normalized Exact Match (EM) (Dua et al., 2019) as the main evaluation metric.

**Comparison of Prompting Strategies.** Table 1 shows that the simple use of Chain-of-Thought (CoT) on *one-table* questions provides the best performance. In particular, the average performance of Program-of-Thoughts (PoT) and NormPoT decreases by 7 and 7.9 EM points respectively, demonstrating that overly complex prompting strategies can lead to sub-optimal performance on simple questions. The only *one-table* scenario where the performance of PoT and NormPoT improves compared to CoT is for the quant dataset, where performance increases by 13.1 and 13.5 points respectively. This indicates that for questions requiring mathematical calculations, performing the calculations through a Python interpreter leads to better results. For the *multi-table*

	<i>GRI-QA one-table</i>						<i>GRI-QA multi-table</i>									
	extra	hier	rel	quant	step	avg	rel2	rel3	rel5	quant2	quant3	quant5	step2	step3	step5	avg
CoT	84.2	80.9	92.7	72.6	33.1	72.7	56.6	34.3	19.5	58.7	20.8	0.0	43.7	32.7	25.5	32.4
PoT	62.4	66.8	89.9	85.7	23.5	65.7 <sup>-7.0</sup>	63.0	41.0	26.4	65.3	36.1	12.0	58.9	37.0	30.9	41.2 <sup>+8.8</sup>
NormPoT	63.9	62.2	91.5	86.1	20.5	64.8 <sup>-7.9</sup>	68.5	59.0	39.1	69.3	50.0	22.0	56.3	42.8	28.2	48.4 <sup>+16.0</sup>

Table 1: Performance of different prompting strategies in one-table and multi-table settings, and for the question categories defined in GRI-QA. In multi-table tasks, the number beside each category indicates the tables involved (e.g., rel5 = 5 tables). Superscripts denote average performance differences relative to the CoT baseline.

		rel2	rel3	rel5	quant2	quant3	quant5	step2	step3	step5	avg
ChatGPT 5.1		70.0	50.0	30.0	50.0	16.0	28.0	60.0	40.0	36.0	42.2
CLARIESG	gpt-4o-mini	60.0	56.0	44.0	72.0	46.0	22.0	72.0	44.0	26.0	49.1 <sup>+6.9</sup>
	gpt-4o-mini + noisy tables	44.0	28.0	30.0	68.0	28.0	0.0	54.0	40.0	30.0	35.8 <sup>-6.4</sup>
	gpt-5-mini	90.0	88.0	88.0	92.0	74.0	68.0	76.0	70.0	66.0	79.1 <sup>+36.9</sup>
	gpt-5-mini + noisy tables	88.0	90.0	84.0	88.0	74.0	66.0	76.0	72.0	60.0	77.6 <sup>+35.4</sup>

Table 2: Performance comparison between ChatGPT 5.1 and CLARIESG (with different LLMs) on the first 50 questions of each GRI-QA *multi-table* benchmark. In the noisy setting, two irrelevant tables are added for each company report. Superscripts indicate average performance deltas vs. ChatGPT 5.1.

benchmarks, on the other hand, PoT and NormPoT achieve average performance that is 8.8 and 16 EM points higher than CoT, respectively. In particular, the integration of a normalization step prior to executing PoT provides significant advantages, by clarifying the intermediate steps required to compare values from different companies that use different units of measurement. The results indicate that for questions requiring numerical calculation or reasoning across multiple tables, the best strategy to adopt is NormPoT.

*Comparison with ChatGPT 5.1.* To validate the quality of CLARIESG, we compare its performance on the first 50 questions of each *multi-table* benchmark of GRI-QA with ChatGPT 5.1. The systems are compared based on how they would be used to perform ESG analysis of corporate reports. For ChatGPT 5.1, for each question, we manually connect to the ChatGPT website, we load the complete reports of the companies required by the question and we annotate its response. Based on the request, ChatGPT 5.1 itself decides whether to think longer (*ChatGPT Thinking*) or provide an immediate answer (*ChatGPT Instant*). For CLARIESG, we use the clean tables provided by GRI-QA. Although the amount of textual context provided as input differs between ChatGPT 5.1 and CLARIESG, the comparison between the two systems is fair, assuming that the user correctly cleans the tables extracted by CLARIESG. Still, to faithfully

evaluate the performance of CLARIESG, we also test it with two additional *noisy* tables as context for each company needed to answer the question.

Table 2 shows the results. In CLARIESG, the average performance of gpt-4o-mini surpasses the performance of ChatGPT 5.1 by 6.9 EM points when CLARIESG is not provided with additional *noisy* tables, whereas its performance falls 6.4 EM points below that of ChatGPT 5.1 when evaluated under the *noisy* setting. By using gpt-5-mini as backbone LLM, CLARIESG greatly outperforms ChatGPT 5.1 with a respective average performance increase of 36.9 and 35.4 EM points for the *clean* and *noisy* settings. Notably, providing clean tables and correct context allows CLARIESG to outperform ChatGPT 5.1 with both gpt-4o-mini and gpt-5-mini, even if the backbone LLM is much smaller. In general, the tabular data management and reduction of context performed by CLARIESG proves to be crucial in providing accurate responses to ESG-related queries.

## 6 Conclusion

We showcased CLARIESG, an end-to-end system for analyzing corporate sustainability reports. Combining robust table extraction with structured prompting for multi-table normalization and Program-of-Thoughts reasoning, CLARIESG provides precise, auditable analytics for ESG benchmarking and claim verification. Experiments on

GRI-QA show that this specialized workflow outperforms general-purpose LLMs such as ChatGPT 5.1. Future work will focus on improving robustness to noisy tables and integrating richer domain knowledge.

## Limitations

The system currently focuses exclusively on the management and extraction of information related to ESG data. This represents an essential step for enabling analysts to gain a deeper understanding of companies' environmental behaviour and to compare performance across sectors. However, integrating indicators that combine ESG and financial data would further enhance the analytical value of the system, as investment decisions are often guided by a combination of both dimensions. We plan to address this limitation in future work.

OpenAI models accessed via API calls are known to produce non-deterministic outputs even when the temperature is set to 0. As a result, the results reported in Table 1 and Table 2 may exhibit slight variability across different runs.

## Risks

A potential risk associated with the use of CLARIESG is that analysts may over-rely on the system's responses. Although the performance of CLARIESG is promising (Table 2), it is not flawless. Even though the reasoning process used to generate answers is fully auditable, users may still place trust in the output without verifying the underlying evidence. For this reason, we recommend that analysts consult CLARIESG as a support tool, but cross-check its answers against the original sources to prevent misinterpretations and mitigate the possibility of hallucinations.

## Use of AI Assistants

When writing this paper, we used AI assistants, such as ChatGPT, to improve the flow of writing and the vocabulary of the initial drafts we manually wrote. Each suggestion has been manually validated by the authors.

## Acknowledgments

We acknowledge the CINECA award under the ISCRA initiative, for the availability of high performance computing resources and support. Moreover, this work was partially funded by the RESISTO project (PR-FESR Emilia-Romagna 2021-

2027) through a grant to the AIRI research center at the University of Modena and Reggio Emilia.

## References

2022. Directive (eu) 2022/2464 of the european parliament and of the council of 14 december 2022 on corporate sustainability reporting (csrd). <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32022L2464>. Official Journal of the European Union, L 322, 16 December 2022.
- Saleema Amershi, Maya Cakmak, William Bradley Knox, and Todd Kulesza. 2014. [Power to the people: The role of humans in interactive machine learning](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, pages 483–490. AAAI Press.
- Dogu Araci. 2019. Finbert: Financial sentiment analysis with pre-trained language models. *CoRR*, abs/1908.10063.
- Marco Bronzini, Carlo Nicolini, Bruno Lepri, Andrea Passerini, and Jacopo Staiano. 2024. Glitter or gold? deriving structured insights from sustainability reports via large language models. *EPJ Data Sci.*, 13(1):41.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021. [FinQA: A dataset of numerical reasoning over financial data](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Michele Luca Contalbo, Sara Pederzoli, Francesco Del Buono, Venturelli Valeria, Francesco Guerra, and Matteo Paganelli. 2025. [GRI-QA: a comprehensive benchmark for table question answering over environmental data](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 15764–15779, Vienna, Austria. Association for Computational Linguistics.
- Sebastião Vieira de Freitas Netto, Marcos Felipe Falcão Sobral, Ana Regina Bezerra Ribeiro, and Gleibson Robert da Luz Soares. 2020. [Concepts and forms of greenwashing: a systematic review](#). *Environmental Sciences Europe*, 32(1):19.
- Yang Deng, Wenqiang Lei, Wenxuan Zhang, Wai Lam, and Tat-Seng Chua. 2022. PACIFIC: towards proactive conversational question answering over tabular and textual data in finance. In *EMNLP*, pages 6970–6984. Association for Computational Linguistics.

- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). 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 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- GRI. 2024. [Global reporting initiative website](#).
- A. Heflich and J. Saulnier. 2024. [Potential economic impact of european sustainable finance](#). Technical report, European Parliamentary Research Service.
- Angel Hsu, Mason Laney, Ji Zhang, Diego Manya, and Linda Farczadi. 2024. [Evaluating ChatNet-Zero, an LLM-chatbot to demystify climate pledges](#). In *Proceedings of the 1st Workshop on Natural Language Processing Meets Climate Change (ClimateNLP 2024)*, pages 82–92, Bangkok, Thailand. 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 *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Alexandra Luccioni, Emily Baylor, and Nicolas Duchene. 2020. Analyzing sustainability reports using natural language processing. *CoRR*, abs/2011.08073.
- Risto Luukkonen, Ville Komulainen, Jouni Luoma, Anni Eskelinen, Jenna Kanerva, Hanna-Mari Kupari, Filip Ginter, Veronika Laippala, Niklas Muennighoff, Aleksandra Piktus, Thomas Wang, Nouamane Tazi, Teven Scao, Thomas Wolf, Osmo Suominen, Samuli Sairanen, Mikko Merioksa, Jyrki Heinonen, Aija Vah-tola, and 2 others. 2023. [FinGPT: Large generative models for a small language](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2710–2726, Singapore. Association for Computational Linguistics.
- Mohammad Mahdavi, Ramin Baghaei Mehr, and Tom Debus. 2024. Combat greenwashing with goalspotter: Automatic sustainability objective detection in heterogeneous reports. In *CIKM*, pages 4752–4759. ACM.
- Srishti Mehra, Robert Louka, and Yixun Zhang. 2022. ESGBERT: language model to help with classification tasks related to companies environmental, social, and governance practices. *CoRR*, abs/2203.16788.
- Wayne Moodaley and Arnesh Telukdarie. 2023. [Greenwashing, sustainability reporting, and artificial intelligence: A systematic literature review](#). *Sustainability*, 15(2).
- Gaku Morio, Soh Young In, Jungah Yoon, Harri Rowlands, and Christopher D. Manning. 2024. Report-parse: A unified NLP tool for extracting document structure and semantics of corporate sustainability reporting. In *IJCAI*, pages 8749–8753. ijcai.org.
- Noémi Nemes, Stephen J. Scanlan, Pete Smith, Tone Smith, Melissa Aronczyk, Stephanie Hill, Simon L. Lewis, A. Wren Montgomery, Francesco N. Tubiello, and Doreen Stabinsky. 2022. [An integrated framework to assess greenwashing](#). *Sustainability*, 14(8).
- Vincent Nguyen, Willow Hallgren, Ashley Harkin, Mahesh Prakash, and Sarvnaz Karimi. 2025. [My climate CoPilot: A question answering system for climate adaptation in agriculture](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 62–70, Vienna, Austria. Association for Computational Linguistics.
- Jingwei Ni, Julia Anna Bingler, Chiara Colesanti Senni, Mathias Kraus, Glen Gostlow, Tobias Schimanski, Dominik Stammach, Saeid Ashraf Vaghefi, Qian Wang, Nicolas Webersinke, Tobias Wekhof, Tingyu Yu, and Markus Leippold. 2023. CHATREPORT: democratizing sustainability disclosure analysis through llm-based tools. In *EMNLP (Demos)*, pages 21–51. Association for Computational Linguistics.
- Tobias Schimanski, Julia Anna Bingler, Mathias Kraus, Camilla Hyslop, and Markus Leippold. 2023. Climatebert-netzero: Detecting and assessing net zero and reduction targets. In *EMNLP*, pages 15745–15756. Association for Computational Linguistics.
- Kuldeep Singh, Simerjot Kaur, and Charese Smiley. 2024. Finqapt: Empowering financial decisions with end-to-end llm-driven question answering pipeline. In *ICAIF*, pages 266–273. ACM.
- R. Smith. 2007. [An overview of the tesseract ocr engine](#). In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, volume 2, pages 629–633.
- Aida Usmanova and Ricardo Usbeck. 2024. [Structuring sustainability reports for environmental standards with LLMs guided by ontology](#). In *Proceedings of the 1st Workshop on Natural Language Processing Meets Climate Change (ClimateNLP 2024)*, pages 168–177, Bangkok, Thailand. Association for Computational Linguistics.
- Saeid Ashraf Vaghefi, Qian Wang, Veruska Muccione, Jingwei Ni, Mathias Kraus, Julia Anna Bingler, Tobias Schimanski, Chiara Colesanti Senni, Nicolas Webersinke, Christian Huggel, and Markus Leippold. 2023. chatclimate: Grounding conversational AI in climate science. *CoRR*, abs/2304.05510.
- Nicolas Webersinke, Mathias Kraus, Julia Anna Bingler, and Markus Leippold. 2021. Climatebert: A pretrained language model for climate-related text. *CoRR*, abs/2110.12010.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Marco Wrzalik, Florian Faust, Simon Sieber, and Adrian Ulges. 2024. **NetZeroFacts: Two-stage emission information extraction from company reports**. In *Proceedings of the Joint Workshop of the 7th Financial Technology and Natural Language Processing, the 5th Knowledge Discovery from Unstructured Data in Financial Services, and the 4th Workshop on Economics and Natural Language Processing*, pages 70–84, Torino, Italia. Association for Computational Linguistics.
- Tianyi Wu, Wei Fan, Junjie Wu, and Hui Xiong. 2022. **A survey on human-in-the-loop machine learning: Challenges and opportunities**. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 16(6):1–31.
- Xianjie Wu, Jian Yang, Linzheng Chai, Ge Zhang, Jiaheng Liu, Xeron Du, Di Liang, Daixin Shu, Xianfu Cheng, Tianzhen Sun, Tongliang Li, Zhoujun Li, and Guanglin Niu. 2025. **Tablebench: A comprehensive and complex benchmark for table question answering**. In *AAAI*, pages 25497–25506. AAAI Press.
- Lei Xia, Mingming Yang, and Qi Liu. 2024. **Using pre-trained language model for accurate ESG prediction**. In *Proceedings of the Eighth Financial Technology and Natural Language Processing and the 1st Agent AI for Scenario Planning*, pages 1–22, Jeju, South Korea. -.
- Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. 2023. **Pixiu: A comprehensive benchmark, instruction dataset and large language model for finance**. In *Advances in Neural Information Processing Systems*, volume 36, pages 33469–33484. Curran Associates, Inc.
- Yilun Zhao, Yunxiang Li, Chenying Li, and Rui Zhang. 2022. **Multihirtt: Numerical reasoning over multi hierarchical tabular and textual data**. In *ACL (1)*, pages 6588–6600. Association for Computational Linguistics.
- Yuchen Zhou and Alexander Perzylo. 2023. **Ontosustain: Towards an ontology for corporate sustainability reporting**. In *ISWC (Posters/Demos/Industry)*, volume 3632 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. **TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance**. 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 3277–3287, Online. Association for Computational Linguistics.
- Jun-Peng Zhu, Peng Cai, Kai Xu, Li Li, Yishen Sun, Shuai Zhou, Haihuang Su, Liu Tang, and Qi Liu. 2025. **UNITQA: A unified automated tabular question answering system with multi-agent large language models**. In *SIGMOD Conference Companion*, pages 279–282. ACM.
- Yi Zou, Mengying Shi, Zhongjie Chen, Zhu Deng, ZongXiong Lei, Zihan Zeng, Shiming Yang, HongXiang Tong, Lei Xiao, and Wenwen Zhou. 2023. **Esgreveal: An llm-based approach for extracting structured data from ESG reports**. *CoRR*, abs/2312.17264.

## 7 Prompts and screenshots

The screenshot in [Figure 4](#) shows the component used to upload and refine the tables extracted from the reports, which could not be included in the main paper due to space constraints.

Below, instead, we provide details on the prompts used to instruct the underlying LLM of CLARIESG to perform the different tasks required by the system. We distinguish between prompts employed for data preparation and those used for analysing the extracted data.

### 7.1 Prompts for the data management and preparation layer

During the pre-processing phase, the LLM is responsible for (i) extracting metadata about the reporting company (such as the legal name and industrial sector), and (ii) accurately identifying the tables contained in the document. Specifically, [Figure 5](#) shows the prompt used to identify the company’s legal name from the front pages of the report. The extracted name is then used to query **Wiki-data** and retrieve the company’s industrial sectors. The SPARQL query used for this retrieval is provided in [Figure 6](#). To ensure robust table extraction, OCR output is further processed through a two-step LLM-based pipeline. This includes (i) verifying the relevance of the extracted content with respect to GRI indicators (see the prompt in [Figure 7](#)), and (ii) refining the structural and formatting consistency of the resulting tables (see the prompt in [Figure 8](#)).

### 7.2 Prompts for the analytics layer

To support comparative analysis over tables extracted from multiple reports, CLARIESG orchestrates a sequence of LLM-guided reasoning steps. These include: (i) *table filtering*, to select rows and

columns relevant to the user query (see the prompt in Figure 9); (ii) *normalization*, to harmonize units, labels, and formatting conventions across heterogeneous sources (see the prompt in Figure 10); (iii) *program generation*, which synthesizes a Python function encoding the required logical and arithmetic operations in a PoT-style workflow (see the prompt in Figure 12); and (iv) *code execution*, to run the generated program and obtain the final computed values. In case of Python execution error, CLARIESG falls back to standard CoT (see the prompt in Figure 11).

Additionally, CLARIESG uses the prompt in Figure 13 to create the *company scorecards*.

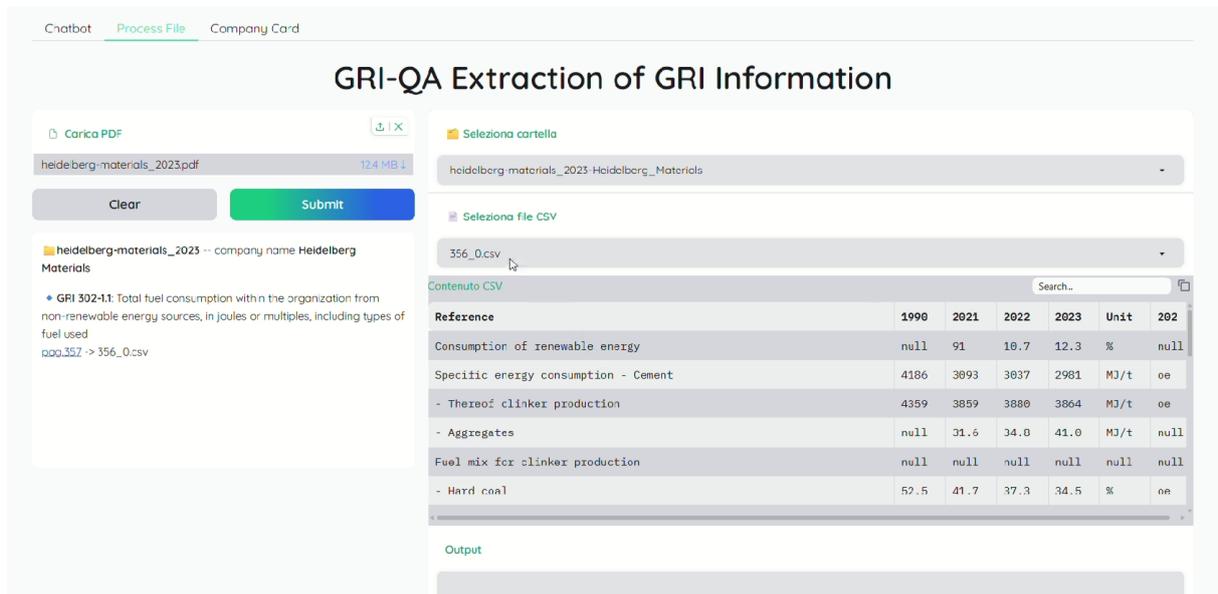


Figure 4: Screen of Tab 2, Upload and processing documents.

**Get company name**

You are an assistant that extracts the name of the main company mentioned in a PDF document.

Read the following text and return **ONLY** the company's full name — no explanations, no punctuation, no additional text.

Figure 5: Prompt to obtain the company name.

### Get company sectors

```
SELECT ?company ?companyLabel ?industry ?industryLabel WHERE {{
  SERVICE wikibase:mwapi {{
    bd:serviceParam wikibase:endpoint "www.wikidata.org";
    wikibase:api "EntitySearch";
    mwapi:search "company_name";
    mwapi:language "en".
    ?company wikibase:apiOutputItem mwapi:item .
  }}
  OPTIONAL {{ ?company wdt:P452 ?industry. }}
  SERVICE wikibase:label{{ bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en". }}
}}
LIMIT 10
```

Figure 6: Wikidata query to retrieve the industry sectors associated with a company.

### Evaluation of extracted tables

You are an expert in sustainability reporting (GRI Standards).

I will give you:

1. A GRI code and its description.
2. The content of a CSV table extracted from a company report.

Task: Decide if this CSV table is relevant to the GRI code.

Answer with ONLY one word: "YES" if the CSV contains information that matches or supports the GRI description, otherwise "NO".

GRI code: gri\_code

Description: gri\_desc

CSV content (partial preview): csv\_preview

Figure 7: Prompt to evaluate the extracted tables.

### Formatted table

You are given the content of a CSV file automatically extracted from a table.  
Your task is to clean and reformat it into a valid table, ensuring that **all rows have the same number of columns**.

Follow these rules strictly:

- Use ; as the column separator in the final output.
- Determine the **maximum number of fields** present in any row, and expand all rows to that length.
- If a row has missing cells, fill them with NaN.
- Keep numeric values as-is, including negative percentages and decimals.
- Fix broken or merged cells, misplaced values, or incorrect headers.
- **Do not add or remove data rows** except for lines that are completely empty or contain only NaN.
- Standardize headers:
  - Create clear, readable names.
  - Avoid duplicates (rename automatically if needed).
  - Do not lose or shorten the meaning of headers.
- Ensure consistent formatting:
  - Align numeric and text values properly.
  - Remove symbols or characters that are clearly OCR or extraction noise.
- Output **only the cleaned CSV content** no explanations or comments.

**REMEMBER THAT ALL ROWS MUST HAVE THE SAME NUMBER OF FIELDS!**

Figure 8: Prompt to format the extracted tables.

### Extract values

You will be given a question and a table.

ONLY IF there are relevant rows and columns, you must indicate the indices of the rows and columns that could be relevant to answer the question. OTHERWISE, if for a certain table there are no relevant rows and columns, write an empty list for both "rows" and "columns". You must not try to answer the question, you must only retrieve the relevant rows if there are. Use the values in the "index" column to refer to the relevant rows.

Additionally, for each selected row include the corresponding row name in the table: use the value from the first non-index column (immediately to the right of "index") as the row's name. Align "row\_names" with "rows". If no such column exists, use an empty string ("").

For column indices, write the number (starting from 0, left to right), not the column name. First reason step-by-step. Then write "Final answer: " followed exclusively by a Python dictionary:

```
{  
  "rows": [row_index1,...,row_indexn],  
  "columns": [column_index1,...,column_indexn],  
  "row_names": [row_name1,...,row_namen]  
}
```

If no relevant rows/columns, return empty lists. Do not write anything else after "Final answer:". Do not use Markdown syntax.

Question: {question}

Table: {table}

Let's think step-by-step.

Figure 9: Prompt to extract relevant rows and columns from a table.

### Normalization

Given multiple tables and a question, decide the unit of measure to use for the final answer. Then, align table values by converting needed values to a unique unit.

If the question specifies a unit, convert values to it. Otherwise, decide the unit and convert. Do not rewrite the tables. Only provide a list of rules/formulas indicating the needed transformations. Transformations must only handle units. Do not discuss solving the question.

Sample rule: 1. 1000 meters = 1 kilometer

First reason step-by-step. Then write "Final answer: " followed exclusively by the list of rules/formulas.

Do not write anything else after "Final answer:". Do not use Markdown syntax.

Question: {question}

Tables: {tables}

Let's think step-by-step.

Figure 10: Normalize units in multiple tables.

### Fallback Python

Consider the following question and content. First reason step-by-step, then provide the answer.

Question: {question}

Content: {content}

Let's think step-by-step.

Figure 11: Prompt to reason step-by-step and provide Python answer.

### Python Prompt

You need to create Python code that answers the following question, taking into account the tables provided and the fact that NOT ALL rows are always useful for generating the answer. Write your reasoning first. Then, at the end, write 'Final answer:' followed by the Python code and nothing else. The Python code must be executable 'as is', so include relevant imports. At the end, print the result with print(). If not already done, specify `python` before the code and ````` at the end.

If the question is Boolean, the output must be exclusively 'yes' or 'no'. If a list of values is required, respond with a comma-separated list. Write numerical values with exactly 2 decimal places.

Ensure the final answer is in the expected form. Do not write anything else after 'Final answer:'. Do not use Markdown syntax.

Question: {question}

Tables: {paragraph}

Let's think step by step.

Figure 12: Prompt to generate Python code considering relevant rows/columns in tables.

### Company Summary

You are an expert assistant in sustainability and GRI standards.

Your task is to analyze data extracted from a company's PDFs in the form of CSV tables related to specific GRI indicators, and provide a clear, concise summary of the company's performance.

Instructions:

- Base your summary strictly on the data provided in the CSV tables.
- Highlight trends, improvements, or regressions in the company's performance where possible.
- Do not add assumptions or information not present in the tables.
- For each key point, reference the row, cell, page, and table number used from the CSV context.
- Make the summary concise, well-structured, and readable for stakeholders.
- If there is no context, reply clearly that you have not received any information. Nothing else.

Here are the CSV tables extracted from the company's PDFs related to GRI indicators:

{context}

Please provide a concise summary of the company's performance based strictly on this data.

Figure 13: Prompt to generate a concise summary of company performance from GRI-related CSV tables as company card.

# Fact Finder - Enhancing Domain Expertise of Large Language Models by Incorporating Knowledge Graphs

Daniel Steinigen<sup>1</sup>  
Roman Teucher<sup>1</sup>  
Timm Heine Ruland<sup>1</sup>  
Max Rudat<sup>1</sup>  
Nicolas Flores-Herr  
Fraunhofer IAIS  
Sankt Augustin, Germany

<first-name>.<last-name>@iaais.fraunhofer.de

<sup>1</sup>All authors contributed equally to this work

Peter Fischer<sup>2</sup>  
Nikola Milosevic<sup>2</sup>  
Christopher Schymura<sup>2</sup>  
Angelo Ziletti<sup>2</sup>  
Bayer AG

<first-name>.<last-name>@bayer.com

<sup>2</sup>Shared senior authorship

## Abstract

Recent advancements in Large Language Models (LLMs) have showcased their proficiency in answering natural language queries. However, their effectiveness is hindered by limited domain-specific knowledge, raising concerns about the reliability of their responses. We introduce a hybrid system that augments LLMs with domain-specific knowledge graphs (KGs), thereby aiming to enhance factual correctness using a KG-based retrieval approach. We focus on a medical KG to demonstrate our methodology, which includes (1) pre-processing, (2) Cypher query generation, (3) Cypher query processing, (4) KG retrieval, and (5) LLM-enhanced response generation. We evaluate our system on a curated dataset of 69 samples, achieving a precision of 78% in retrieving correct KG nodes. Our findings indicate that the hybrid system surpasses a standalone LLM in accuracy and completeness, as verified by an LLM-as-a-Judge evaluation method. This positions the system as a promising tool for applications that demand factual correctness and completeness, such as target identification — a critical process in pinpointing biological entities for disease treatment or crop enhancement. Moreover, its intuitive search interface and ability to provide accurate responses within seconds make it well-suited for time-sensitive, precision-focused research contexts. We publish the source code together with the dataset and the prompt templates used<sup>1</sup>.

## 1 Introduction

Recently, Large Language Models (LLMs) have enabled sophisticated question-answering systems, revolutionizing the landscape of natural language processing (OpenAI, 2023; Jiang et al., 2023a;

<sup>1</sup><https://github.com/chrschy/fact-finder>

Bubeck et al., 2023; Park et al., 2023; Touvron et al., 2023; Katz et al., 2023). These advanced models, with their ability to understand and generate human-like text, have shown great potential in various domains, including life sciences (Nori et al., 2023; Waisberg et al., 2023; Bašaragin et al., 2024). However, LLMs are limited by the timeframe of their training data and can produce incorrect statements, known as hallucinations (Ji et al., 2023), or incomplete answers by providing only a few relevant entities while missing others not included in their internal knowledge.

In domains such as life sciences, obtaining answers with current and factual information is paramount for many use cases (Malaviya et al., 2023; Ljajić et al., 2024). Factually correct AI-generated reviews can aid researchers in information retrieval and hypothesis building. For example, target identification requires up-to-date knowledge of the latest literature. Target identification involves pinpointing a biological entity, such as a gene or protein, that can be manipulated to achieve a desired effect, like treating a disease in humans (pharmaceuticals) or improving crop resilience in plants (crop sciences). Similarly, designing effective field or clinical trials requires considering up-to-date information. For field trials in crop sciences, this may include environmental and climatic conditions, market developments, and regulatory requirements. For clinical trials and medical writing, relevant information includes details about the drug under development, market conditions, planned sites, and regulatory requirements. Current and timely information is crucial for competitive intelligence, including insights on competitor products, disease epidemiology, and market size. AI-based solutions can assist in identifying concurrent disease occurrences and help researchers develop hypotheses

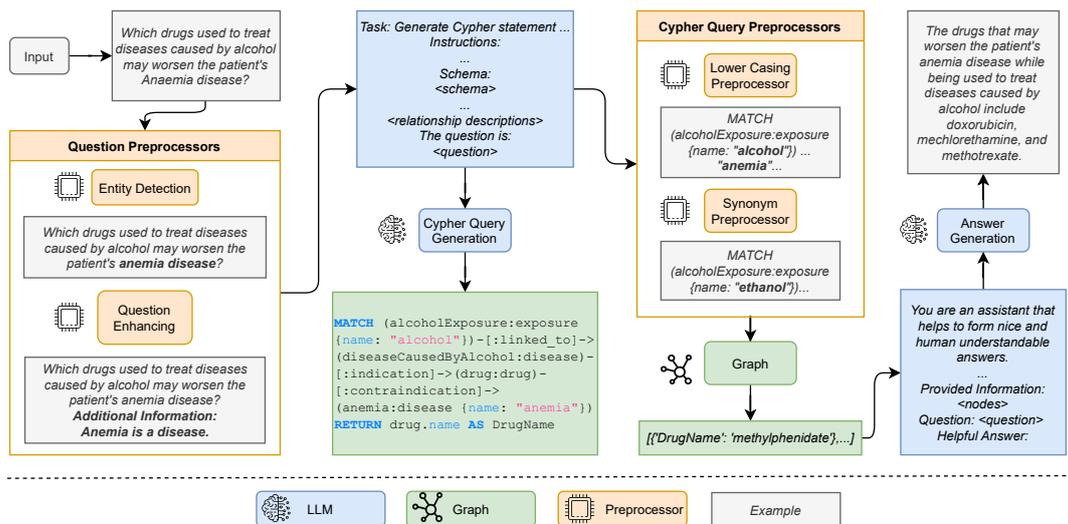


Figure 1: Overview of the FactFinder pipeline using large language models and knowledge graphs to answer scientific questions.

using real-world data.

Knowledge Graphs (KGs) represent a promising strategy for improving factual correctness in LLMs (Jiang et al., 2023b; Baek et al., 2023; Sen et al., 2023), including the life science domain (Feng et al., 2024). By organizing entities such as drugs, diseases, and genes, along with their relationships, into a structured network, KGs provide useful additional context for LLMs for precise and relevant information retrieval (Chandak et al., 2023; Milošević and Thielemann, 2023; Badenes-Olmedo and Corcho, 2023). This organized data framework allows LLMs to produce more factually accurate and comprehensive responses (Pan et al., 2024). In addition, KGs enable systems to leverage current and comprehensive information, including recent data not available during the LLMs’ training phase. From an organizational standpoint, integrating KGs enhances top-tier LLMs with proprietary or specialized knowledge. This integration facilitates the inclusion of unique organizational data sources, such as historical and ongoing lab experiments or licensed datasets.

In this paper, we present FactFinder - a hybrid question answering (QA) system - which leverages both KG and LLM to provide answers to scientific questions. We use the Neo4j graph database<sup>2</sup> to provide the KG and Cypher<sup>3</sup> as query language, and build our pipeline borrowing components from

Langchain<sup>4</sup>. Fig. 1 depicts the system’s architecture, which is structured as a pipeline with several subcomponents. Our main contributions are:

- We provide an easy-to-use system that answers scientific questions combining LLMs and KGs, with automatic subgraph visualizations that enhance interpretability by showing graph-native evidence for each answer.
- We release a dataset of manually annotated text-to-Cypher query pairs, which could serve as benchmark for validating text-to-Cypher conversion system.
- We present a methodology showing that current state-of-the-art LLMs are able to generate satisfactory Cypher queries for the life science domain.
- We share our dataset, source code and prompt templates<sup>5</sup>.

## 2 Data

**Knowledge graph.** We use PrimeKG (Chandak et al., 2023) as our source of fact-based background knowledge. PrimeKG integrates 20 high-quality resources to describe 17,080 diseases with 4,050,249 relationships, including over 100,000 nodes and 29 types of edges that densely connect disease nodes with drugs, genes, exposures, and phenotypes. We preprocess the graph data by mapping names to

<sup>2</sup><https://neo4j.com/>

<sup>3</sup><https://opencypher.org/>

<sup>4</sup><https://www.langchain.com/>

<sup>5</sup><https://github.com/chrschy/fact-finder>

their preferred terms, as described in Section 3.1, and converting all entries to lowercase.

**Text-to-Cypher dataset.** We manually generated a ground-truth dataset containing 69 text-to-Cypher query pairs specifically designed for medical questions. These queries are complex, often involving multiple hops in the graph, aggregation, and boolean question structures. They require deep knowledge of Cypher and the KG. Each entry includes a natural language question, the corresponding Cypher query, the expected answer, and relevant nodes and relationships. This dataset provides a benchmark to evaluate the ability of text-to-Cypher systems to interpret and execute complex queries. While specialized to the PrimeKG graph, the dataset leverages PrimeKG’s extensive applicability, making it a valuable resource for various medical information retrieval tasks. Examples from the dataset include simpler questions like *Which drugs have pterygium as a side effect?* and more complex ones such as *Which medications have more off-label uses than approved indications?* and *Which diseases have only treatments that have no side effects at all?*

### 3 System Description

#### 3.1 Cypher Query Generation

Generating code to query structured databases from natural language inquiries used to be a complex process, involving steps such as entity and relation extraction, entity and relation linking, query type classification, template-based or compositional query generation (Srivastava et al., 2021; Chakraborty et al., 2021). Retrieval-based approaches like Graph-RAG aim to simplify this by partitioning the graph into communities of nodes and edges, which are then retrieved and summarized using a LLM to generate an answer (Edge et al., 2024). However, Graph-RAG struggles with complex queries that span multiple graph communities and require exact graph operations.

With the advent of LLMs, however, QA systems can now understand domain-specific questions and generate valid queries directly, allowing for more flexible approaches. Much of the research has been centered on text-to-SQL generation, where LLMs have demonstrated considerable effectiveness (Gao et al., 2023; Chang and Fosler-Lussier, 2023), including in the medical domain (Ziletti and D’Ambrosi, 2024), and have often performed better than specialized models (Pourreza and Rafiei,

2023). Conversely, the area of text-to-Cypher query generation remains relatively under-explored, with prior research primarily focused on sequence-to-sequence models (Zhao et al., 2024; Guo et al., 2022). Only recently has the application of LLMs to this task begun to emerge (Feng et al., 2023). To bridge this gap, our work evaluates the capabilities of LLMs to produce robust Cypher queries for scientific QA in the medical domain (see Sec. 4.1).

We prompt LLMs with questions and graph schemas, including node and relationship types and their properties, to generate Cypher queries. When a graph relation is not self-explanatory, we add its natural language description to the prompt. For instance, for the relation *ppi*, we add *"Temporary, non-covalent binding between protein molecules. Protein-protein interactions occur..."*. During instruction prompting, we also identify questions that cannot be answered by the given graph schema. In such cases, the LLM returns the string `SCHEMA_ERROR` along with a brief explanation of why it could not generate an answer. FactFinder detects this marker using a regex and returns the explanation to the user.

We include an entity extraction model to align entity names in questions with those in the KG, based on Linnaeus (Gerner et al., 2010) and developed set of vocabularies for entity types. This step ensures consistency by replacing detected entities with their preferred KG terms (e.g., *alcohol* to *ethanol*) and generating sentences linking each entity to its category (see Fig. 1 left), reducing the LLM’s reliance on domain-specific knowledge.

#### 3.2 Query Pre-Processors

Before querying the graph with the generated Cypher query, we preprocess it to increase the system’s robustness. This leverages the structured natural language understanding provided by the translation of questions to Cypher queries. Various regular expression-based methods target specific query elements.

**Formatting.** First, we format the query to improve readability and consistency. This includes adding indentation, line breaks, and ensuring consistent naming conventions, which simplifies the application of regular expressions in subsequent steps.

**Lowercasing Property Values.** We convert property values in the Cypher query to lowercase, matching the pre-lowered graph properties.

**Synonym Selection.** We map entities in the Cypher query to preferred terms used in the graph. If no

mapping is available, we use external tools (e.g., `skos:altLabel` queries against Wikidata) to find synonyms and match them to graph terms.

**Deprecated Code Handling.** We correct deprecated code generated by the LLM, such as replacing the obsolete `SIZE()` keyword with the current `COUNT()` keyword.

**Child to Parent Node Mapping.** In PrimeKG (Chandak et al., 2023), some node types are connected by parent-child relationships. We replace child nodes with their parent nodes in the Cypher query to ensure completeness.

### 3.3 Graph Question Answering and Verbalization

The pre-processed Cypher query is executed on the graph, returning a unique set of nodes that may include names, properties, IDs, and other elements. Next, the question and graph results are incorporated into a prompt template and sent to an LLM. The prompt instructs the LLM to rely solely on the graph information to formulate an answer, which is then provided as the final natural language output.

### 3.4 Explainability through Evidence

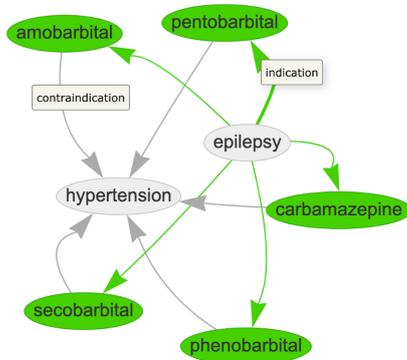


Figure 2: Example of the evidence subgraph for *Which drugs against epilepsy should not be used by patients with hypertension?*

To ensure transparency and explainability, the system provides various forms of evidence alongside the natural language answer. These include intermediate results, such as Cypher generation prompts and graph results, as well as explicitly created information like the underlying subgraph for a query. These evidences are displayed in the frontend to track the system’s behavior.

**Cypher Query Evidence.** The primary evidence is the Cypher query generated by the model, demonstrating how the given question maps to the graph

structure. This enables expert users to evaluate the model’s understanding of the question and the quality of the generated query.

**Graph Response.** The system also provides the actual response from the graph, consisting of the nodes and relationships returned by the executed query.

**Subgraph Visualization.** To enhance interpretability, we provide a subgraph as part of the evidence. This subgraph visually displays (via Pyvis<sup>6</sup>) the relevant nodes and edges, illustrating the subset of the main graph that contributed to the specific answer, as shown in Fig. 2.

**Sub-Graph Generator.** A notable challenge is that the original Cypher query usually returns only the nodes required to answer the question, omitting the edges that connect the question’s entity to the answer nodes. To address this, we generate a new Cypher query that fetches both the answer nodes and the connecting edges. This is achieved by submitting the original Cypher query to a LLM and instructing it to return all nodes and edges present in the query.

For example, if our original Cypher query is:

```

MATCH (g:gene_or_protein name:"pink1")-
[:associated_with]->(d:disease)
RETURN d.id AS ID, d.name AS Name

```

our subgraph Cypher query is:

```

MATCH (g:gene_or_protein name:"pink1")-
[a:associated_with]->(d:disease)
RETURN g, d, a

```

Note that this last subgraph Cypher query returns all relevant nodes and edges.

### 3.5 User Interface and Example of Usage

Our target audience includes researchers in the life sciences, such as those working in medical research and crop science, who are interested in exploring connections between drugs, genes, proteins, and other biological entities to discover new research directions. To make our pipelines accessible, we developed a graphical user interface using Streamlit<sup>7</sup>. This interface allows users to input questions and view generated answers. Users can select different pipelines, such as LLM-only or those incorporating KGs or documents. For instance, a medical researcher might ask, *Which drugs against epilepsy*

<sup>6</sup><https://pyvis.readthedocs.io/>

<sup>7</sup><https://streamlit.io/>

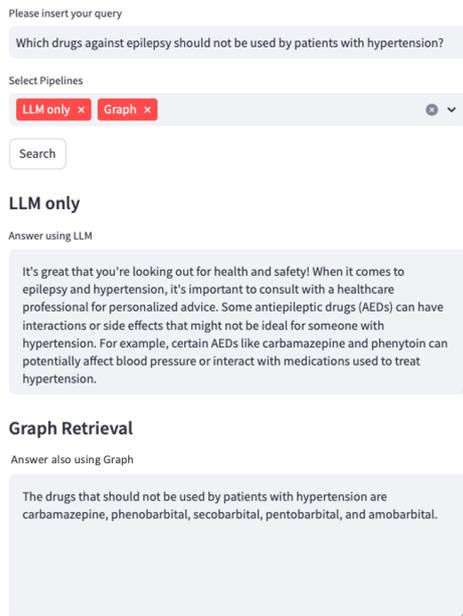


Figure 3: User Interface with question and answers of the standalone LLM and our graph-based hybrid system.

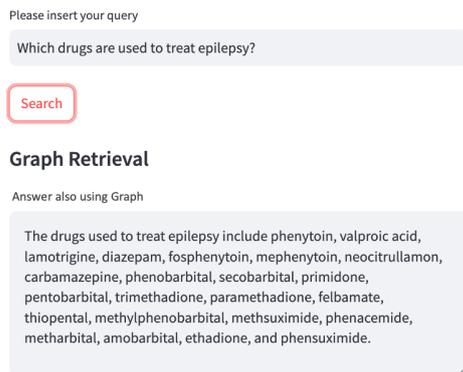


Figure 4: Answer for exploring drugs used to treat epilepsy.

*should not be used by patients with hypertension?* (Fig. 3). To delve deeper, they could follow up with *Which drugs are used to treat epilepsy?* (Fig. 4). Going further, with the question *Which genes are targeted by amobarbital but not lamotrigine?*, they could take a first step towards understanding the genetic interactions that differentiate drug effects, resulting in the identification of four genes exclusively targeted by amobarbital (Fig. 5). Additionally, users can interactively visualize the relevant subgraph, generated Cypher query, and graph response, enabling them to verify response accuracy and understand the underlying data (Fig. 6). These features (see Section 3.4) enhance transparency and foster trust in the system.

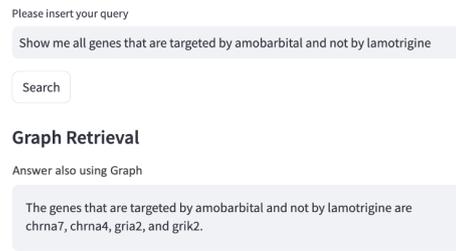


Figure 5: Answer for exploring genes targeted by amobarbital but not lamotrigine.

## 4 Evaluation

### 4.1 Graph Retrieval Evaluation

To quantify the graph retrieval step, we assess the returned nodes from graph queries by comparing result nodes from executing the ground truth queries with those from the generated queries. This enables a quantitative evaluation of the text-to-Cypher step. We use the ground truth text-to-Cypher dataset described in Sec. 2 for this evaluation. We compute intersection over union (IoU), precision, and recall for the expected and generated graph result sets, as shown in Table 1.

The results demonstrate robust performance across models, with the highest accuracy exceeding 75%. GPT-4o demonstrates superior performance across all evaluation metrics. Contrary to expectations, Entity Enhancement (EE, Fig.1 left and Sec.3.1) decreased performance in all models except GPT-4-Turbo. This decline stems from two main factors: (1) PrimeKG’s merged gene-protein nodes create ambiguity when EE is applied, leading to incorrect relations, and (2) a terminology mismatch where the synonym endpoint returns formal medical terms while PrimeKG uses non-medical terminology. The results suggest that leveraging LLMs’ internal knowledge is more effective for Cypher query generation than external entity enhancement. The following more detailed analysis focuses only on the most strongly performing models, GPT-4o and GPT-4-turbo.

### 4.2 Evaluating Correctness and Completeness

To assess the quality of LLM-generated answers, we conduct two evaluations using the LLM-as-a-Judge approach (Zheng et al., 2024): (1) comparing the answers from our KG-LLM-based system to those from an LLM-only system, and (2) evaluating the reliability of LLM verbalization of information provided by the KG. In both cases, correctness



Figure 6: User Interface with the generated Cypher query and the graph response as evidence.

Model	EE	IoU	Precision	Recall
gpt-4-turbo	True	71.3	73.6	71.6
gpt-4-turbo	False	62.7	65.3	64.9
gpt-4o	True	74.9	77.0	77.6
gpt-4o	False	<b>75.2</b>	<b>77.5</b>	<b>77.8</b>
Llama-3.1-70B	True	46.6	46.9	53.3
Llama-3.1-70B	False	52.8	55.4	59.7
Mixtral-8x7B	True	31.4	32.3	33.6
Mixtral-8x7B	False	37.6	39.5	40.1

Table 1: Results for the graph retrieval evaluation (metrics in %). IoU stands for intersection over union and EE for Entity Enhancement, see Section 4.1.

is defined as the inclusion of only facts from the graph nodes, and completeness as the inclusion of all such facts.

**Hybrid system vs. LLM-only.** We compare the hybrid KG-based system against a standalone LLM. The hybrid system (GPT-4o without entity enhancement, Sec. 4.1) is evaluated to produce more correct (complete) answers in 94.12% (96.08%) of cases, demonstrating its superior performance in providing accurate and complete responses.

**LLM verbalization.** We evaluate the verbalization of natural language answers from graph results. In this evaluation, 89.13% of answers are deemed correct, and 80.43% complete, indicating high accuracy in verbalization.

### 4.3 Handling Incorrect Graph Responses

Finally, we evaluate FactFinder’s ability to handle incorrect or incomplete information in graph responses. The system should be able to refuse to answer if the Cypher query generation step produces a wrong query, thus retrieving the correct data from the KG. We test this by disabling Cypher query generation and supplying incorrect Cypher queries for each question, resulting in incorrect graph results.

The results in Table 2 show that both GPT-4-turbo and GPT-4-o can detect irrelevant informa-

(in %)	gpt-4o	gpt-4-turbo
Answer Denied	65/69 (94.2)	63/69 (91.3)
Uncertain Answer	1/69 (1.5)	1/69 (1.5)
Full Answer	3/69 (4.3)	5/69 (7.3)

Table 2: Handling irrelevant information in graph responses.

tion and correctly respond with "I don’t know" in over 90% of cases, demonstrating that the LLMs can reason and understand when the knowledge passed to them is not relevant. This highlights FactFinder’s ability to enhance reliability by leveraging both structured and world knowledge.

Manual analysis revealed that in one case, the LLM expressed uncertainty with the phrase "The provided information does not mention ..." indicating potential inaccuracies. In a few cases, the models provided full answers despite unrelated KG results, especially for count, boolean, and lengthy responses, highlighting detection challenges in these scenarios.

## 5 Conclusion

This work demonstrates the value of integrating structured, factual knowledge into a user-friendly chat system, providing researchers with a reliable tool for answering scientific questions while minimizing hallucinations. We show that LLMs can generate valid Cypher queries to retrieve relevant data from a KG, informing accurate answers. By grounding every answer in the underlying knowledge graph, FactFinder provides an interpretable and reproducible response pathway that pure LLM systems cannot match. The creation of such a robust system is crucial for enhancing research capabilities. Future work will focus on expanding the evaluation dataset, quantifying system uncertainty, and enabling access to multiple KGs, possibly through agent-based retrieval.

## References

- Carlos Badenes-Olmedo and Oscar Corcho. 2023. Lessons learned to enable question answering on knowledge graphs extracted from scientific publications: A case study on the coronavirus literature. *Journal of Biomedical Informatics*, 142:104382.
- Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 78–106, Toronto, Canada. Association for Computational Linguistics.
- Bojana Bašaragin, Adela Ljajić, Darija Medvecki, Lorenzo Cassano, Miloš Košprdić, and Nikola Milošević. 2024. How do you know that? teaching generative language models to reference answers to biomedical questions. *BioNLP Workshop 2024*.
- 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*.
- Nilesh Chakraborty, Denis Lukovnikov, Gaurav Maheshwari, Priyansh Trivedi, Jens Lehmann, and Asja Fischer. 2021. Introduction to neural network-based question answering over knowledge graphs. *WIREs Data Mining and Knowledge Discovery*, 11(3):e1389.
- Payal Chandak, Kexin Huang, and Marinka Zitnik. 2023. Building a knowledge graph to enable precision medicine. *Scientific Data*, 10(1):67.
- Shuaichen Chang and Eric Fosler-Lussier. 2023. How to prompt llms for text-to-sql: A study in zero-shot, single-domain, and cross-domain settings. *Preprint*, arXiv:2305.11853.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*.
- Guandong Feng, Guoliang Zhu, Shengze Shi, Yue Sun, Zhongyi Fan, Sulin Gao, and Jun Hu. 2023. Robust nl-to-cypher translation for kbqa: Harnessing large language model with chain of prompts. In *Knowledge Graph and Semantic Computing: Knowledge Graph Empowers Artificial General Intelligence*, pages 317–326, Singapore. Springer Nature Singapore.
- Yichun Feng, Lu Zhou, Yikai Zheng, Ruikun He, Chao Ma, and Yixue Li. 2024. Knowledge graph-based thought: a knowledge graph enhanced llms framework for pan-cancer question answering. *bioRxiv*.
- Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2023. Text-to-sql empowered by large language models: A benchmark evaluation. *Preprint*, arXiv:2308.15363.
- Martin Gerner, G. Nenadic, and Casey M. Bergman. 2010. Linnaeus: A species name identification system for biomedical literature. *BMC Bioinformatics*, 11:85 – 85.
- Aibo Guo, Xinyi Li, Guanchen Xiao, Zhen Tan, and Xiang Zhao. 2022. Spcql: A semantic parsing dataset for converting natural language into cypher. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM '22*, page 3973–3977, New York, NY, USA. Association for Computing Machinery.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12).
- 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. 2023a. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin Zhao, and Ji-Rong Wen. 2023b. StructGPT: A general framework for large language model to reason over structured data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9237–9251, Singapore. Association for Computational Linguistics.
- Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. 2023. Gpt-4 passes the bar exam. *Available at SSRN 4389233*.
- Adela Ljajić, Miloš Košprdić, Bojana Bašaragin, Darija Medvecki, Lorenzo Cassano, and Nikola Milošević. 2024. Scientific qa system with verifiable answers. *The 6th International Open Search Symposium*.
- Chaitanya Malaviya, Subin Lee, Sihao Chen, Elizabeth Sieber, Mark Yatskar, and Dan Roth. 2023. Expertqa: Expert-curated questions and attributed answers. *arXiv preprint arXiv:2309.07852*.
- Nikola Milošević and Wolfgang Thielemann. 2023. Comparison of biomedical relationship extraction methods and models for knowledge graph creation. *Journal of Web Semantics*, 75:100756.
- Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

## A Appendix

- Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*.
- 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.
- Mohammadreza Pourreza and Davood Rafiei. 2023. [DIN-SQL: Decomposed in-context learning of text-to-SQL with self-correction](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Priyanka Sen, Sandeep Mavadia, and Amir Saffari. 2023. [Knowledge graph-augmented language models for complex question answering](#). In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 1–8, Toronto, Canada. Association for Computational Linguistics.
- Saurabh Srivastava, Mayur Patidar, Sudip Chowdhury, Puneet Agarwal, Indrajit Bhattacharya, and Gautam Shroff. 2021. [Complex question answering on knowledge graphs using machine translation and multi-task learning](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3428–3439, 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, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Ethan Waisberg, Joshua Ong, Mouayad Masalkhi, Sharif Amit Kamran, Nasif Zaman, Prithul Sarker, Andrew G Lee, and Alireza Tavakkoli. 2023. Gpt-4: a new era of artificial intelligence in medicine. *Irish Journal of Medical Science (1971-)*, 192(6):3197–3200.
- Ziyu Zhao, Wei Liu, Tim French, and Michael Stewart. 2024. Cyspider: A neural semantic parsing corpus with baseline models for property graphs. In *AI 2023: Advances in Artificial Intelligence*, pages 120–132, Singapore. Springer Nature Singapore.
- 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.
- Angelo Ziletti and Leonardo D’Ambrosi. 2024. [Retrieval augmented text-to-sql generation for epidemiological question answering using electronic health records](#). *Preprint*, arXiv:2403.09226.



# Fact Finder

## Factual Question-Answering

Select Example

- Insert Query
- What are the phenotypes associated with cardiofacial dysplasia?
- What are the genes responsible for psoriasis?
- Which diseases involve PINK1?
- How many drugs against epilepsy are available?
- Which medications have more off-label uses than approved indications?

Please insert your query

Which drugs are used to treat ocular hypertension?

Select Pipelines

LLM only x Graph x

Search

### LLM only

Answer using LLM

Ocular hypertension is a condition where the pressure inside the eye (intraocular pressure) is higher than normal, which can lead to glaucoma if not managed properly. There are several types of medications that are commonly used to treat ocular hypertension:

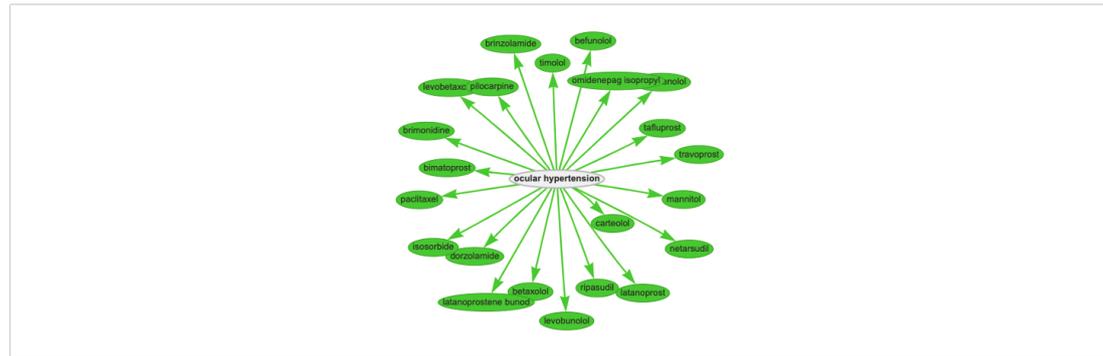
- Prostaglandin Analogues**: These help to increase the outflow of fluid from the eye. Examples include latanoprost (Xalatan), bimatoprost (Lumigan), and travoprost (Travatan).
- Beta Blockers**: These reduce the production of fluid within the eye. Examples include timolol (Timoptic) and betaxolol (Betoptol).

### Graph Retrieval

Answer using Graph

The drugs used to treat ocular hypertension include tafiprost, paclitaxel, bimatoprost, betaxolol, timolol, carteolol, levobunolol, metipranolol, befunolol, levobetaxolol, pilocarpine, brinzolamide, dorzolamide, brimonidine, isosorbide, netarsudil, travoprost, latanoprost, latanoprostene bunod, ripasudil, mannitol, and omdenepag isopropyl.

Relevant Subgraph:



Show Evidence

Cypher Query:

```
MATCH (d:disease {name: "ocular hypertension"})-[:indication]->(drug:drug)
RETURN DISTINCT drug.name AS drug_name
```

Cypher Response:

```
[{'drug_name': 'tafiprost'}, {'drug_name': 'paclitaxel'}, {'drug_name': 'bimatoprost'}, {'drug_name': 'betaxolol'}, {'drug_name': 'timolol'}, {'drug_name': 'carteolol'}, {'drug_name': 'le
```

Cypher Prompt

Task: Generate Cypher statement to query a graph database described in the following schema.  
Instructions:  
Use only the provided relationship types and properties in the schema.  
Do not use any other relationship types or properties that are not provided.  
If there is no sensible Cypher statement for the given question and schema, state so and prepend SCHEMA\_ERROR to your answer.  
Any variables that are returned by the query must have readable names.  
Remove modifying adjectives from the entities queried to the graph.

Answer Prompt

You are an assistant that helps to form nice and human understandable answers.  
The information part contains the provided information that you must use to construct an answer.  
The provided information is authoritative, you must never doubt it or try to use your internal knowledge to correct it.  
Make the answer sound as a response to the question. Do not mention that you based the result on the given information.  
If the provided information is a list, include all entries in your response.  
If the provided information is empty, say that you don't know the answer.  
Provided information:

Figure 7: User interface of Fact Finder for the question *Which drugs are used to treat ocular hypertension?*. The answers of the standalone LLM and our graph-based hybrid system are compared as output. In addition, the relevant subgraph is displayed as evidence together with the generated Cypher query and the answer from the graph.



## Fact Finder

### Factual Question-Answering

Select Example

- Insert Query  What are the phenotypes associated with cardiofacial dysplasia?  What are the genes responsible for psoriasis?  Which diseases involve PINK1?  How many drugs against epilepsy are available?  
 Which medications have more off-label uses than approved indications?

Please insert your query

Which drugs against epilepsy should not be used by patients with hypertension?

Select Pipelines

LLM only x Graph x

Search

### LLM only

Answer using LLM

It's great that you're looking out for health and safety! When it comes to epilepsy and hypertension, it's important to consult with a healthcare professional for personalized advice. Some antiepileptic drugs (AEDs) can have interactions or side effects that might not be ideal for someone with hypertension. For example, certain AEDs like carbamazepine and phenytoin can potentially affect blood pressure or interact with medications used to treat hypertension.

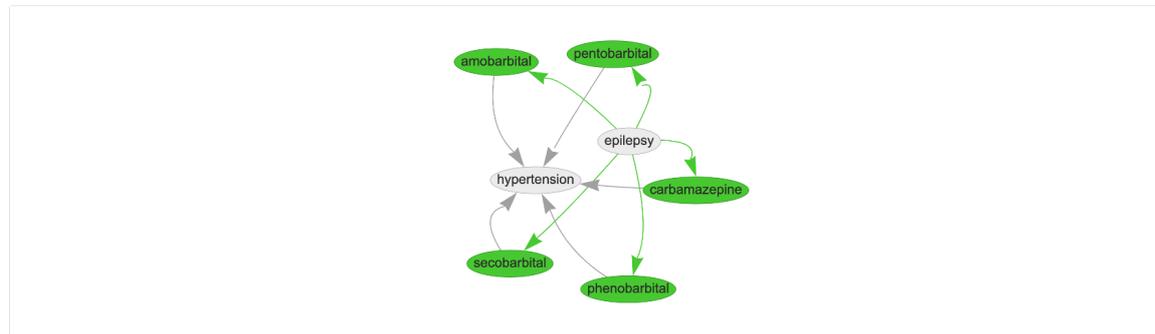
However, the best course of action is always to discuss your specific situation with your doctor or a pharmacist. They can provide guidance tailored to your individual health needs. Is there anything else you're curious about or any other topic you'd like to chat about?

### Graph Retrieval

Answer using Graph

The drugs that should not be used by patients with hypertension are carbamazepine, phenobarbital, secobarbital, pentobarbital, and amobarbital.

Relevant Subgraph:



Show Evidence

Cypher Query:

```
MATCH (epilepsy:disease {name: "epilepsy"})-[:indication]->(drug:drug)-[:contraindication]->(hypertension:disease {name: "hypertension"})
RETURN DISTINCT drug_name AS drug_name
```

Cypher Response:

```
[{'drug_name': 'carbamazepine'}, {'drug_name': 'phenobarbital'}, {'drug_name': 'secobarbital'}, {'drug_name': 'pentobarbital'}, {'drug_name': 'amobarbital'}]
```

Cypher Prompt

Task: Generate Cypher statement to query a graph database described in the following schema.  
Instructions:  
Use only the provided relationship types and properties in the schema.  
Do not use any other relationship types or properties that are not provided.  
If there is no sensible Cypher statement for the given question and schema, state so and prepend SCHEMA\_ERROR to your answer.  
Any variables that are returned by the query must have readable names.  
Remove modifying adjectives from the entities queried to the graph.

Answer Prompt

You are an assistant that helps to form nice and human understandable answers.  
The information part contains the provided information that you must use to construct an answer.  
The provided information is authoritative, you must never doubt it or try to use your internal knowledge to correct it.  
Make the answer sound as a response to the question. Do not mention that you based the result on the given information.  
If the provided information is a list, include all entries in your response.  
If the provided information is empty, say that you don't know the answer.  
Provided Information:

Figure 8: User interface of Fact Finder for the question *Which drugs against epilepsy should not be used by patients with hypertension?*.

# Simplifying Outcomes of Language Model Component Analyses with ELIA

Aaron Louis Eidt<sup>1,2</sup> Nils Feldhus<sup>1,3</sup>

<sup>1</sup>Technische Universität Berlin

<sup>2</sup>Fraunhofer Heinrich Hertz Institute

<sup>3</sup>BIFOLD – Berlin Institute for the Foundations of Learning and Data

aaron.eidt@hhi.fraunhofer.de, feldhus@tu-berlin.de

## Abstract

While mechanistic interpretability has developed powerful tools to analyze the internal workings of Large Language Models (LLMs), their complexity has created an accessibility gap, limiting their use to specialists. We address this challenge by designing, building, and evaluating ELIA (Explainable Language Interpretability Analysis), an interactive web application that simplifies the outcomes of various language model component analyses for a broader audience. The system integrates three key techniques – Attribution Analysis, Function Vector Analysis, and Circuit Tracing – and introduces a novel methodology: using a vision-language model to automatically generate natural language explanations (NLEs) for the complex visualizations produced by these methods. The effectiveness of this approach was empirically validated through a mixed-methods user study, which revealed a clear preference for interactive, explorable interfaces over simpler, static visualizations. A key finding was that the AI-powered explanations helped bridge the knowledge gap for non-experts; a statistical analysis showed no significant correlation between a user’s prior LLM experience and their comprehension scores, suggesting that the system reduced barriers to comprehension across experience levels. We conclude that an AI system can indeed simplify complex model analyses, but its true power is unlocked when paired with thoughtful, user-centered design that prioritizes interactivity, specificity, and narrative guidance.

## 1 Introduction

The growing capabilities of LLMs are coupled with a proportional increase in their inscrutability. While the field of mechanistic interpretability has made major strides in developing tools to reverse-engineer the internal algorithms of these black-box systems (Bereska and Gavves, 2024; Ferrando et al., 2024), a new challenge has emerged:

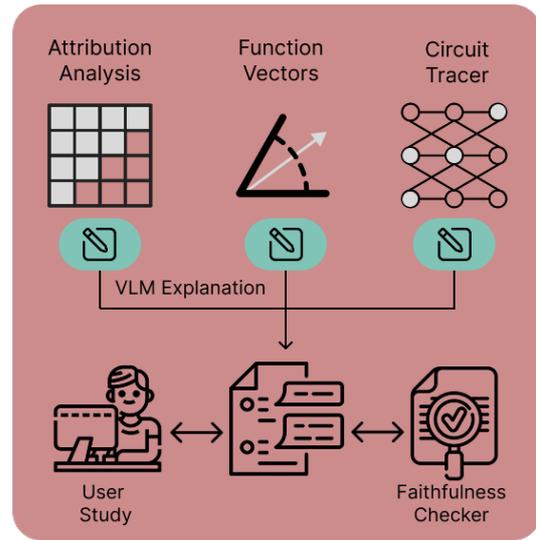


Figure 1: ELIA system overview, including three analysis methods (Attribution Analysis, Function Vectors, and Circuit Tracer) and the explanation generation workflow using VLMs to transform complex interpretability analyses into accessible NLEs. The system is evaluated using a Faithfulness Checker and a user study.

the outputs of these analyses are often as complex as the models they seek to explain. Techniques such as attribution analysis, which traces predictions to input tokens (Sarti et al., 2023), or circuit tracing (Lindsey et al., 2025), which maps specific computational pathways, produce visualizations and data that require specialized expertise to decipher. This creates an accessibility gap, limiting the vital conversation around AI safety and reliability (Weidinger et al., 2023) to a small circle of specialists and excluding developers, domain experts, and policymakers who could benefit most.

To bridge this gap, we introduce ELIA (Explainable Language Interpretability Analysis), an interactive web application<sup>1</sup> designed to make the outcomes of complex model analyses accessible to a broader audience. ELIA integrates three

<sup>1</sup>Demo: <https://hf.co/spaces/aaron0eidt/ELIA>  
GitHub: <https://github.com/aaron0eidt/ELIA>

powerful interpretability techniques, Attribution Analysis (Sarti et al., 2023), Function Vector Analysis (Todd et al., 2024), and Circuit Tracing (Lindsey et al., 2025), within a user-centered interface. A vision-language model then generates NLEs for the intricate visualizations produced by these analyses.

Through a mixed-methods user study, we demonstrate the effectiveness of this approach. Our findings show that the AI-generated explanations helped reduce the knowledge gap, enabling non-experts to comprehend complex model behaviors at levels approaching those of users with prior LLM experience. Furthermore, the study revealed a strong user preference for interactive, explorable interfaces over static visualizations. This work provides empirical evidence that the strategic combination of AI-powered explanation and thoughtful, interactive design can significantly lower the barrier to understanding the internal workings of LLMs.

## 2 Background and Related Work

The field of **NLP interpretability** has progressed through three interconnected streams: moving from correlational to causal analysis, shifting focus from input-output attribution to internal component analysis, and developing methods to communicate these complex findings to a broader audience (Saphra and Wiegrefe, 2024; Calderon and Reichart, 2025).

Early interpretability work adapted **attribution techniques** from computer vision, such as Integrated Gradients (Sundararajan et al., 2017), to create saliency heatmaps that identify influential input tokens. However, the “attention is not explanation” debate and critical sanity checks (Jain and Wallace, 2019; Wiegrefe and Pinter, 2019; Adebayo et al., 2018) revealed the limitations of these correlational methods, pushing the field toward more rigorous, intervention-based approaches.

Techniques like activation patching and causal tracing now allow researchers to establish causal links between specific model components and their behavior by **intervening in the computational process** (Zhang and Nanda, 2024). Landmark findings include the identification of induction heads that perform in-context learning (Nanda et al., 2022) and the discovery that entire tasks can be represented by abstract Function Vectors within the model’s activation space (Todd et al., 2024). These vectors can be extracted and even composed, demonstrating that models learn structured, portable representations of functions.

Despite these powerful analytical tools, **communicating the findings remains a significant bottleneck**. The raw outputs, complex graphs, heatmaps, and high-dimensional plots, are often inscrutable to non-experts (Colin et al., 2022; Schuff et al., 2022). To address this, interactive visualization tools like LIT, BertViz, Inseq, and LM Transparency Tool provide explorable interfaces for experts (Tenney et al., 2020; Vig, 2019; Sarti et al., 2023; Tufanov et al., 2024). More recently, the focus has shifted to **automated explanation** systems that use explainer models to generate natural language descriptions for neuron activity or attention patterns (Bills et al., 2023; Feldhus and Kopf, 2025), agents using vision-language models for end-to-end interpretability experiment design (Shaham et al., 2024; Kim et al., 2025), and discovering circuits that represent a particular higher-level function of a model (Wang et al., 2023; Hanna et al., 2025). However, this automation introduces a critical trade-off between the faithfulness of an explanation (how accurately it reflects the model’s process) and its simplicity (Feldhus et al., 2023; Parcalabescu and Frank, 2024). Our work is situated at this frontier, aiming to bridge the gap between complex, faithful analyses and simple, accessible explanations through a combination of interactive visualization and AI-generated narrative.

## 3 ELIA

### 3.1 System Architecture

ELIA is an interactive web application designed to make the internal mechanisms of LLMs more transparent and understandable. The system is built using Streamlit<sup>2</sup>, a Python-based framework chosen for its ability to rapidly create data-centric, interactive user interfaces. The backend leverages the scientific Python ecosystem, with PyTorch and the Transformers library for model handling, Plotly<sup>3</sup> for dynamic visualizations, and the *inseq* toolkit<sup>4</sup> for attribution analyses (Sarti et al., 2023).

The architecture is centered around two core models: a **subject model**, the 7-billion parameter OLMo-2, whose behavior is being analyzed (Groeneveld et al., 2024); and a vision-enabled **explanation model** (Qwen2.5-VL-72B) tasked with simplifying the analytical outputs. When a user interacts with one of ELIA’s three analysis pages

<sup>2</sup><https://streamlit.io>

<sup>3</sup><https://plotly.com>

<sup>4</sup><https://inseq.org>

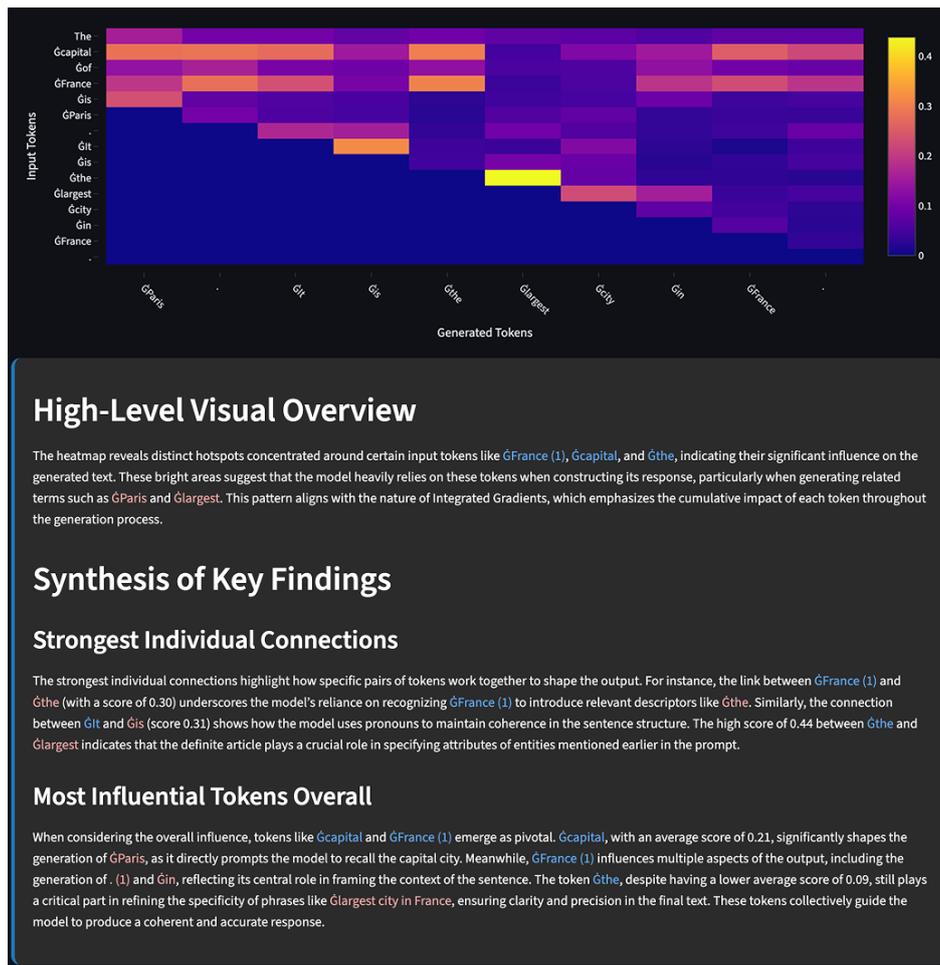


Figure 2: The interactive Attribution Heatmap using Integrated Gradients with an AI-generated natural language explanation. The heatmap visualizes the influence of input tokens on the generated output, and the explanation interprets these results in an accessible narrative.

(Attribution Analysis, Function Vector Analysis, or Circuit Tracing), the subject model’s internal activations and outputs are visualized (Figure 1). These visualizations, along with structured textual data, are passed to the explanation model, following prior work on verbalizing explanations (Feldhus et al., 2023). This generates a structured, natural language summary of the key insights in an accessible narrative (Figure 2).

To ensure consistency, API calls to the explanation model are made with a low temperature and a fixed seed, making the generated text largely deterministic. The entire application is internationalized, with full support for both English and German to broaden its accessibility.

### 3.2 Faithfulness Verification

A key component of ELIA’s architecture is an automated faithfulness verification system, designed to ensure the reliability of the AI-generated explanations. This system leverages the same explanation

model in a multi-step process. First, after generating the initial narrative, the explanation model is prompted again, this time to act as a claim extraction agent, parsing its own text to identify all verifiable, factual statements and structure them as a JSON list, following a similar approach to atomic fact extraction in FActScore (Min et al., 2023). These claims range from specific quantitative statements (e.g., “Layer 12 had the highest activation.”) to more abstract semantic assertions (e.g., “Early layers handle syntax.”). In the second stage, a verification module programmatically checks each claim against the ground-truth data from the underlying analysis. For quantitative claims, this is a direct data comparison. For more abstract semantic claims, the explanation model is called a third time, now tasked to act as a fact-checker to assess the logical plausibility of the claim against the data. The outcome, a *verified* or *contradicted* status for each claim and the supporting evidence, is then presented to the user.

To mitigate the circularity risk of using Qwen2.5-VL-72B for both generation and verification, the verification module operates deterministically (temperature 0.0, fixed seed) and relies heavily on programmatic grounding rather than purely LLM-based judgments. When LLM-based semantic verification is required, the explainer model is constrained by hard-coded rules, negative constraints, and exact synonym-mapping directives, effectively preventing the model from self-affirming its own hallucinations.

### 3.3 Attribution Analysis

The Attribution Analysis page provides a granular view of the model’s decision-making by quantifying the influence of individual input tokens on the generated output. It integrates two key features: core attribution methods and an influence tracer.

The primary analysis is grounded in three established **feature attribution** techniques, Saliency, Integrated Gradients, and Occlusion, which are implemented using the *Inseq* toolkit (Sarti et al., 2023). After the subject model generates text from a user’s prompt, the chosen method computes an attribution matrix that is visualized as an interactive heatmap. To translate this complex data into an accessible narrative, the explanation model is given a multi-modal prompt. This prompt combines the heatmap image with a rule-based textual summary that highlights key data points, such as the most influential input tokens and the most affected output tokens, guiding the model to generate a comprehensive, structured explanation of the token-level interactions (Figure 2). The faithfulness of these explanations is also verified (Figure 6 in Appendix).

To further contextualize the model’s output, the page includes an **Influence Tracing** feature that identifies similar documents from the model’s training data by performing a  $k$ -nearest neighbors search against a pre-computed Faiss index of the Dolma dataset, the training corpus for OLMo (Douze et al., 2024). The user’s prompt is embedded into the same vector space as the training documents, and the most similar examples are retrieved. This allows users to perform a form of data archaeology, exploring potential sources that may have influenced the model’s response (Figure 7).

### 3.4 Function Vector Analysis

The Function Vector Analysis page offers a high-level, semantic view of the model’s behavior. It allows users to explore how the model represents dif-

ferent instructions by comparing a user’s prompt to a pre-computed, high-dimensional space of Function Vectors (Todd et al., 2024).

The core of this analysis is a custom dataset of instructional prompts (see Appendix B for details), organized into a hierarchy of broad “function types” (e.g., abstractive tasks) and specific “function categories” (e.g., summarization). The function vector for each category is pre-computed by averaging the final-layer, final-token activations of all its example prompts. When a user enters a new prompt, its own activation vector is computed and compared against this space using cosine similarity.

The results are presented through a suite of interactive visualizations (Figure 3, left). A 3D scatter plot, generated using Principal Component Analysis (PCA), shows the geometric relationship between the user’s prompt and the function vector clusters, providing an intuitive map of the model’s functional space. This is complemented by a bar chart of the top-scoring function types and a hierarchical sunburst chart that visualizes the similarity scores for all categories. For each of these visualizations, a targeted, AI-powered explanation is generated, synthesizing the key quantitative findings into an accessible, natural-language summary. The faithfulness of the PCA explanation is also verified (see Figure 8 in the Appendix).

### 3.5 Circuit Trace Analysis

This page offers the most granular view of the model’s internal workings, building on the circuit tracing framework by (Lindsey et al., 2025). The analysis is centered around a small autoencoder, called a Cross-Layer Transcoder (CLT) (Dunefsky et al., 2024), which is pre-trained to learn a simplified, sparse representation of the OLMo model’s internal activations. This CLT is trained on a diverse corpus from the Dolma dataset using  $L_1$  sparsity regularization, gradient clipping, and cosine annealing learning rate scheduling (Figure 9 for training dynamics). The CLT is trained to reconstruct the main model’s signals while being penalized for using too many features, forcing it to identify the most functionally significant patterns.

These learned features are then given semantic meaning through an automated interpretation step. For each feature, the top-activating input tokens are passed to the explanation model, which generates a concise functional label (e.g., “identifying JSON syntax”). These interpretable features become the nodes in the main visualization: a layer-



Figure 3: Function Vector and Circuit Trace Analysis visualizations. The 3D PCA plot (left) places the user’s prompt in a semantic functional space, while the Circuit Graph (right) traces the flow of information through interpretable features across layers. Both are accompanied by AI-generated explanations.

by-layer Circuit Graph (Figure 3, right). This graph shows the flow of information from the input tokens, through the activated features, to the final output, with node size and color indicating activation strength and edge thickness representing influence. Additionally, the system provides a view of local path ablations (Figure 10 in the Appendix), which demonstrates what happens when specific paths in the top feature graph are ablated. To make this complex graph accessible, a multi-modal prompt containing the graph image and a summary of key feature activity is used to generate a structured, AI-powered narrative of the information flow. The page also includes interactive "Subnetwork" and "Feature" explorers (see Figures 12 and 13 in the Appendix), allowing users to drill down into the

behavior of individual features and their local computational pathways, each augmented with its own targeted, AI-generated explanation. The faithfulness of these explanations is verified and displayed to users (Figure 11 in the Appendix).

## 4 Faithfulness Analysis

Maintaining fidelity between the automatically generated narratives and the underlying model behavior requires dedicated instrumentation on every analysis page. We implement specific verification pipelines that inspect the data powering each visualization, run deterministic checks, and produce aggregate diagnostics that we summarize in Table 1 and describe in the following for every explanation.

Component	Feature	Verified	Faith. (%)
Attribution	Saliency	59/68	86.8%
	Int. Gradients	56/61	91.8%
	Occlusion	66/75	88.0%
Func. Vectors	Placement	23/24	95.8%
	Func. Type	47/47	100.0%
	Categories	44/48	91.7%
	Layer Evol.	36/36	100.0%
Circuits	Overview	45/46	97.8%
	Subnetwork	132/137	96.4%
	Features	120/125	96.0%

Table 1: Faithfulness verification results of each explanation type across all three analysis pages. All explanations are generated by Qwen2.5-VL-72B and verified against ground-truth data from the underlying analyses.

**Attribution Analysis** Each time a user runs any attribution method, we collect the raw attribution matrix, compute per-token peak and mean contributions, and identify the strongest interactions between input and output tokens. These statistics drive the automatically generated explanation while also feeding the Faithfulness Checker.

**Function Vector Analysis** The Function Vector workflow exports three independent checkpoints. First, the similarity rankings for function types and categories are recomputed from the cached activations, ensuring that any statement about top matches reflects the actual cosine ordering. Second, the PCA narrative is compared with the cluster centroids that underpin the 3D plot, and a verifier must agree that the textual summary follows from those coordinates. Third, descriptions of layer evolution are cross-checked against the activation norms and change magnitudes from the forward pass.

**Circuit Trace Analysis** Circuit tracing explanations operate over graphs extracted from the Cross-Layer Transcoder. For every narrative, we therefore repeat three stages: we confirm structural statements by inspecting the underlying graph (e.g., upstream/downstream connectivity, active features), validate numeric assertions by reading the stored activation values, and run a semantic check that ensures qualitative summaries remain consistent with the graph evidence. Since the same pipeline is applied to the main circuit view, the feature explorer, and the subnetwork explorer, we obtain uniformly perfect scores for the benchmark prompts.

To further validate the causal relevance of the discovered circuits, we perform intervention experiments (Figure 4). We use a set of exemplary prompts covering knowledge retrieval, code generation, and literary analysis, and by ablating the

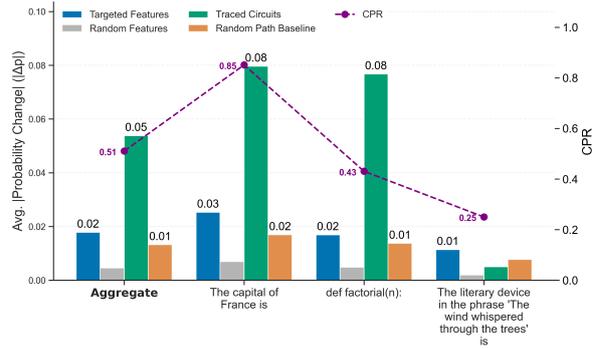


Figure 4: Impact of intervention on model output probability ( $|\Delta p|$ ). We compare the effect of ablating top- $k$  targeted features and traced circuits against random baselines (ablating random features or edges).

features and paths identified by the CLT, there is a substantially larger impact on the model’s output probability compared to random baselines. This confirms that the system successfully isolates functionally critical components. We also compute the Circuit Performance Ratio (CPR) metric (Mueller et al., 2025), which quantifies how well a circuit recovers model performance as a function of the fraction of circuit components included.

## 5 Use Case: Synergistic Analysis of Knowledge Retrieval

To demonstrate how ELIA’s three components can work in concert, we consider a typical knowledge retrieval prompt: “The capital of France is”.

First, the *Attribution Analysis* identifies the token “France” as having the highest saliency score, indicating that the model attends to the subject.

Second, the *Function Vector Analysis* projects the prompt’s activation into the pre-computed semantic space, locating it within the “Abstractive Tasks” cluster. Specifically, it scores highly on the “Next Item” and “Country Capitals” categories, suggesting that the model’s internal state aligns with the high-level task type of factual completion.

Finally, the *Circuit Trace Analysis* reveals the mechanism. In early layers, features related to “article usage” and “country-related information” are active, processing the basic syntax and geographical context. In middle layers, features for “country-related terms” become prominent, linking the context to specific terminology. In late layers, the model synthesizes this information, with high activation in features related to “Geographical knowledge” and “country-related phrases”. This progression shows a systematic buildup from basic grammar to sophisticated geographical knowledge,

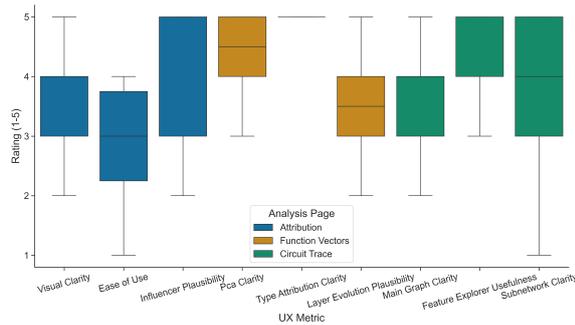


Figure 5: Grouped boxplots of UX ratings for all participants across the three analysis pages.

explaining how the model generates the answer.

This multi-layered approach allows users to build a more comprehensive understanding by combining evidence from different interpretability methods: establishing input dependence, checking semantic alignment, and inspecting components.

## 6 User Study

To empirically evaluate the effectiveness of the explanations and the overall usability of ELIA, a within-subjects user study was conducted with 18 undergraduate computer science students with mostly novice or intermediate experience with LLMs, and designed to assess subjective user experience (UX) and objective comprehension gains. The average completion time was approx. 1h.

### 6.1 Quantitative Results

Participants rated each of the three analysis pages on a 5-point Likert scale according to the following properties: visual clarity, ease of use, plausibility. As shown in Figure 5, the Function Vectors and Circuit Trace pages received significantly more positive UX ratings than the Attribution Analysis page (Kruskal-Wallis H-test,  $p = 0.006$ ). The ‘PCA Clarity’ metric on the Function Vectors page received a perfect median score of 5, while the ‘Ease of Use’ for the Attribution Analysis heatmap received the lowest median score of the study (3), indicating it was confusing for users.

To measure objective comprehension, participants answered three multiple-choice questions per page. While there was a positive trend, a Spearman correlation test found no statistically significant relationship between a user’s prior LLM experience and their correctness score ( $\rho = 0.30, p = 0.23$ ). The average correctness scores were high across all groups: *Experts* achieved a perfect score of 1.00, *Intermediates* scored 0.98, and even *Novices*

reached 0.95. This minimal performance gap suggests that the system helped reduce barriers to comprehension, enabling users with little to no prior knowledge to understand complex model behaviors at levels approaching those of more experienced users.

### 6.2 Qualitative Findings

Analysis of user interview transcripts revealed several key themes. A primary motivation for users was the desire to understand the black box, and the tool was most effective when it provided clear, causal explanations. However, a central challenge was the tension between detail and simplicity; users were often overwhelmed by visualizations that presented too much information at once, such as the main circuit graph.

There is also a clear user preference for interactive visualizations with strong conceptual metaphors. The 3D PCA plot and the Subnetwork Explorer were consistently praised as “very clear”, while more abstract, static visuals like heatmaps were found to be confusing. Finally, a recurring theme was the need for more integrated, automated guidance. Users frequently relied on the facilitator for context and suggested that adding summaries, tooltips, and adaptive complexity would significantly improve the tool’s accessibility.

## 7 Conclusion

ELIA shows that the tools of mechanistic interpretability can become approachable, interactive experiences without giving up analytic depth. By combining the complementary views from attribution heatmaps, function vector projections, and circuit tracing graphs with structured AI narratives and automated faithfulness checks, the platform closes a long-standing accessibility gap: participants in our study reached similar comprehension regardless of prior LLM experience and clearly preferred the explorable interfaces that ELIA provides. The quantitative gains, qualitative feedback, and high verification scores together suggest that interpretability workflows can feel like guided investigations instead of expert-only diagnostics. We hope to encourage the interpretability community to treat usability and faithfulness as co-equal concerns, nudging the field toward tools that invite participation rather than gatekeep expertise.

## Limitations

ELIA is currently limited to two languages (English, German), OLMo as the explained model, and the Qwen-VL model as the explainer model. Richer intervention tools (e.g., counterfactual editing or causal scrubbing) might be necessary to provide a comprehensive user-centric view on interpreting language model behavior. The user study was limited to subjective ratings and short-term interactions, while studying longer-term usage in professional settings remains a necessary future direction to prove the advantages of ELIA we have so far recorded.

## Ethical Statement

The participants in the user study were compensated at or above the minimum wage, in accordance with the standards of our host institutions' regions.

## Acknowledgments

We thank the reviewers of the EACL 2026 System Demonstrations track for their valuable feedback. This research is funded by the Berlin Institute for the Foundations of Learning and Data (BIFOLD, ref. 01IS18037A)

## References

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. 2018. [Sanity checks for saliency maps](#). In *Advances in Neural Information Processing Systems*, volume 31, pages 9505–9515. Curran Associates, Inc.
- Leonard Bereska and Stratis Gavves. 2024. [Mechanistic interpretability for AI safety - a review](#). *Transactions on Machine Learning Research*. Survey Certification, Expert Certification.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. 2023. [Language models can explain neurons in language models](#). OpenAI technical report.
- Nitay Calderon and Roi Reichart. 2025. [On behalf of the stakeholders: Trends in NLP model interpretability in the era of LLMs](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 656–693, Albuquerque, New Mexico. Association for Computational Linguistics.
- Julien Colin, Thomas Fel, Rémi Cadène, and Thomas Serre. 2022. [What I cannot predict, I do not understand: A human-centered evaluation framework for explainability methods](#). In *Advances in Neural Information Processing Systems*, volume 35. Curran Associates, Inc.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. [The faiss library](#). *arXiv Preprints*.
- Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. 2024. [Transcoders find interpretable LLM feature circuits](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Nils Feldhus, Leonhard Hennig, Maximilian Dustin Nasert, Christopher Ebert, Robert Schwarzenberg, and Sebastian Möller. 2023. [Saliency map verbalization: Comparing feature importance representations from model-free and instruction-based methods](#). In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 30–46, Toronto, Canada. Association for Computational Linguistics.
- Nils Feldhus and Laura Kopf. 2025. [Interpreting language models through concept descriptions: A survey](#). In *Proceedings of the 8th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 149–162, Suzhou, China. Association for Computational Linguistics.
- Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R. Costa-jussà. 2024. [A primer on the inner workings of transformer-based language models](#). *arXiv*, abs/2405.00208.
- Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, and 24 others. 2024. [OLMo: Accelerating the science of language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15789–15809, Bangkok, Thailand. Association for Computational Linguistics.
- Michael Hanna, Mateusz Piotrowski, Jack Lindsey, and Emmanuel Ameisen. 2025. [Circuit-tracer: A new library for finding feature circuits](#). In *Proceedings of the 8th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 239–249, Suzhou, China. Association for Computational Linguistics.
- Sarthak Jain and Byron C. Wallace. 2019. [Attention is not Explanation](#). 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 3543–3556, Minneapolis, Minnesota. Association for Computational Linguistics.
- Been Kim, John Hewitt, Neel Nanda, Noah Fiedel, and Oyvind Tafjord. 2025. [Because we have llms, we](#)

- can and should pursue agentic interpretability. *arXiv*, abs/2506.12152.
- Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L Turner, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, and 1 others. 2025. **On the biology of a large language model.** *Anthropic*.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. **Factscore: Fine-grained atomic evaluation of factual precision in long form text generation.** In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 12076–12100. Association for Computational Linguistics.
- Aaron Mueller, Atticus Geiger, Sarah Wiegrefe, Dana Arad, Iván Arcuschin, Adam Belfki, Yik Siu Chan, Jaden Fried Fiotto-Kaufman, Tal Haklay, Michael Hanna, Jing Huang, Rohan Gupta, Yaniv Nikankin, Hadas Orgad, Nikhil Prakash, Anja Reusch, Aruna Sankaranarayanan, Shun Shao, Alessandro Stolfo, and 4 others. 2025. **MIB: A mechanistic interpretability benchmark.** In *Forty-second International Conference on Machine Learning*.
- Neel Nanda, Nelson Elhage, Catherine Olsson, and 1 others. 2022. **In-context learning and induction heads.** Transformer Circuits thread.
- Letitia Parcalabescu and Anette Frank. 2024. **On measuring faithfulness or self-consistency of natural language explanations.** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6048–6089, Bangkok, Thailand. Association for Computational Linguistics.
- Naomi Saphra and Sarah Wiegrefe. 2024. **Mechanistic?** In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 480–498, Miami, Florida, US. Association for Computational Linguistics.
- Gabriele Sarti, Nils Feldhus, Ludwig Sickert, and Oskar van der Wal. 2023. **Inseq: An interpretability toolkit for sequence generation models.** In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 421–435, Toronto, Canada. Association for Computational Linguistics.
- Hendrik Schuff, Alon Jacovi, Heike Adel, Yoav Goldberg, and Ngoc Thang Vu. 2022. **Human interpretation of saliency-based explanation over text.** In *2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT '22*, page 611–636, New York, NY, USA. Association for Computing Machinery.
- Tamar Rott Shaham, Sarah Schwettmann, Franklin Wang, Achyuta Rajaram, Evan Hernandez, Jacob Andreas, and Antonio Torralba. 2024. **A multimodal automated interpretability agent.** In *Forty-first International Conference on Machine Learning*.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. **Axiomatic attribution for deep networks.** In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328. PMLR.
- Ian Tenney, James Wexler, Jasmijn Bastings, Tolga Bolukbasi, Andy Coenen, Sebastian Gehrmann, Ellen Jiang, Mahima Pushkarna, Carey Radebaugh, Emily Reif, and Ann Yuan. 2020. **The language interpretability tool: Extensible, interactive visualizations and analysis for NLP models.** In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 107–118, Online. Association for Computational Linguistics.
- Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2024. **Function vectors in large language models.** In *The Twelfth International Conference on Learning Representations*.
- Igor Tufanov, Karen Hambardzumyan, Javier Ferrando, and Elena Voita. 2024. **Lm transparency tool: Interactive tool for analyzing transformer language models.** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (System Demonstrations)*, pages 51–60. Association for Computational Linguistics. Open-source code available at <https://github.com/facebookresearch/11m-transparency-tool>.
- Jesse Vig. 2019. **A multiscale visualization of attention in the transformer model.** In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 37–42, Florence, Italy. Association for Computational Linguistics.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023. **Interpretability in the wild: a circuit for indirect object identification in GPT-2 small.** In *The Eleventh International Conference on Learning Representations*.
- Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-Garcia, Stevie Bergman, Jackie Kay, Conor Griffin, Ben Bariach, Iason Gabriel, Verena Rieser, and William Isaac. 2023. **Sociotechnical safety evaluation of generative ai systems.** *arXiv*, abs/2310.11986.
- Sarah Wiegrefe and Yuval Pinter. 2019. **Attention is not not explanation.** 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 11–20, Hong Kong, China. Association for Computational Linguistics.
- Fred Zhang and Neel Nanda. 2024. **Towards best practices of activation patching in language models: Metrics and methods.** In *The Twelfth International Conference on Learning Representations*.

## A Attribution Analysis

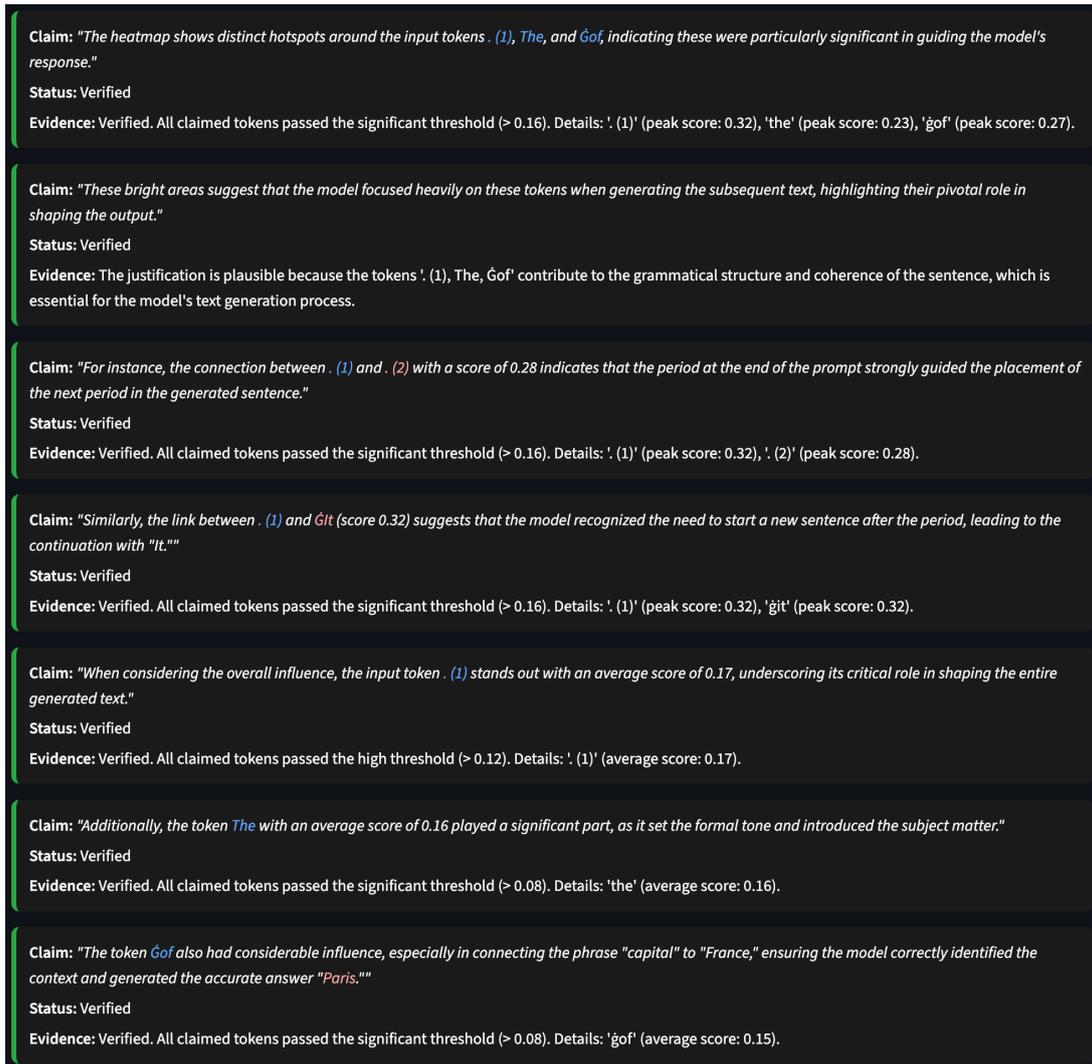


Figure 6: Attribution Analysis: Faithfulness verification of the AI-generated explanation.

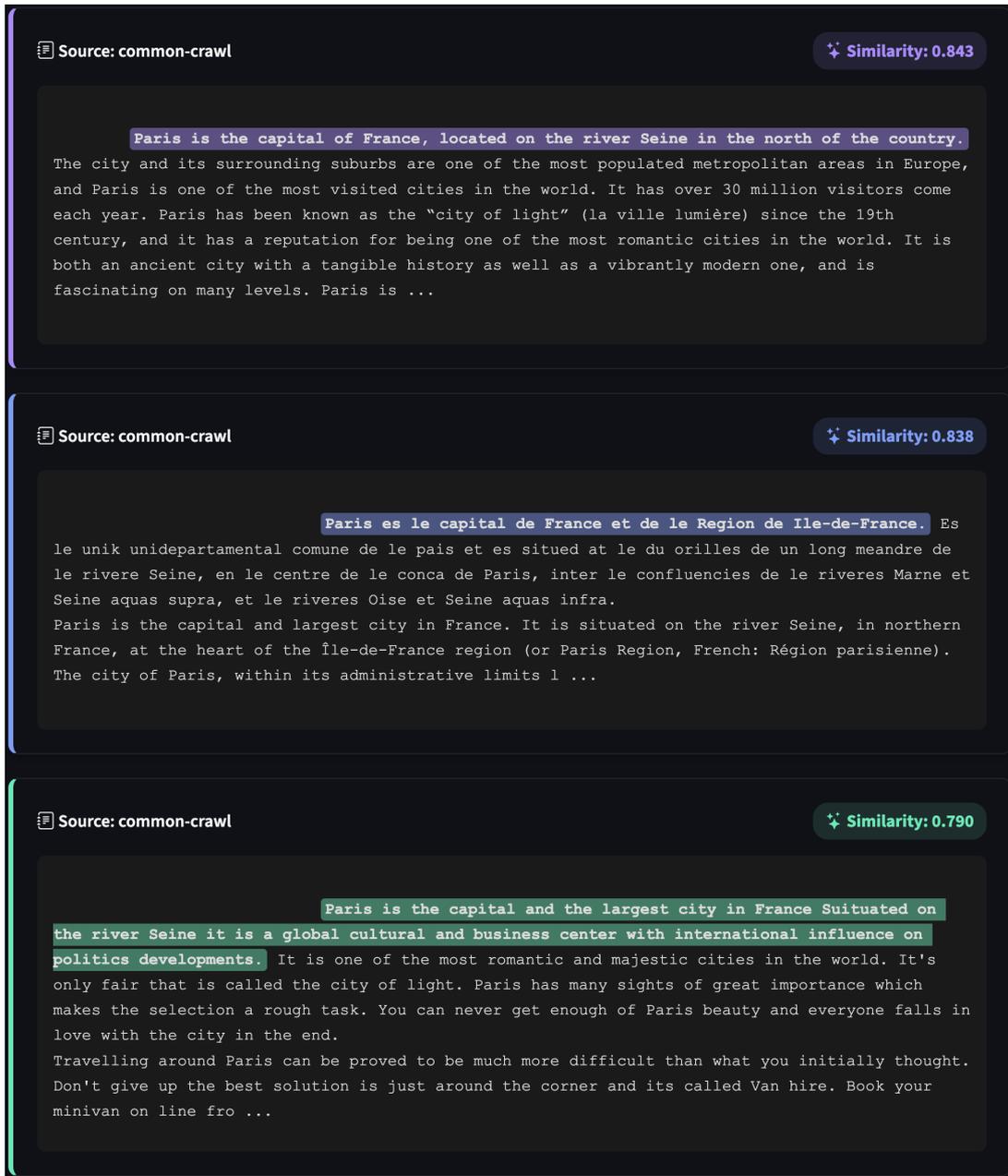


Figure 7: Influence Tracer: Retrieves and displays training documents similar to the prompt.

## B Function Vector Analysis

The Function Vector space is constructed using a custom bilingual dataset (English and German) to ensure robust, cross-lingual functional mappings. For each category, we pass the instructional prompts through OLMo-2-1124-7B and extract the activation vector from the final token position of the last hidden layer. The definitive function vector for a given category is computed as the mean of these activations. To maintain interactivity and minimize computational overhead in the live application, this high-dimensional space is pre-computed offline. During real-time usage, analyzing a new prompt requires only a single forward pass to extract its activation vector, followed by a highly efficient cosine similarity comparison against the pre-cached semantic space.

Function Type	Category	Example Prompts
Abstractive Tasks	country_capital	The capital of France is What city serves as the capital of Japan?
	translation_german	Translating 'hello' into German gives What would a German speaker say for 'world'?
	next_item	After 'Monday' comes The number following 5 is
Multiple Choice QA	commonsense_qa	What happens when you mix red and blue? Why do people wear coats in winter?
	math_qa	What is 15 multiplied by 8? Calculate the area of a square with side 5
	geography_qa	Which is the largest ocean? What is the longest river in the world?
Text Classification	sentiment_analysis	Is this text positive or negative? What emotion does this express?
	language_detection	What language is this text written in? Identify the language of this sentence
	spam_detection	Is this message spam? Classify this email as spam or legitimate
Extractive Tasks	adjective_vs_verb	Is 'running' an adjective or verb? Classify 'beautiful' as adjective or verb
	living_vs_nonliving	Is 'tree' living or non-living? Classify 'car' as living or non-living
	concrete_vs_abstract	Is 'happiness' concrete or abstract? Classify 'table' as concrete or abstract
Named Entity Recognition	ner_person	Identify the person name in this text Extract all person names mentioned
	ner_location	What location is mentioned here? Extract all place names from the text
	ner_organization	Find the organization name Extract company or institution names
Text Generation	complete_sentence	Complete this sentence: "The weather today is" Finish the thought: "In the future, we will"
	continue_story	Continue the story: "Once upon a time..." What happens next in this story?
	question_generation	Generate a question about this topic Create a question based on this text

Table 2: Overview of the Function Vector dataset structure. The dataset contains 6 function types with 120 total categories. Each category includes 5 example prompts. This table shows representative examples from each function type.

Faithfulness Check

**How This Works:** The faithfulness checker verifies three types of claims from the AI's explanation:

- **Ranking Claims:** Checks if a claimed 'most similar' function type or category is actually within the top 3 matches based on cosine similarity scores.
- **Positional Claims:** Semantically verifies if the AI's description of the input's position (e.g., "near text classification") is a plausible summary of the actual top-ranked functions.
- **Justification Claims:** Semantically analyzes whether the reasoning provided for a category's relevance is plausible and logically consistent with the input prompt.

**Claim:** "The user's prompt, "Summarize the plot of 'Hamlet' in one sentence," is positioned within a functional neighborhood dominated by tasks related to text processing and comprehension."

**Status:** Verified

**Evidence:** The claimed functional neighborhood aligns well with 'Abstractive Tasks' as both involve manipulating and understanding textual content. Additionally, tasks like 'Text Classification' inherently require comprehension of text, supporting the relevance of the claimed neighborhood.

**Claim:** "Specifically, it falls into a region characterized by abstractive tasks, text classification, and text generation."

**Status:** Verified

**Evidence:** The claimed functional neighborhood directly corresponds to the actual top-ranked functions, including 'Abstractive Tasks'. This alignment indicates that the claim accurately reflects the primary functionalities without introducing unrelated concepts.

**Claim:** "This placement indicates that the prompt is closely associated with functions that require understanding, summarizing, and generating textual information."

**Status:** Verified

**Evidence:** The justification is plausible because summarizing the plot of 'Hamlet' in one sentence directly involves understanding and generating textual information, which aligns with the specified functional category. The connection between the task of summarization and the category's functions is clear and relevant.

Figure 8: Function Vector Analysis: Faithfulness verification of the AI-generated PCA explanation.

## C Circuit Trace Analysis

The Cross-Layer Transcoder (CLT) is trained offline to learn a sparse, simplified representation of the model’s internal activations. The model was trained on text samples from the Dolma dataset using a batch size of 16 over 1,500 training steps. Our architecture maps the hidden dimension to 512 interpretable features per layer, utilizing a JumpReLU activation (threshold = 0.0) alongside an  $L_1$  sparsity penalty ( $\lambda = 1e^{-3}$ ). Optimization was performed using Adam (learning rate:  $3e^{-4}$ ) with cosine annealing and gradient clipping (max norm: 1.0) to ensure stable convergence. The training dynamics are visualized in Figure 9.

By performing this resource-intensive training offline, the live computational cost of ELIA is drastically reduced. Real-time operations are limited to standard forward passes and asynchronous API calls for the natural language explanations, ensuring the interface remains highly interactive without requiring extensive local compute resources.

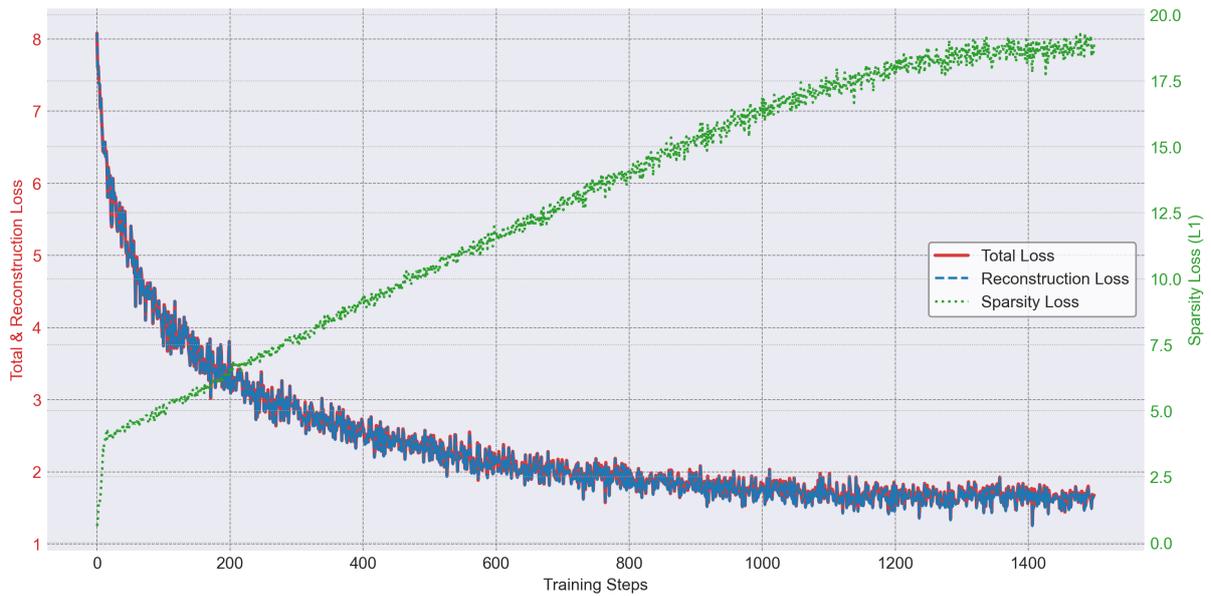


Figure 9: Cross-Layer Transcoder training dynamics showing Total Loss, Reconstruction Loss, and Sparsity Loss ( $L_1$ ) over 1500 training steps. The CLT is trained on the Dolma dataset with  $L_1$  sparsity regularization, gradient clipping, and cosine annealing learning rate scheduling.

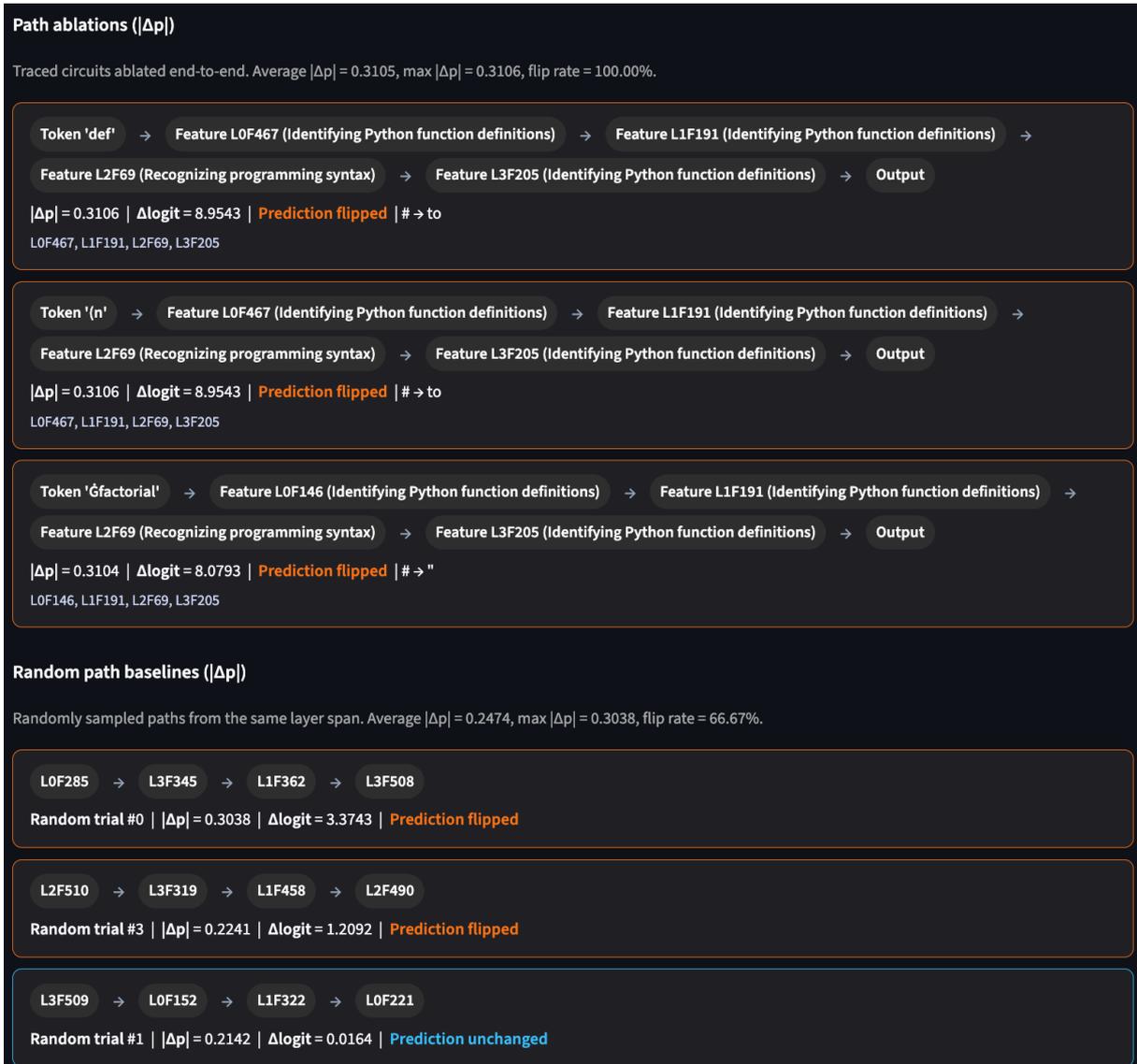


Figure 10: Circuit Trace Analysis: Local path ablations. This view shows the effect of ablating specific paths within the top feature graph shown in the main circuit trace visualization.

<p><b>Claim:</b> "In the early layers of the model, the primary role is to dissect and encode the fundamental elements of the input text."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Summary 'dissect and encode the fundamental elements of the input text': The claimed summary accurately reflects the early layer's role in handling syntax, grammar, and basic patterns, which is a well-established general principle.</p>
<p><b>Claim:</b> "For instance, in Layer 10, a feature focused on 'country-related terms' becomes active, suggesting that the model begins to identify and categorize words associated with countries. Similarly, Layer 4 sees the activation of a feature related to 'country capitals,' indicating an early recognition of the query's intent. Features dealing with 'article usage' and 'country-related phrases' in Layers 7 and 8 further refine the grammatical and contextual understanding of the input."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Verified: 'country-related terms' in L10 matched 'Identifying country-related terms'. Verified: 'country capitals' in L4 matched 'Identifying country capitals'. Verified: 'article usage' in L7 matched 'Identifying article usage'. Verified: 'country-related phrases' in L8 matched 'Identifying country-related phrases'. Semantic reasoning for early layers: The claimed summary accurately reflects the early layers' focus on syntax, grammar, and basic patterns while also acknowledging the identification of country-related terms and geographical context, which aligns with the general principles and provided data points.</p>
<p><b>Claim:</b> "These layers essentially lay the groundwork by recognizing key terms and basic syntactic structures, setting the stage for deeper analysis in subsequent layers."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Summary 'lay the groundwork by recognizing key terms and basic syntactic structures': The claimed summary accurately reflects the early layer's role in recognizing key terms and basic syntactic structures, which aligns with the general principle that early layers handle syntax, grammar, and basic patterns.</p>
<p><b>Claim:</b> "In Layer 21, multiple features related to 'article structures' show high activation, implying that the model is now constructing more complex grammatical frameworks. Additionally, another feature in the same layer focuses on 'geographical references,' which helps to contextualize the country-related terms previously identified."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Verified: 'article structures' in L21 matched 'Identifying article structures'. Verified: 'geographical references' in L21 matched 'Identifying geographical references'. Semantic reasoning for middle layers: The claimed summary accurately reflects the development of thematic connections and abstract meaning by integrating context and identifying article structures and geographical references, which aligns with the general principle of middle layers.</p>
<p><b>Claim:</b> "Layers 19 and 20 continue this trend with features that recognize 'country-related phrases,' building upon the earlier activations and forming a more coherent understanding of the geographical and political context."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Verified: 'country-related phrases' in L19 matched 'Identifying country-related phrases'. Verified: 'country-related phrases' in L20 matched 'Identifying country-related phrases'. Semantic reasoning for middle layers: The claimed summary accurately reflects the general principle of middle layers developing thematic connections and abstract meaning, as well as being semantically consistent with the actual data points.</p>
<p><b>Claim:</b> "This phase marks a transition from basic term recognition to a more nuanced comprehension of the relationships between those terms."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Summary 'transition from basic term recognition to nuanced comprehension': The claimed summary accurately reflects the general principle of middle layers developing thematic connections and abstract meaning, as well as the data indicating a transition from identifying basic terms to more nuanced understanding.</p>
<p><b>Claim:</b> "In the late layers, the model synthesizes all the gathered information to generate the final output."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Summary 'synthesize information for final output': The claimed summary accurately reflects the general principle that late layers synthesize all information to finalize the output, which is supported by the data point stating 'This layer family synthesizes accumulated abstractions to finalize consolidated, coherent outputs ready for downstream use.'</p>
<p><b>Claim:</b> "Layer 31, in particular, exhibits a significant activation of a feature dedicated to 'country names,' which likely plays a crucial role in identifying 'Paris' as the correct answer. Other highly activated features in this layer include 'country-related terms,' 'geographical context,' and 'geographical references,' all contributing to the model's ability to provide a precise and contextually accurate response."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Verified: 'country names' in L31 matched 'Identifying country names'. Verified: 'country-related terms' in L31 matched 'Identifying country-related terms'. Verified: 'geographical context' in L31 matched 'Identifying geographical context'. Verified: 'geographical references' in L31 matched 'Identifying geographical references'. Semantic reasoning for late layers: The claimed summary accurately reflects the late layer's role in synthesizing information to finalize the output, highlighting relevant features like 'country names' and 'geographical context' that contribute to a precise and contextually accurate response.</p>
<p><b>Claim:</b> "These layers act as the decision-making hub, where all the previously processed data converges to produce the final prediction."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Summary 'decision-making hub': The claimed summary 'decision-making hub' aligns with the general principle that late layers synthesize all information to finalize the output, which is supported by the data indicating synthesis and coherent output preparation.</p>
<p><b>Claim:</b> "The high activation levels suggest a strong confidence in the model's output."</p> <p><b>Status:</b> Verified</p> <p><b>Evidence:</b> Summary 'indicate strong confidence': The claimed summary 'indicate strong confidence' is verified as true because it aligns with the general principle that late layers synthesize all information to finalize the output, which inherently suggests a strong confidence in the synthesized results.</p>

Figure 11: Circuit Trace Analysis: Faithfulness verification of the circuit explanation.

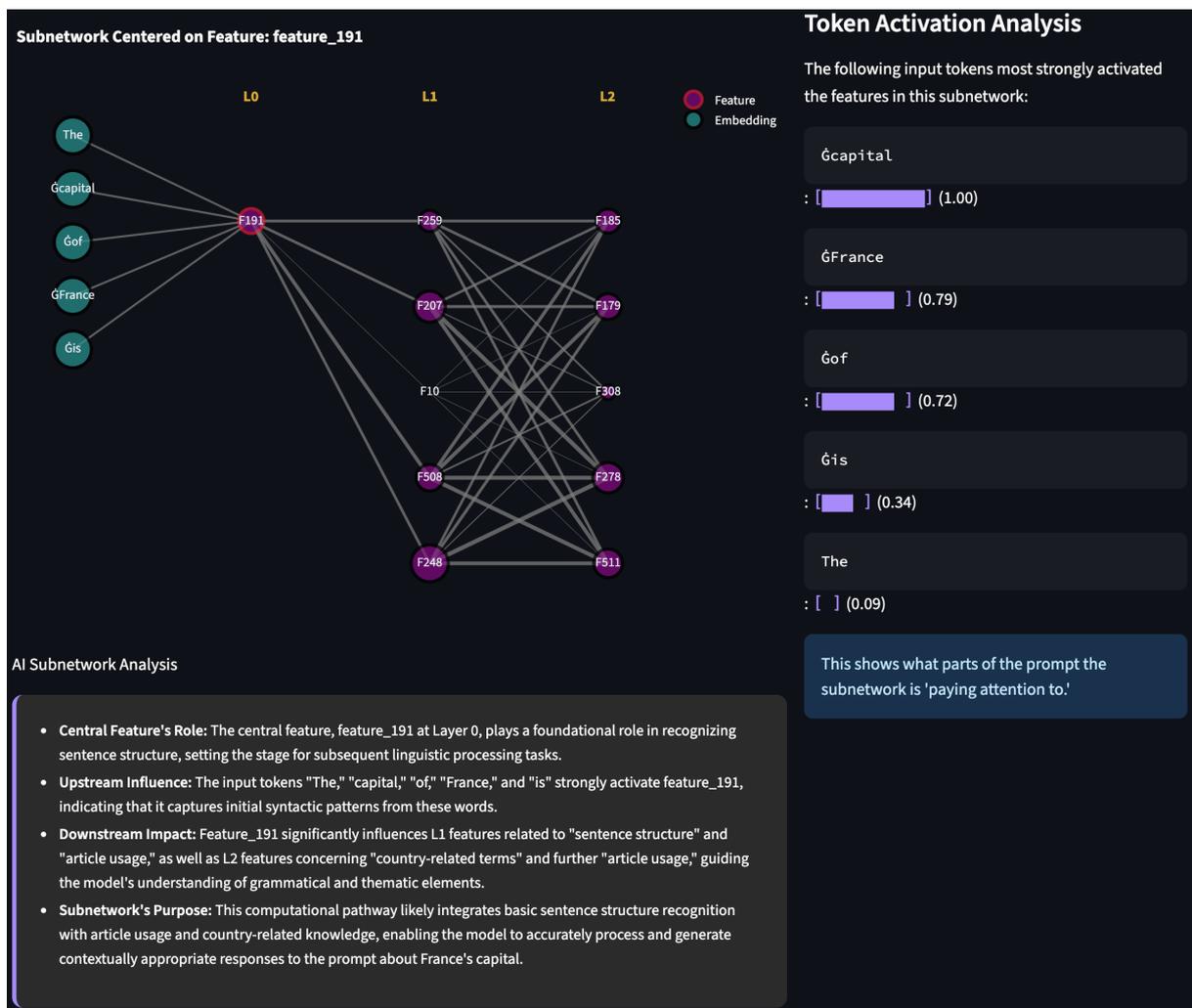


Figure 12: Circuit Trace Analysis: Subnetwork Explorer. This interactive tool allows users to isolate and visualize the specific computational pathways connected to a chosen feature, revealing its upstream influences and downstream effects.

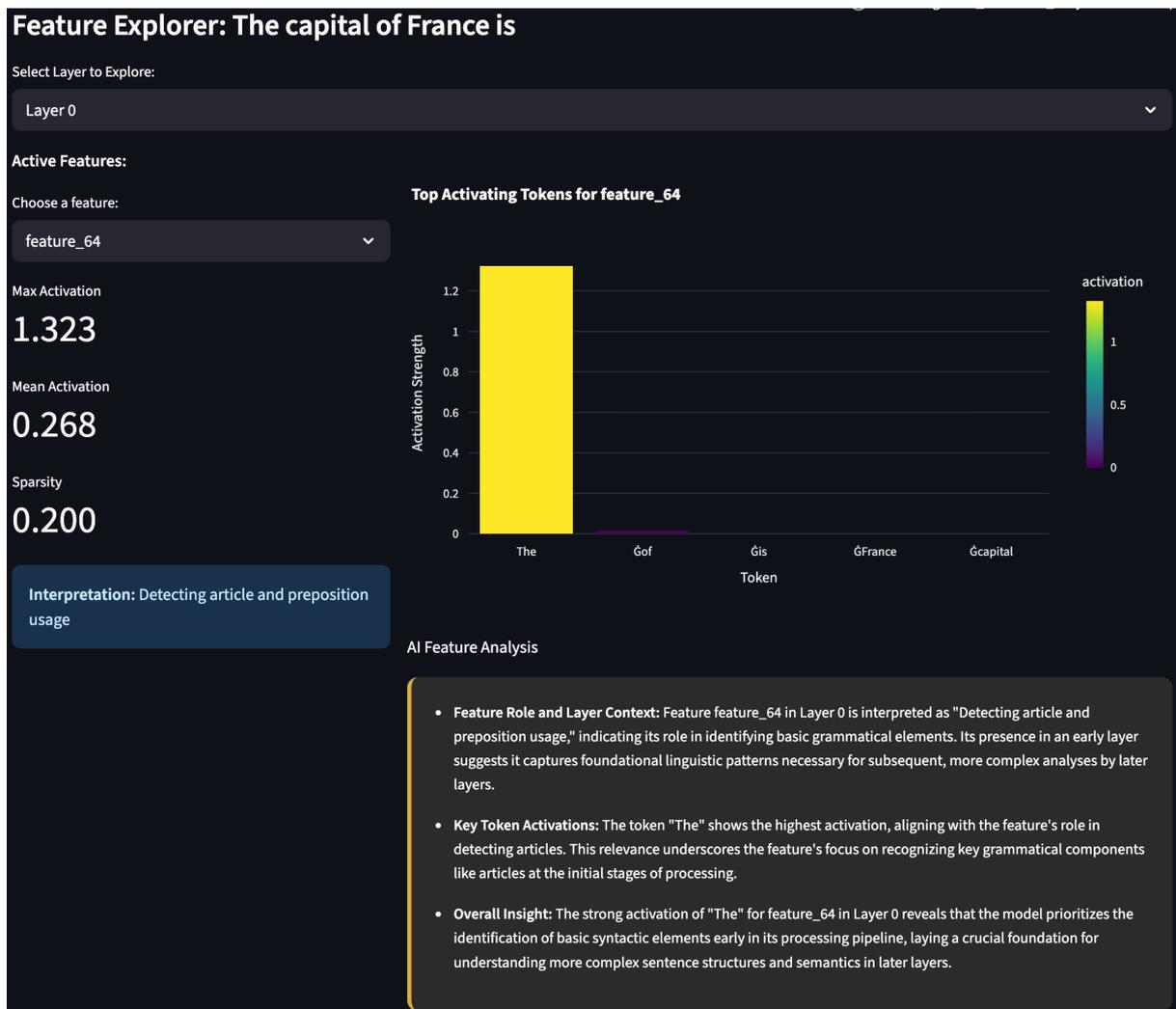


Figure 13: Circuit Trace Analysis: Feature Explorer. Users can inspect individual features in detail, viewing their top activating tokens, sparsity statistics, and AI-generated functional interpretations.

# IntelliCode: A Multi-Agent LLM Tutoring System with Centralized Learner Modeling

Jones David<sup>1</sup>, Shreya Ghosh<sup>2</sup>

<sup>1</sup>School of Computer Science and Engineering, VIT-AP University, India

<sup>2</sup>School of Electrical and Computer Sciences, IIT Bhubaneswar, India

Correspondence: jones.22bce8135@vitapstudent.ac.in, shreya@iitbbs.ac.in

## Abstract

LLM-based tutors are typically single-turn assistants that lack persistent representations of learner knowledge, making it difficult to provide principled, transparent, and long-term pedagogical support. We introduce **IntelliCode**, a multi-agent LLM tutoring system built around a centralized, versioned learner state that integrates mastery estimates, misconceptions, review schedules, and engagement signals. A StateGraph Orchestrator coordinates six specialized agents: skill assessment, learner profiling, graduated hinting, curriculum selection, spaced repetition, and engagement monitoring, each operating as a pure transformation over the shared state under a single-writer policy. This architecture enables auditable mastery updates, proficiency-aware hints, dependency-aware curriculum adaptation, and safety-aligned prompting. Our demo showcases an end-to-end tutoring workflow: a learner attempts a DSA problem, receives a conceptual hint when stuck, submits a corrected solution, and immediately sees mastery updates and a personalized review interval. We report validation results with simulated learners, showing stable state updates, improved task success with graduated hints, and diverse curriculum coverage.<sup>12</sup>

## 1 Introduction

Large Language Models (LLMs) have rapidly expanded the possibilities of automated tutoring, yet most existing systems remain fundamentally reactive: each query is treated in isolation, with little continuity or awareness of a learner’s evolving knowledge (Wu et al., 2025). Human tutors, in contrast, maintain rich, persistent models of student understanding that support targeted feedback, curriculum scaffolding, and long-term learning trajectories (VanLEHN, 2011). This discrepancy limits

<sup>1</sup>Live System: <https://intellicode.redomic.in>,

<sup>2</sup>Video Demo: <https://youtu.be/o08bZfele0U>

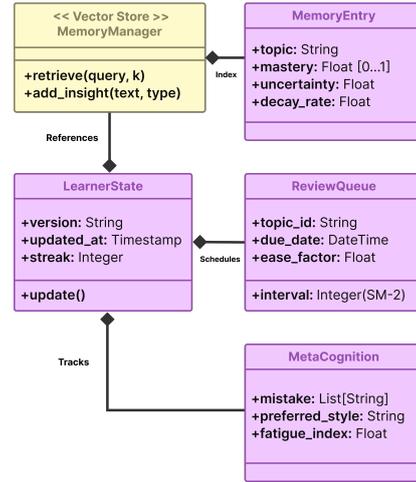


Figure 1: Unified Learner State Schema.

the pedagogical reliability of current LLM tutors, which often provide inconsistent hints, overlook dependencies between concepts, and fail to account for systematic misconceptions (Chu et al., 2025; Mukherjee and Ghosh, 2025). Recent frameworks such as GenMentor (Wang et al., 2025) and SocraticLM (Liu et al., 2024) demonstrate the promise of multi-agent orchestration and structured dialogue in stabilizing tutoring behavior. However, these systems typically rely on ephemeral or implicit memory, lacking an explicit, auditable learner model shared across agents. At the same time, decades of work on learner modeling from Bayesian Knowledge Tracing (BKT) (Corbett and Anderson, 1995) and Deep Knowledge Tracing (DKT) (Piech et al., 2015) to memory networks (Zhang et al., 2017) and PFA (Pavlik Jr et al., 2009), underscore the importance of accurate mastery estimation for personalization. Yet, few LLM-based tutors integrate such formal models with generative reasoning in a persistent, updateable state.

**Our goal is to bridge this gap.** We introduce **IntelliCode**, a multi-agent LLM tutoring system built around a centralized, versioned learner state that

serves as the single source of truth for all pedagogical decisions. Unlike prior systems, IntelliCode enforces a single-writer policy through a *StateGraph Orchestrator*, ensuring that every mastery update, hint intervention, or curriculum choice results from a coherent, formally validated transformation of the learner model. This design mitigates drift, prevents conflicting agent outputs, and enables transparent, multi-turn personalization.

IntelliCode integrates well-established instructional principles: mastery estimation, graduated hinting, dependency-aware curriculum planning, and spaced repetition; with modern LLM capabilities. The learner state encodes mastery vectors with uncertainty, misconceptions, review schedules, and behavioral signals. Six specialized agents (Skill Assessment, Learner Profiler, Pedagogical Feedback, Content Curator, Progress Synthesizer, and Engagement Orchestrator) each operate as pure functions over this shared state. For example, if a learner repeatedly omits a base case in recursion problems, the Profiler records a misconception, the Feedback agent adapts its hint level, and the Curator adjusts upcoming tasks accordingly.

**Contributions.** This work makes the following contributions: (1) a centralized, versioned learner state and single-writer orchestration mechanism enabling consistent, auditable multi-turn tutoring; (2) a suite of pedagogical agents that implement mastery estimation, graduated hinting, curriculum planning, and spaced repetition as pure transformations over the shared state; and (3) a fully functional, end-to-end LLM tutoring system demonstrating stable interaction, interpretable decision-making, and robust content coverage through simulated learner studies.

Figure 1 highlights the unified learner state that drives our adaptation policies. By grounding agent behavior in this structured representation while leveraging LLM reasoning for high-variance tasks such as hinting and code analysis, IntelliCode offers a transparent and pedagogically consistent alternative to memory-less conversational tutoring.

## 2 System Architecture

We frame adaptive personalized education as a Partially Observable Markov Decision Process (POMDP). At each timestep, the learner state  $\mathbf{S}_t$  maintains mastery vectors, SM-2–based review schedules, engagement metrics, and metacognitive memory. Observations  $\mathcal{O}_t$  reflect noisy behavioral

signals such as submissions, errors, and hint requests, while actions  $\mathcal{A}_t$  correspond to pedagogically meaningful interventions including content recommendation, graduated hinting, and schedule adjustments. The reward function  $R_t$  trades off mastery gains with penalties for excessive hint usage and inefficient solve times (details in Appendix B). This formulation enables principled adaptation while supporting the modular multi-agent design showcased in the demo.

### 2.1 Orchestrator Overview

At the core of IntelliCode is the **StateGraph Orchestrator**, the only component permitted to write to the persistent learner record. It maintains a synchronized in-memory copy of the learner state and coordinates all interactions among the six pedagogical agents. When an event occurs, the orchestrator routes it to the appropriate agents, aggregates their outputs, validates the proposed state changes, and then commits them as an atomic update. When multiple agents produce overlapping deltas, the orchestrator applies a fixed priority ordering—Skill Assessment, then Learner Profiler, then remaining agents—and validates each proposed update against the learner state schema before committing. Malformed or conflicting deltas are rejected and logged for audit. Similar to the coordination strategy in GenMentor (Wang et al., 2025), this mechanism prevents conflicting writes, enforces safety and schema constraints, and ensures that the system behaves predictably over multi-turn, long-term tutoring sessions. The orchestrator thus provides the reliability and auditability necessary for principled learner modeling.

### 2.2 Trigger Types and Routing

The orchestrator reacts to pedagogically meaningful events and dispatches them to the relevant agents. These triggers operationalize the full workflow demonstrated in the system:

- **on\_submission:** A code or answer submission triggers Skill Assessment, followed by the Learner Profiler and Pedagogical Feedback.
- **on\_hint\_request:** A learner request for help triggers the Pedagogical Feedback agent, informed by current proficiency and hint history.
- **on\_session\_check:** A daily check-in triggers

Table 1: Roles and responsibilities of the six pedagogical agents.

Agent	Responsibility
Pedagogical Feedback	Provides proficiency-aware, five-level graduated hinting without solution disclosure.
Content Curator	Selects personalized problems based on mastery, dependencies, and the 40/50/10 curriculum policy.
Engagement Orchestrator	Monitors motivation, pacing, and disengagement signals to issue supportive nudges.
Skill Assessment	Performs hybrid evaluation using test-case execution and semantic code review.
Learner Profiler	Estimates mastery deltas, identifies misconceptions, and infers behavioral trends.
Progress Synthesizer	Schedules reviews using an enhanced SM-2 mechanism with context-aware adjustments.

the Content Curator and Engagement Orchestrator.

- **on\_daily\_generation:** The system generates the day’s personalized problem set via the Content Curator.
- **on\_review\_due:** When an SM-2 review is due, the Progress Synthesizer and Content Curator are invoked.

These triggers allow IntelliCode to respond adaptively to the learner’s evolving behavior, maintaining pedagogical continuity across sessions.

### 2.3 Overview of System Agents

Each agent operates as a pure transformation over the shared learner state, producing structured outputs that the orchestrator validates and integrates. Together, the six agents support assessment, personalization, pacing, hinting, and review scheduling. Table 1 summarizes their responsibilities, and Figure 2 shows how the orchestrator mediates their communication. This design ensures that all instructional decisions are grounded in a consistent, auditable learner state and remain traceable throughout the tutoring trajectory.

The overall data flow is illustrated in Figure 2. The StateGraph Orchestrator mediates all communication between agents and the persistent learner state, ensuring that instructional decisions remain traceable and consistent across sessions.

## 3 Learner State and Agent Adaptation

IntelliCode maintains a centralized, versioned learner state that governs all pedagogical decisions. The state is initialized from historical activity and updated using a Bayesian Knowledge Tracing (BKT)–inspired mechanism that incorporates difficulty, recency, hint usage, and solve-time effects. Spaced repetition is managed by an enhanced SM-2 scheduler that computes personalized review intervals based on recall quality and

interaction history. This shared representation enables consistent, long-term adaptation rather than isolated single-turn responses.

To illustrate the update process, consider a learner who solves a recursion problem correctly but requests several hints and exceeds the expected solve time. The BKT-inspired update assigns a modest mastery gain, attenuated by the heavier reliance on hints, while the Progress Synthesizer schedules an earlier review to reinforce retention. The Content Curator then interprets recursion as lying in the learner’s “growth” region and adjusts future problem selection accordingly.

### 3.1 Agent Behaviors and Pedagogical Logic

All six agents operate as pure transformations over the learner state, producing structured outputs validated by the orchestrator before being committed as atomic updates. While the architecture supports fully generative agents, deterministic logic is used for the Learner Profiler and Content Curator for reproducibility, whereas higher-variance components (e.g., hinting, code analysis) leverage LLMs.

**Learner Profiler** The Profiler acts as the diagnostic backbone of the system, identifying mastery deltas, misconceptions, and behavioral trends such as fatigue or decreasing velocity. It consumes correctness, topic tags, error patterns, time-on-task, hint usage, and the current mastery map. For example, if a learner repeatedly omits base cases in recursion problems, the Profiler records a misconception related to termination conditions, which later guides both hinting and content selection. For new learners, the Profiler initializes mastery with uninformative Beta priors ( $\alpha = 1, \beta = 1$ ) and conservative estimates, allowing the system to warm up within 2–3 submissions as observed outcomes rapidly concentrate the posterior.

**Skill Assessment** This agent performs hybrid evaluation by executing test cases and conducting semantic code review. When tests fail, errors alone

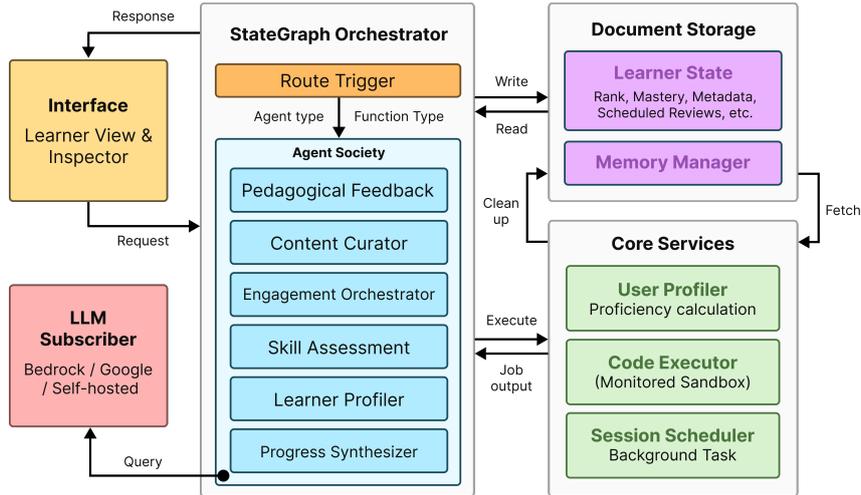


Figure 2: System Architecture: The StateGraph Orchestrator manages the flow between six specialized agents and the persistent learner state. Arrows indicate data flow and trigger events.

are surfaced; when they pass, the agent provides improvement suggestions across categories such as time complexity, space usage, readability, and edge-case coverage. For instance, after a successful merge sort implementation, the agent may recommend reducing auxiliary memory to improve space efficiency.

**Pedagogical Feedback** Informed by Socratic Playground (Zhang and Others, 2024), the Pedagogical Feedback agent employs a five-level graduated hinting protocol:

- **Metacognitive:** prompt the learner to reflect (“What did you try, and what happened?”).
- **Conceptual:** surface the key idea (“This problem relies on identifying a recurrence relation.”).
- **Strategic:** suggest an approach (“Consider breaking the input and solving the halves recursively.”).
- **Structural:** highlight missing logic (“Your solution lacks a base case for empty input.”).
- **Targeted:** point to a region of interest (“Inspect the condition near line 14; termination may not be guaranteed.”).

## System Demonstration

The specificity of hints scales with the learner’s proficiency estimate  $\hat{p}$ . Beginners receive simple analogies and single-step cues, intermediate learners

receive pattern-oriented guidance, and advanced learners receive concise nudges with edge-case emphasis. For the same recursion bug, a beginner might be told to “think of recursion like climbing down a ladder,” while an advanced learner might be prompted to “check whether the termination condition is reachable.”

**Content Curator** The Curator operationalizes the learner state into task selection using a dependency-aware 40/50/10 policy:

$$\text{selection} = 0.4 \times \text{due\_reviews} + 0.5 \times \text{growth\_zone} + 0.1 \times \text{challenge}. \quad (1)$$

Growth-zone items correspond to mastery levels between 0.3 and 0.7, while challenge items target skills below 0.3. The Curator enforces prerequisite dependencies, avoids repetition within a  $k$ -day window, and ensures topic diversity. For example, a learner showing intermediate recursion mastery but weak dynamic programming mastery may receive (i) a recursion review task, (ii) a medium-difficulty DP subproblem, and (iii) a lightweight DP challenge problem.

**Progress Synthesizer** The Progress Synthesizer governs spaced repetition using SM-2 (Wozniak, 1990) and forgetting-curve theory (Ebbinghaus, 1913), augmented with contextual features (Settles and Meeder, 2016; Reddy et al., 2016). Review intervals shrink when hints are heavily used, expand when solutions are fast and confident, and tighten when predicted recall drops near the due date. If the recall probability for a graph traversal concept falls below threshold, the agent prepones

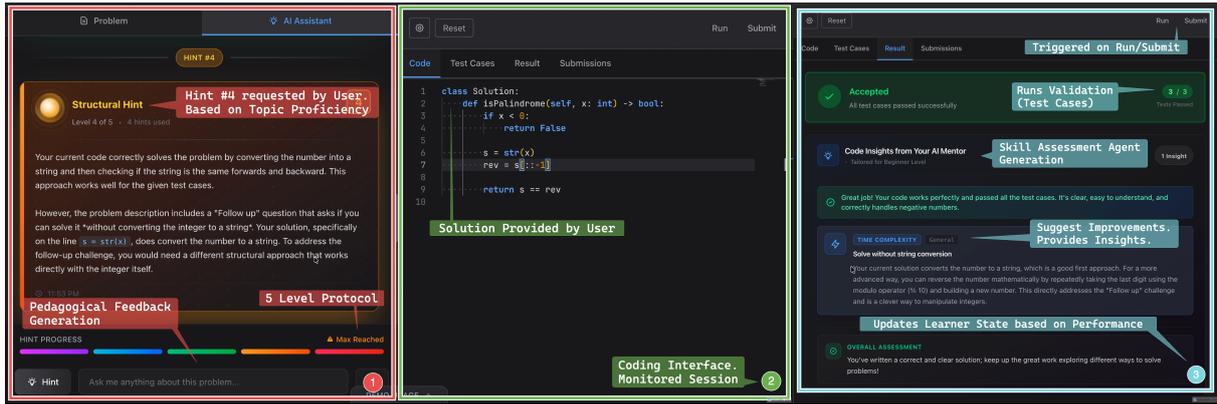


Figure 3: Learner interface featuring the graduated hinting mechanism (Left), the Code Editor (Middle), and proficiency-aware code analysis (Right).

the review—even if the learner has not recently interacted.

**Engagement Orchestrator** Finally, the Engagement Orchestrator monitors motivational signals. It issues supportive prompts after broken streaks, encourages re-engagement after periods of inactivity, and suggests simpler variants when failure streaks accumulate. All interventions are rate-limited and phrased non-judgmentally. For instance, after multiple failed attempts on tree problems, the system may suggest: “Would you like to revisit the easier ‘binary tree basics’ exercise before trying again?”

Collectively, these agents form IntelliCode’s adaptation engine, grounding every hint, problem selection, and review decision in a unified, auditable learner state.

The IntelliCode platform is implemented using FastAPI for the backend, React for the frontend interface, and LangGraph for multi-agent orchestration. A persistent, graph-structured learner model is maintained in ArangoDB, enabling the system to track mastery, misconceptions, and review schedules across sessions. All LLM-powered agents (Pedagogical Feedback, Skill Assessment, Engagement Orchestrator) use Google Gemini 2.5 Flash via API without fine-tuning. Median end-to-end latency for a single agent call is approximately 2–3 seconds, dominated by LLM inference. The backend serves concurrent users via asynchronous request handling in FastAPI; API throughput is rate-limited to 5 requests per second per user to comply with provider quotas.

The demo showcases the full adaptive tutoring loop. A learner begins at a curriculum roadmap and is assigned a data structures and algorithms (DSA) problem selected by the Content Curator using the

40/50/10 policy. If the learner struggles, the Pedagogical Feedback agent produces a proficiency-aligned hint—for example, a Level 2 conceptual cue such as “This problem requires identifying the recurrence pattern.” After incorporating the hint, the learner submits a correct solution. This submission triggers a sequence of coordinated agent behaviors. The Skill Assessment agent validates correctness and offers semantic feedback; the Learner Profiler updates the learner’s mastery estimate for recursion; and the Progress Synthesizer schedules a spaced-repetition review two days later, reflecting the hint usage and solve-time profile. The system interface allows the learner to view mastery trajectories, upcoming reviews, and past interactions, making the adaptation process transparent. This end-to-end interaction exemplifies how IntelliCode integrates real-time assessment, graduated hinting, curriculum adaptation, and spaced repetition into a coherent, state-driven teaching cycle.

## 4 Evaluation Protocols

We evaluate IntelliCode along offline, online, and fairness dimensions to demonstrate the reliability of its learner modeling, content adaptation, and multi-agent interactions.

**Offline Metrics** Our offline analysis focuses on validating the fidelity of the learner model. We measure **Mastery Calibration** using Expected Calibration Error (ECE), which bins predicted mastery values and compares them against observed success rates to assess whether the model’s confidence is well-aligned with actual outcomes. We omit per-problem discriminative metrics such as Brier score and AUROC, as individual problem success in our simulation is influenced by stochastic code gener-

ation and error injection; aggregate mastery is not designed to predict single-problem outcomes, but rather to track long-term learning trajectories. The **Content Policy** is assessed by tracking topic coverage and diversity to ensure balanced curriculum exposure. To isolate the contribution of the multi-agent architecture, we compare against a **stateless baseline** in which learners receive random problem selection, a single attempt per problem, and no agent support (no hints, orchestrator, or code analysis).

**Online Metrics** We also track performance in live interactions. **Learning Gains** are estimated through pre/post mastery changes on held-out assessments. **Engagement** is monitored through streak retention, inactivity gaps, and voluntary practice rates. System **Efficiency** is assessed via median end-to-end latency, with a target of under 500 ms. **Safety** is evaluated through hint acceptance rates (aiming for  $\geq 70\%$ ) and verification that the system never discloses full solutions.

**Fairness Analysis** To ensure equitable support across learner profiles, we compare learning gains, hint levels, and pacing behaviors across proficiency deciles, targeting an interquartile range within 15% of the median. Component-level ablation studies isolating the individual contributions of the Content Curator, Pedagogical Feedback agent, and SM-2 scheduler are planned for future work.

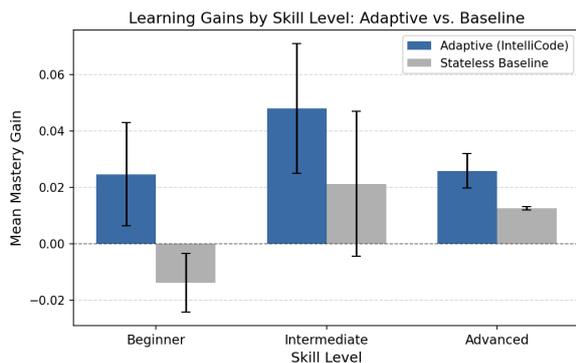


Figure 4: Adaptive IntelliCode consistently outperforms the stateless baseline across all skill levels. Baseline beginners show negative mastery growth in the absence of scaffolding.

### Validation with Simulated Learners

To assess the responsiveness and stability of the multi-agent architecture, we conducted controlled simulations using agent-based learner personas (Wu et al., 2025), drawing on methodolo-

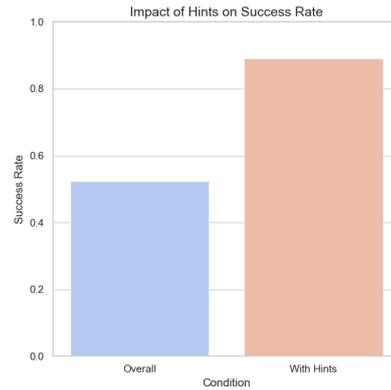


Figure 5: Success rates with and without hint utilization.

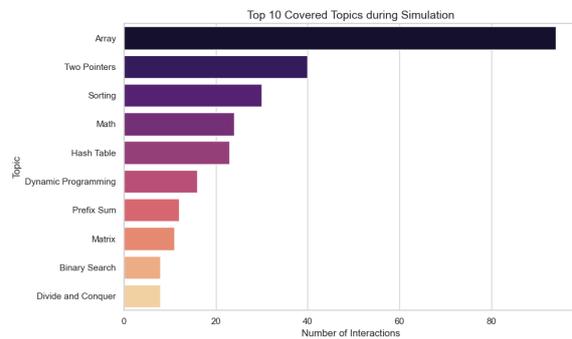


Figure 6: Top 10 topics covered during the simulation.

gies from generative agent societies (Park et al., 2023). While these simulations cannot substitute for human studies in evaluating educational efficacy, they provide a controlled comparison of the system’s adaptation capabilities across diverse cognitive profiles. Ten synthetic learners (4 beginner, 4 intermediate, 2 advanced) each completed a 14-day learning trajectory under two conditions: (1) the full adaptive IntelliCode system with mastery-based problem selection, graduated hinting, and up to three attempts per problem; and (2) a stateless baseline with random problem selection, a single attempt, and no agent interactions.

### 4.1 Comparative Results

Table 2 summarizes the results. The adaptive system achieved a mean mastery gain of 3.42%, compared to 0.54% for the stateless baseline—a 6.3 $\times$  improvement (Figure 4). Learning velocity was similarly differentiated: adaptive learners gained 0.0029 mastery points per active day, versus 0.0004 for the baseline. Notably, baseline beginners experienced *negative* mastery growth ( $-1.39\%$ ), suggesting that without adaptive scaffolding, weaker learners regress due to poorly matched problem

Table 2: Adaptive vs. stateless baseline across 10 simulated learner trajectories over a 14-day horizon.

Metric	Adaptive	Baseline
Mastery Gain (Mean)	0.0342	0.0054
ECE ( $\downarrow$ )	0.043	0.286
Success Rate	89.1%	52.4%
Learning Velocity (gain/day)	0.0029	0.0004
Topics Covered	38	36

difficulty and the absence of corrective feedback. In contrast, adaptive beginners achieved positive gains (+2.46%).

Mastery calibration, measured by ECE, showed a large separation: 0.043 (adaptive) versus 0.286 (baseline). This indicates that the adaptive system’s mastery estimates are well-aligned with observed success rates, whereas the baseline’s predictions are poorly calibrated. The graduated hinting mechanism contributed meaningfully: tasks in which simulated learners requested hints exhibited a success rate of 89.1%, compared with 52.4% without hints (Figure 5), confirming that the Pedagogical Feedback agent provides effective conceptual guidance without disclosing solutions.

## 4.2 Content Coverage

The Content Curator maintained strong diversity across topics (Figure 6), demonstrating the ability of the 40/50/10 policy to avoid topic starvation while respecting prerequisite relationships. Even in extended sessions, the system preserved balanced coverage across skill areas, validating the orchestrator’s ability to coordinate long-term learning arcs and manage curriculum progression.

## 5 Conclusion

In this paper, we presented IntelliCode, a principled multi-agent LLM framework for adaptive education built around a persistent, auditable learner state. By integrating formal mastery update rules, proficiency-aware graduated hinting, dependency- and fairness-aware curriculum adaptation, and safety-aligned prompting, the system offers transparent and consistent multi-turn tutoring capabilities. Controlled comparison against a stateless baseline demonstrates that the adaptive multi-agent architecture produces substantially higher learning gains and better-calibrated mastery estimates, while maintaining robust content coverage and effective hinting interventions. Future work will extend IntelliCode through larger-scale human studies, offline policy evaluation, and federated learn-

ing integrations to strengthen privacy guarantees. We envision IntelliCode as a foundation for next-generation educational systems that blend modern LLM reasoning with established principles from learner modeling, instructional design, and cognitive science.

## Limitations

While IntelliCode demonstrates promising architectural and pedagogical capabilities, there are a few limitations. First, the mastery estimates rely on BKT/DKT-inspired proxies that require sufficient interaction scale and careful calibration; cold-start learners, in particular, necessitate conservative priors and may experience reduced personalization in early sessions. To mitigate this, the system initializes mastery with uninformative priors and applies conservative problem selection during the first few sessions, progressively refining estimates as interaction data accumulates. Second, LLM-driven components introduce variability due to model drift, occasional refusals, and cost constraints. Our guardrails and validation schemas mitigate these issues but cannot fully eliminate them. Finally, rigorous fairness evaluation requires diverse and representative datasets, and remains vulnerable to selection bias, behavior-signal noise, LLM drift, data leakage, and survivorship bias. These considerations underscore the need for large-scale, longitudinal human studies in future work.

## Ethical Considerations

The deployment of LLM-based agents in educational settings necessitates careful attention to accuracy, dependency, and privacy. While IntelliCode integrates verifiers like the *Skill Assessment* agent to validate code logic, generative components such as the *Pedagogical Feedback* agent remain susceptible to hallucinations or plausible but incorrect explanations. Consequently, the system is designed to function as a supplemental tutor rather than a replacement for formal instruction, and we recommend its use under the guidance of human educators who can monitor for potential deviations.

To mitigate the risk of learner over-dependence on AI assistance, we implemented strict graduated hinting protocols. However, we acknowledge that prolonged reliance on automated scaffolding may impact unassisted problem-solving capabilities. Our design prioritizes metacognitive prompting over direct solution disclosure to foster genuine

skill acquisition.

Regarding data privacy, all learner interactions are processed with strict minimization principles. Personally Identifiable Information (PII) is redacted prior to agent ingestion, and the *Learner Profiler* operates on anonymized mastery integers rather than raw user profiles. Finally, while our curriculum policy incorporates fairness constraints to ensure equitable topic coverage, the underlying datasets used for cold-start calibration may inherently reflect historical biases, requiring ongoing monitoring of learning outcomes across diverse demographic groups.

## Acknowledgments

SPARC (Scheme for Promotion of Academic and Research Collaboration) Phase-III (Project ID: 3385) and TIH IIT Tirupati (IITNiF/TPD/2024-25/P16) partially supported this research work.

## References

- Zhendong Chu, Shen Wang, Jian Xie, Tinghui Zhu, Yibo Yan, Jingheng Ye, Aoxiao Zhong, Xuming Hu, Jing Liang, Philip S. Yu, and Qingsong Wen. 2025. [LLM agents for education: Advances and applications](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 13782–13810, Suzhou, China. Association for Computational Linguistics.
- Albert T Corbett and John R Anderson. 1995. Knowledge tracing: Modeling the acquisition of procedural knowledge. *User modeling and user-adapted interaction*, 4(4):253–278.
- Hermann Ebbinghaus. 1913. *Memory: A contribution to experimental psychology*. Teachers College, Columbia University.
- Jiayu Liu, Zhenya Huang, Tong Xiao, Jing Sha, Jinze Wu, Qi Liu, Shijin Wang, and Enhong Chen. 2024. [SocraticLM: Exploring socratic personalized teaching with large language models](#). In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Arka Mukherjee and Shreya Ghosh. 2025. [mmJEE-eval: A bilingual multimodal benchmark for evaluating scientific reasoning in vision-language models](#). In *Proceedings of the 14th International Joint Conference on Natural Language Processing and the 4th Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics*, pages 2268–2290, Mumbai, India. The Asian Federation of Natural Language Processing and The Association for Computational Linguistics.
- Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. [Generative agents: Interactive simulacra of human behavior](#). *Preprint*, arXiv:2304.03442.
- Philip I Pavlik Jr, Hao Cen, and Kenneth R Koedinger. 2009. Performance factors analysis—a new alternative to knowledge tracing. In *Artificial Intelligence in Education*, pages 531–538. IOS Press.
- Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J Guibas, and Jascha Sohl-Dickstein. 2015. Deep knowledge tracing. In *Advances in neural information processing systems*, pages 505–513.
- Sashank Reddy, Igor Labutov, Siddhartha Banerjee, and Thorsten Joachims. 2016. Unbounded human learning: Optimal scheduling for spaced repetition. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1815–1824.
- Burr Settles and Brendan Meeder. 2016. A trainable spaced repetition model for language learning. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1248–1258.
- KURT VanLEHN. 2011. [The relative effectiveness of human tutoring, intelligent tutoring systems, and other tutoring systems](#). *Educational Psychologist*, 46(4):197–221.
- Xinyi Wang, Jiacheng Li, Yu Zhang, and Others. 2025. Llm-powered multi-agent framework for goal-oriented learning. In *Proceedings of the AAI Conference on Artificial Intelligence*.
- Piotr A Wozniak. 1990. Optimization of learning. *Master’s thesis, University of Technology in Poznan*.
- Tao Wu, Jingyuan Chen, Wang Lin, Mengze Li, Yumeng Zhu, Ang Li, Kun Kuang, and Fei Wu. 2025. [Embracing imperfection: Simulating students with diverse cognitive levels using llm-based agents](#). *Preprint*, arXiv:2505.19997.
- Jiani Zhang, Xingjian Shi, Irwin King, and Dit-Yan Yeung. 2017. Dynamic key-value memory networks for knowledge tracing. In *Proceedings of the 26th international conference on World Wide Web*, pages 765–774.
- Yifan Zhang and Others. 2024. SPL: Socratic playground for learning with large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*.

## A Resources and Availability

To support reproducibility and further research, we provide open access to our platform components and simulation data under the **MIT License**:

- **Live Demo:** <https://intellicode.redomic.in>

- **Frontend:** <https://github.com/Redomic/Intellicode-frontend>
- **Backend:** <https://github.com/Redomic/Intellicode-backend>
- **Simulation Framework:** [https://github.com/Redomic/intellicode\\_student\\_sim](https://github.com/Redomic/intellicode_student_sim)

## B Mathematical Formulation Details

We formulate adaptive personalized education as a Partially Observable Markov Decision Process (POMDP):

$$\text{POMDP} = (\mathcal{S}, \mathcal{A}, \mathcal{O}, T, R, \gamma, b_0) \quad (2)$$

### B.1 State Space

The learner state at time  $t$  is:

$$\mathbf{S}_t = \{m_t, r_t, e_t, p_t, M_t, v_t\} \quad (3)$$

where:

- $m_t$ : mastery vector,  $m_{t,i} \in [0, 1]$  for topic  $i \in T$
- $r_t$ : review schedule, items with  $(q_{id}, \text{topics}, d_{\text{due}}, \text{interval}, \text{EF}, n_{\text{reviews}})$
- $e_t$ : engagement state, streak, last-seen timestamp, recent activity windows
- $p_t$ : preferences, skill level, modality, time budget, opt-outs
- $M_t$ : long-term memory, structured text sections on trends, misconceptions, insights
- $v_t$ : version, timestamp for auditing

We also maintain uncertainty  $u_{t,i}$  per topic, encoded as Beta parameters  $(\alpha_{t,i}, \beta_{t,i})$ .

### B.2 Observation Space

Observations are partial, noisy signals of state:

$$\mathcal{O}_t \in \{ \text{submission}, \text{hint\_request}, \text{session\_start}, \text{due\_review} \} \quad (4)$$

For each submission, we observe:

$$o_t = (q_{id}, y, \tau, h_{\text{cnt}}, \text{errors}, t_{\text{solve}}) \quad (5)$$

where  $y \in \{0, 1\}$  (pass/fail),  $\tau$  (timestamp),  $h_{\text{cnt}}$  (hints used), errors (semantic signals),  $t_{\text{solve}}$  (time on task).

### B.3 Action Space

The orchestrator (via agents) selects actions:

$$\mathcal{A}_t \in \{ \text{recommend\_item}, \text{hint}(l), \text{adjust\_schedule}, \text{intervene}, \text{feedback}(d) \} \quad (6)$$

where  $l \in \{1, 2, 3, 4, 5\}$  is hint level,  $d$  is feedback detail level.

### B.4 Reward Function

The reward proxies learning progress while penalizing inefficiency:

$$R_t = \underbrace{w_m \Delta m_t}_{\text{mastery gain}} + \underbrace{w_r \mathcal{K}[\text{review\_success}]}_{\text{retention}} - \underbrace{w_h h_{\text{cnt}}}_{\text{hint penalty}} - \underbrace{w_t \max(0, t_{\text{solve}} - \mu_t)}_{\text{time penalty}} \quad (7)$$

with regularizers for fairness (topic coverage) and engagement.

## C Learner State Updates

### C.1 State Initialization

From historical submissions, we compute initial mastery using a recency-weighted exponential moving average:

$$m_{t,i}^{(0)} = 0.6 \text{ success\_rate}_i + 0.4 \text{ recent\_success\_rate}_i + \mathcal{N}(0, \sigma_0^2) \quad (8)$$

We initialize Beta parameters as  $(\alpha_0, \beta_0) = (1, 1)$  (uninformative prior), and review queue empty with ease factor  $\text{EF}_0 = 2.5$ .

### C.2 Mastery Update Rule

Upon outcome  $y \in \{0, 1\}$  on a question tagged with topics  $Q$ :

$$m_{t,i} \leftarrow \begin{cases} \min(1, m_{t,i} + \alpha w_d w_r (1 - m_{t,i})) & (y=1), \\ \max(0, m_{t,i} - \beta w_d^{-1} w_r m_{t,i}) & (y=0) \end{cases} \quad (9)$$

where:

- $w_d \in \{0.8, 1.0, 1.2\}$  maps difficulty  $\in \{\text{Easy}, \text{Medium}, \text{Hard}\}$  (inverted for failure penalties)
- $w_r = \exp(-\Delta t / \tau_{\text{upd}})$  decays with recency

- Hint/time penalties:  $m_{t,i} \leftarrow m_{t,i} - \eta_h h_{\text{used}} - \eta_t \max(0, t_{\text{solve}} - \mu_i)$
- Momentum smoothing:  $m_{t,i} \leftarrow (1 - \lambda)m_{t,i} + \lambda m_{t,i}^{\text{new}}$  reduces jitter

### C.3 Proficiency Composite

Overall proficiency is a weighted composite:

$$\hat{p} = \sum_k w_k s_k \quad (10)$$

where  $s_k$  includes: topic mastery average (0.40), expertise rank (0.25), self-reported skill (0.20), recent success rate (0.10), streak normalization (0.05).

### C.4 Spaced Repetition Updates (SM-2)

For a review item with ease factor EF, we derive quality score  $q \in \{0, 1, 2, 3, 4, 5\}$  from:

$$q = \begin{cases} 5 & \text{fast, no hints} \\ 4 & \text{solved, minor delay} \\ 3 & \text{solved with hints} \\ \leq 2 & \text{failed/forgot} \end{cases} \quad (11)$$

Update:

$$EF' = \max(1.3, EF - 0.8 + 0.28q - 0.02q^2) \quad (12)$$

Intervals:  $I_1 = 1, I_2 = 6, I_n = \text{round}(I_{n-1} \cdot EF')$  days. Predicted recall:

$$R(\Delta t) = \exp(-\Delta t / \tau) \quad (13)$$

where  $\tau \propto EF'$ .

## D Reproducibility Details

1. Seeded random splits; fixed topic-graph snapshot.
2. Frozen prompt versions and role texts.
3. Logged hyperparameters:  $\alpha, \beta, w_d, \tau_{\text{upd}}, \lambda$ , proficiency weights.
4. Validation schemas for agent outputs (JSON specs).
5. Behavior signal preprocessing with masking rules.
6. Ablation code and offline evaluation scripts.

# FiMMIA: scaling semantic perturbation-based membership inference across modalities

**Anton Emelyanov**  
SberAI  
login-const@mail.ru

**Sergei Kudriashov**  
Sber, HSE University  
sakudryashov@hse.ru

**Alena Fenogenova**  
SberAI  
alenush93@gmail.com

## Abstract

Membership Inference Attacks (MIAs) aim to determine whether a specific data point was included in the training set of a target model. Although there have been numerous methods developed for detecting data contamination in large language models (LLMs), their performance on multimodal LLMs (MLLMs) falls short due to the instabilities introduced through multimodal component adaptation and possible distribution shifts across multiple inputs. In this work, we investigate multimodal membership inference and address two issues: first, by identifying distribution shifts in the existing datasets, and second, by releasing an extended baseline pipeline to detect them. We also generalize the perturbation-based membership inference methods to MLLMs and release **FiMMIA** — a modular **F**ramework for **M**ultimodal **M**IA.<sup>1</sup> We propose to train a neural networks to analyze the target model’s behavior on perturbed inputs, capturing interactions between semantic domains and loss values on members and non-members in the local neighborhood of each sample. Comprehensive evaluations on various fine-tuned multimodal models demonstrate the effectiveness of our perturbation-based membership inference attacks in multimodal settings.

## 1 Introduction

The development of MLLMs has exceeded expectations (Liu et al., 2023a; Lin et al., 2023), showcasing extraordinary performance on various multimodal benchmarks (Chervyakov et al., 2025; Lu et al., 2022; Liu et al., 2023b; Song et al., 2024), even surpassing human performance. However, due to the partial obscurity associated with MLLMs training or fine-tuning (OpenAI, 2023; Reid et al., 2024), it remains challenging to definitively ascer-

tain the impact of training data on model performance, despite some works showing the employment of the training set of certain datasets (Liu et al., 2023a; Chen et al., 2023; Bai et al., 2023). The issue of data contamination occurs when training or test data of benchmarks is exposed during the model training or fine-tuning phase (Xu et al., 2024) and could potentially instigate inequitable performance comparisons among models.

Although numerous works in the field of LLMs have proposed methods for detecting data contamination (Mozaffari and Marathe, 2024; Hu et al., 2022a; Song et al., 2025; Li et al., 2024b), MLLMs, due to their various modalities that, in most implementations, lack corresponding target tokens for multimodal inputs, while multiple training phases, common for MLLM training, complicate an inference when one tries to apply these methods directly. Therefore, there is a necessity in a multimodal contamination detection framework specifically tailored for MLLMs. Our main contributions can be summarized as follows:

- We extended the work of Das et al. (2024) to multimodal data and assessed image as well as recent text MIA benchmarks (Fu et al., 2025; Hallinan et al., 2025) for distribution shifts. We have found that even the *most recent proposed benchmarks are subject to distribution shifts between member and non-member data.*
- We *release a baseline attack pipeline for text, image, video and audio data*, that collects various statistics from the dataset distribution and trains a classifier on top to distinguish members from non-members without any signal from the target model.
- We *extend perturbation-based MIA methods to MLLMs*, revealing their effectiveness and transferability even at the scale of billion-parameter models.
- We *release a modular framework FiMMIA* supporting diverse datasets, modalities, and

<sup>1</sup>The source code and framework have been made publicly available under the MIT license via [link](#). The video demonstration is available on [YouTube](#).

neighbor generation methods. Our pipelines support MIA in multiple settings: when only text, multimodal or both parts are assumed to be leaked.

## 2 Related Work

### 2.1 Data contamination and distribution shifts hinder reliable evaluations

Preserving training data confidentiality is critical for LLMs, as their datasets can contain sensitive private information and tests (Yeom et al., 2018; Hu et al., 2022b). Additionally, data contamination between training and test sets undermines benchmark reliability and complicates model comparison (Balloccu et al., 2024; Sainz et al., 2023), driving recent adoption of dynamically updated benchmarks (White et al., 2025).

Distribution shifts pose significant risks as neural networks’ ability to extract subtle correlations makes them vulnerable to adversarial examples (Moayeri et al., 2022), spurious correlations in explanations (Ribeiro et al., 2016), and data poisoning (Souly et al., 2025). Recent studies have also found that modern LLMs are capable of intentional *sandbagging*, i.e., strategically underperforming during the evaluations in the presence of an incentive to do so (van der Weij et al., 2024). In other words, capable LLMs can intentionally manipulate their logprobs, which poses an additional challenge both for capability elicitation and loss-based MIA attacks<sup>2</sup>.

### 2.2 Membership inference attacks aim to solve the problem

Membership Inference Attacks (MIAs) determine whether a data sample was part of a model’s training set (Shokri et al., 2017) or originates from the general distribution. As noted by (Carlini et al., 2022), this constitutes a hypothesis testing task that crucially relies on the i.i.d. assumption.

Membership Inference Attacks have been the subject of considerable research across a variety of machine learning models, including classification models (Long et al., 2018; Song et al., 2019; Choquette-Choo et al., 2021), generative models (Hayes et al., 2017; Hilprecht et al., 2019; Chen et al., 2020), and embedding models (Song and Raghunathan, 2020; Mahloujifar et al., 2021). The

<sup>2</sup>Such behavior is only possible if the evaluation data or environment presents enough evidence to distinguish it from the training environment, even due to subtle cues.

	Dataset / task	Best reported(%)	Our baseline(%)
text	WikiMIA-hard	64.0 (Hallinan et al., 2025)	57.7 ± 2.5
	WikiMIA-24	99.8 (Fu et al., 2025)	99.9 ± 0.1
	VL-MIA-Text (32 tok.)	96.2 (Li et al., 2024c)	84.9 ± 4.0
	VL-MIA-Text (64 tok.)	99.3 (Li et al., 2024c)	95.5 ± 0.9
image	VL-MIA-Flickr	94.2 (Yin et al., 2025)	99.1 ± 0.4
	VL-MIA-Flickr-2k	74.0 (Li et al., 2024c)	98.6 ± 0.4
	VL-MIA-Flickr-10k	NA	99.3 ± 0.1
	VL-MIA-DALL-E	84.0 (Yin et al., 2025)	99.9 ± 0.1
	LAION-MI*	2.42 (Dubinski et al., 2023)	1.11 ± 0.1

Table 1: AUC-ROC Evaluations of image and text MIA datasets for the occurrence of distribution shifts between members and non-members data. \* corresponds to TPR@1FPR instead. Datasets with no distribution shifts between members and non-members should display values of **50%** for AUC-ROC and **0** for TPR@1FPR.

appearance of LLMs has likewise led to numerous studies investigating membership inference attacks against them (Mireshghallah et al., 2022; Fu et al., 2023; Shi et al., 2024; Mattern et al., 2023). However, the field of MIAs for multimodal models is still in its nascent stages and requires further exploration, facing challenges due to the absence of targets for modality-related tokens, instabilities from multimodal adaptation etc. Several methods (Ko et al., 2023; Hu et al., 2022d) proposed to conduct MIAs based on the similarity between an image and its associated text label. However, this technique is limited to the presence of a paired entry (pair image/text), not the presence of a solitary image or text sequence.

MIAs are commonly categorized into metric-based and shadow model-based approaches (Hu et al., 2022b). Metric-based MIAs (Yeom et al., 2018; Salem et al., 2018; Song and Mittal, 2021; Shi et al., 2024) compare model output statistics against a threshold, while shadow model-based methods (Shokri et al., 2017; Salem et al., 2018) require computationally expensive model replication. Recent work has introduced semantic MIAs (Koike et al., 2025; Mozaffari and Marathe, 2024) that exploit local model properties through sample perturbations. We extend this semantic approach to image, audio, video, and text modalities.

## 3 FiMMA

### 3.1 Overview

The system is the first collection of models and pipelines for membership inference attacks against LLMs, built and evaluated initially on the Russian language, and extendable to any other language or MMLM dataset. The pipeline supports differ-

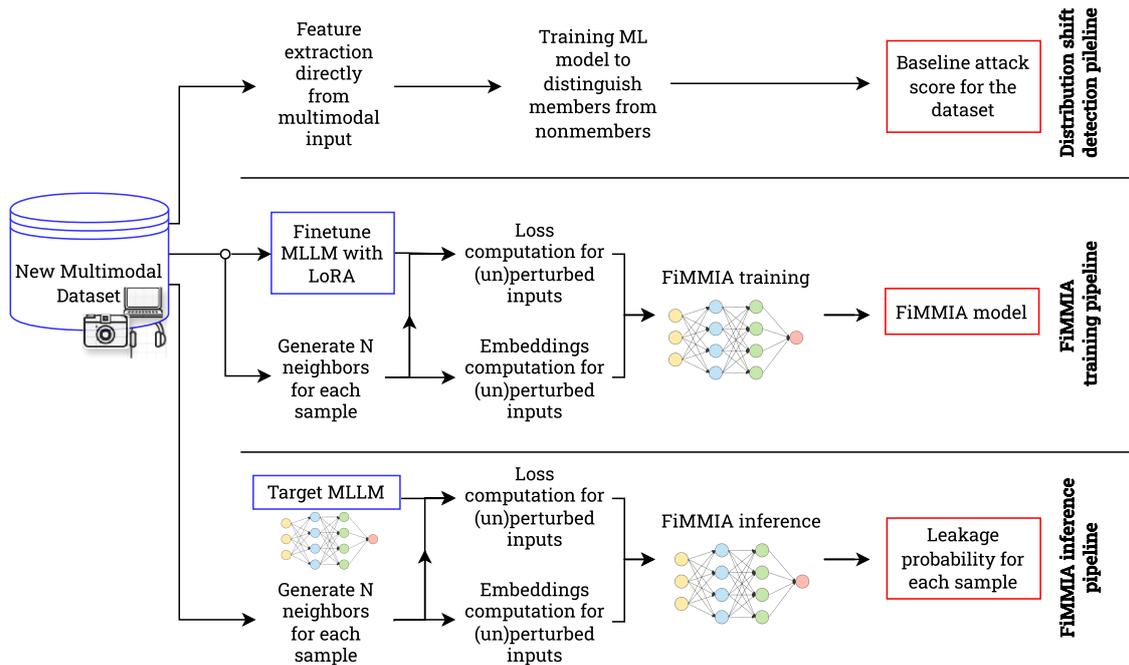


Figure 1: Overview of FiMMIA Inference pipeline for MLLMs. Inputs to the pipeline are shown in blue. Outputs of the pipeline are shown in red.

ent modalities: text, image, audio and video and is fully open source<sup>3</sup>. In order to allow for membership inference in cases, when only the text or multimodal part is assumed to be leaked, we support separate neighbor and embedding generation for both parts of the input, thus providing an option to disentangle their contribution to the final MIA score.

We release pretrained FiMMIA models to promote experiments within the community<sup>4</sup>. Although in our experiments we focus on MERA datasets (Chervyakov et al., 2025) to ensure independence in the split between members and non-members, *the presented pipeline is built with the idea of supporting modular extension and is intended to be easily adopted.*

Overall, the system is a set of models and Python scripts in a GitHub repository that supports three major functionalities: 1) a baseline attack based on distribution statistics, intended to ensure the reliability of multimodal MIA baselines; 2) inference scripts for the FiMMIA model; 3) a training pipeline for new datasets. Main system components are shown at Figure 1. We describe the gen-

<sup>3</sup>[https://github.com/ai-forever/data\\_leakage\\_detect](https://github.com/ai-forever/data_leakage_detect)

<sup>4</sup><https://huggingface.co/collections/ai-forever/fimmia>

eral pipeline for multimodal MIA in subsection 3.4.

### 3.2 Multimodal membership inference benchmarks suffer from distribution shifts

Recently, (Das et al., 2024) have evaluated common textual membership inference benchmarks using blind statistical methods, and have found that they suffer from distribution shifts, with baseline methods independent of any target model’s output outperforming best membership inference attacks on these datasets. An introduction of embedding model into the pipeline (Mozaffari and Marathe, 2024; Hu et al., 2022c) obviously makes the matter even worse, as they shine in tasks related to the separation of different distributions. This fact has, e.g. been recently utilized by (Miyamoto et al., 2025), who have also acknowledged the problem, and used a DINO-V2 (Oquab et al., 2023) to extract image features to show that VL-MIA member and non-member data suffer from a distributional mismatch introduced by the generative nature of non-member samples with an AUC-ROC of **94.9%** using their method. There are reasons for us to argue against this approach. Foremost, the usage of advanced deep learning model still poses threats alike the ones outlined above. Thus, we extend the work of (Das et al., 2024) to multimodal data and, to our surprise, find that attacks that directly

use features obtained from the dataset samples in absence of any information from the target model outperform best known attacks on most multimodal MIA benchmarks.

### 3.3 Distribution shift detection & baseline attacks

Essentially, for each input sample from the dataset with specified members and non-members we extract common heuristic (e.g. SIFT, LBP histogram) or spectral features, and them as inputs to a shallow ML model (e.g. logistic regression or gradient boosting)<sup>5</sup>. The model is trained on 5-fold cross-validation splits with the final attack score for each dataset taken as an average of ones obtained across folds. We assume that if both members and non-members come from the same distribution, i.e. the assumption of i.i.d. samples is valid, then this type of attacks should fail, showing AUC-ROC around 50%. Otherwise, if data collection method was biased (e.g. due to temporal differences, different data generation processes or other factors), these baseline attacks should serve as a lower bound for the proposed membership inference approaches.

We evaluated recently proposed MIA benchmarks in text (Fu et al., 2025; Hallinan et al., 2025) and image (Li et al., 2024c) modalities using the proposed method, and found that most of them suffer from severe distribution shifts, making them hardly useful to evaluate MIAs, with only LAION-MI (Dubiński et al., 2023) being mostly unaffected. See Table 1. Thus, in order to ensure credible results, we aim to use random splits of recently open-sourced multimodal datasets for Russian language (Chervyakov et al., 2025) in our further experiments. Although we are unaware of any common MIA benchmarks for audio or video data, we release both image and audio pipelines and encourage the community to use them prior to the release of new MIA datasets.

### 3.4 Methodology

Membership inference attacks (MIAs) against LLMs aim to determine whether given a target model  $\mathcal{M}$  and a given data point was part of the training dataset used to train the target model. Given a multimodal sample  $x = (t, s)$  from the dataset  $D \sim \mathcal{P}(\mathcal{T} \times \mathcal{S})$  where  $s \in \mathcal{S}$  is some modality (image/video/audio),  $t \in \mathcal{T}$  is the text,

<sup>5</sup>Details on the design of distribution shift detection pipeline and features extracted are available at A.6

estimate  $\mathbb{P}(x \in D|\mathcal{M})$ , probability that a target model was trained on  $x$ .

In accordance with the original article (Mozafari and Marathe, 2024), we divided the training algorithm into the following subsequent steps with some modifications:

1. Neighbor generation
2. Embedding generation
3. Loss computation
4. Training the attack model

#### 3.4.1 Neighbor and embedding generation

For each original data point  $(t, s)$  we generate  $K = 24$  perturbed "neighbors"  $(t_r^k, s_r^k)$ . Recently, there have been increasing attempts to link adversarial theory of neural networks to membership inference, arguing for the special local properties of the loss function in the neighborhood of each input (Xue et al., 2025; Ali et al., 2023). However, there have been several reasons for us to refrain from this approach: generating adversarial examples in discrete domains faces challenges due to non-differentiability (Yang et al., 2020) and generally necessitates to assume a white-box access to the target model, which was against our design principles. Moreover, recently (Gupta et al., 2025) have shown that adversarial examples for MMLMs are not generally transferable, which would additionally limit the applicability of our framework and its transferability across models. Instead, we've performed 4 different structured perturbations to the untokenized input string  $t$ :

1. Random masking and sampling masked words with Fred-T5 model<sup>6</sup>
2. Removing random words
3. Duplication of random words
4. Swapping random words

Each technique is applied to the each text sample  $t$  6 times, resulting in totally 24 "neighbors" per sample. Although, in our experiments we fix  $s = s_r^k, \forall s \in D$ , so the modality data remains unchanged, the pipeline can be modified to support neighbors from different modalities as well.

Then for each original text  $t$  and its neighbors  $t_r^k$  we extract their text embeddings using a fixed encoder:

$$e = \mathcal{E}(t), \quad e'_k = \mathcal{E}(t_r^k)$$

where  $\mathcal{E}$  is `intfloat/e5-mistral-7b-instruct`<sup>7</sup>.

<sup>6</sup>ai-forever/FRED-T5-1.7B, (Zmitrovich et al., 2024)

<sup>7</sup>intfloat/e5-mistral-7b-instruct in our experiments. It used

### 3.4.2 Loss computation

We compute the multimodal loss for both models  $\mathcal{M}$  and  $\mathcal{M}_{leak}$  on both the original and neighbor data points:

$$\mathcal{L} = \mathcal{L}(\mathcal{M}, t, s), \quad \mathcal{L}'_k = \mathcal{L}(\mathcal{M}, t'_i, s'_i)$$

Text input  $t$  is provided to each model, accompanied by the corresponding modality  $s$  (image, video, or audio data in its original, unchanged form).

### 3.4.3 Attack model training

The core of FiMMIA is a binary neural network classifier trained to distinguish between models that have and have not seen the data. For each neighbor  $k$  we create two training examples by computing feature differences<sup>8</sup>:

$$\Delta\mathcal{L} = \mathcal{L} - \mathcal{L}'_k, \quad \Delta e = e - e'_k$$

These feature vectors are paired with labels  $y \in \{0, 1\}$  indicating whether the losses came from  $\mathcal{M}$  (non-leaked) or  $\mathcal{M}_{leak}$  (leaked). However, absolute values of these statistics may vary across datasets and models. To make the system more stable, we apply the z-score normalization technique (Wikipedia, 2025). The values mean  $\mu$  and standard deviation  $\sigma$  of the models' loss differences  $\Delta\mathcal{L}$ , used to normalize input features during training and evaluation are obtained from disjoint train/test splits to mimic real-world scenarios.

$$\Delta\mathcal{L}_{norm} = \frac{\Delta\mathcal{L} - \mu}{\sigma}$$

This process yields random batch training triplets  $(\Delta\mathcal{L}_{norm}, \Delta e, y)$  per original data point. The FiMMIA detector,  $f_{FiMMIA}$  is trained to predict the probability  $p = f_{FiMMIA}(\Delta\mathcal{L}_{norm}, \Delta e)$  that the input features originate from a model that has been trained on the target data. We provide the details of the architecture for FiMMIA model in subsection A.1 and the hyperparameters for training the FiMMIA model in subsection A.2.

It should be noted, that although we suppose a grey-box access to the MLLM in our experiments,

to be SoTA on the MTEB benchmark (Muennighoff et al., 2022) at the time of the model experiments

<sup>8</sup>Similar ideas has been already explored e.g. in (He et al., 2024) where the authors explored both utilizing shadow models and perturbed datasets as calibration data, and found that they are, to a large degree, interchangeable. The idea of using embedding differences as a proxy for difficulty calibration serves as another intuition for our method.

i.e. an attacker has full access to the model's logprobs for loss computation, our setup can be extended to the black-box scenario in presence of compatible APIs, with e.g. only top-k logprobs being released, using approaches from (Finlayson et al., 2024; Bao et al., 2025). We plan to implement such functionality in future releases.

### 3.4.4 Inference

To infer if a target model  $\mathcal{M}'$  has been trained on a specific data point  $(t, s)$ , we compute the loss and embedding differences for this model. We then compute the leakage score  $A$  for the data point by taking the average probability output by the detector over all  $K$  neighbors:

$$A(t, m) = \frac{1}{K} \sum_{k=1}^K f_{FiMMIA}(\Delta\mathcal{L}_{norm}^k, \Delta e^k)$$

## 4 Experiment setup

### 4.1 Data

We evaluate our method on the MERA benchmark (Chervyakov et al., 2025), which comprises 18 audio, video, and image datasets. All tasks in the benchmark are multimodal, taking both a modality input and an instruction, and requiring a text output in a constrained format (e.g., multiple-choice or short-answer). For training phase we fine-tune MLLM  $\mathcal{M}_{leak}$  on each modality separately. Each sample in the training data for the MLLM can be represented as  $x = (s, q, a)$ , a concatenation of the question and the answer as the textual part  $t$ , along with the multimodal input  $s$  (image, video, or audio). In order to ensure credible evaluation of FiMMIA model we split each dataset into train and test parts randomly. The size of the test part is 10% of original dataset. Normalization parameters  $\mu_{D, \mathcal{M}}$  and  $\sigma_{D, \mathcal{M}}$  are calculated from the train part of each of the splitted datasets for each model.

The detailed overview of the benchmark is presented in Table 2.

### 4.2 Models

We evaluate 9 publicly available multimodal models from the most trending model families on HuggingFace, varying in size from 3B to 12B parameters. See Appendix A.3 for detailed model descriptions.

### 4.3 Cross-lingual transfer

This section presents our experimental evaluation, extending the pipeline to English image datasets

	Dataset / task	Size	Answer
audio	ruEnvAQA	596	MC
	RuSLUn	741	OE
	*AQUARIA	738	MC
	*ruTiE-Audio	1500	MC
image	ruCLEVR	1148	OE
	ruCommonVQA	3015	OE
	ruNaturalScienceVQA	363	MC
	WEIRD	814	MC
	*LabTabVQA	339	MC
	*RealVQA	773	OE
	*ruHHH-Image	595	MC
	*ruMathVQA	502	OE
	*ruTiE-Image	1500	MC
	*SchoolScienceVQA	4227	MC
	*UniScienceVQA	7432	OE
video	CommonVideoQA	907	MC
	*RealVideoQA	671	MC
	*ruHHH-Video	911	MC

Table 2: Overview of datasets in MERA benchmark. Those marked with an asterisk were collected from scratch by Chervyakov et al. (2025), while the others are *public datasets* compiled from open-source datasets. **Size** column shows the number of samples in the dataset, and **Answer** column is the task format (MC and OE stand for multiple-choice and open-ended, respectively).

and models. Following the paper by (Song et al., 2025), our analysis leverages two multi-choice datasets: ScienceQA (Lu et al., 2022) and MM-Star (Chen et al., 2024), along with caption dataset: COCO-Caption2017 (Lin et al., 2015). We randomly selected 2000 samples from ScienceQA’s test set, respectively, with 1000 samples from the other datasets. We select Qwen2.5-VL-3B-Instruct as a target fine-tuned MLLM and train FiMMIA as described in section subsection 3.4 only on MERA benchmark (Chervyakov et al., 2025) without fine-tuning or using any English data. We evaluate 4 publicly available multimodal models similar to the paper (Song et al., 2025) that presents MM-DETECT method (see Table 9 for model descriptions). That method calculates  $\Delta$  score for the dataset and if  $\Delta < 0$ , dataset leakage is presumed. In order to make a comparison with this method we calculate % of leaked samples from the dataset, guided by our pipeline.

## 5 Results

We report AUC-ROC for binary classification (leaked vs. clean) as shown in Tables 3, 5, 4. Also we report TPR with low FPR in Tables 12, 10, 11. In order to evaluate the transferability of the trained attack model we also report scores when the origin and test models differ. The  $\mathcal{M}_{\text{origin}}$  is the model used to train FiMMIA, while  $\mathcal{M}_{\text{test}}$  is the model whose losses are used to test FiMMIA

$\mathcal{M}_{\text{origin}}$	$\mathcal{M}_{\text{test}}$	AUC
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-3B-Instruct	<b>96.2</b>
Qwen2.5-VL-3B-Instruct	Qwen2-VL-7B-Instruct	86.0
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-7B-Instruct	88.0
Qwen2.5-VL-3B-Instruct	Llava-Next-8b-hf	<b>90.2</b>
Qwen2.5-VL-3B-Instruct	Gemma-3-4B-it	65.8
Qwen2.5-VL-3B-Instruct	Gemma-3-12b-it	67.9
Qwen2-VL-7B-Instruct	Qwen2.5-VL-3B-Instruct	78.0
Qwen2-VL-7B-Instruct	Qwen2-VL-7B-Instruct	<b>96.2</b>
Qwen2-VL-7B-Instruct	Qwen2.5-VL-7B-Instruct	<b>80.5</b>
Qwen2-VL-7B-Instruct	Llama3-llava-next-8b-hf	78.0
Qwen2-VL-7B-Instruct	Gemma-3-4b-it	77.7
Qwen2-VL-7B-Instruct	Gemma-3-12b-it	73.7
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-3B-Instruct	92.8
Qwen2.5-VL-7B-Instruct	Qwen2-VL-7B-Instruct	93.1
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-7B-Instruct	<b>98.1</b>
Qwen2.5-VL-7B-Instruct	Llama3-llava-next-8b-hf	95.8
Qwen2.5-VL-7B-Instruct	Gemma-3-4b-it	95.4
Qwen2.5-VL-7B-Instruct	Gemma-3-12b-it	94.5
Llama3-llava-next-8b-hf	Qwen2.5-VL-3B-Instruct	94.6
Llama3-llava-next-8b-hf	Qwen2-VL-7B-Instruct	90.0
Llama3-llava-next-8b-hf	Qwen2.5-VL-7B-Instruct	96.6
Llama3-llava-next-8b-hf	Llama3-llava-next-8b-hf	97.7
Llama3-llava-next-8b-hf	Gemma-3-4b-it	99.1
Llama3-llava-next-8b-hf	Gemma-3-12b-it	<b>99.5</b>
Gemma-3-4b-it	Qwen2.5-VL-3B-Instruct	76.0
Gemma-3-4b-it	Qwen2-VL-7B-Instruct	71.5
Gemma-3-4b-it	Qwen2.5-VL-7B-Instruct	85.2
Gemma-3-4b-it	Llama3-llava-next-8b-hf	86.5
Gemma-3-4b-it	Gemma-3-4b-it	<b>99.4</b>
Gemma-3-4b-it	Gemma-3-12b-it	98.7
Gemma-3-12b-it	Qwen2.5-VL-3B-Instruct	84.1
Gemma-3-12b-it	Qwen2-VL-7B-Instruct	81.3
Gemma-3-12b-it	Qwen2.5-VL-7B-Instruct	91.2
Gemma-3-12b-it	Llama3-llava-next-8b-hf	93.3
Gemma-3-12b-it	Gemma-3-4b-it	99.4
Gemma-3-12b-it	Gemma-3-12b-it	<b>99.7</b>

Table 3: AUC-ROC FiMMIA performance metrics for various evaluated **Image** MLLMs.

(predict whether the data sample was used to train  $\mathcal{M}_{\text{test}}$  or not).

$\mathcal{M}_{\text{origin}}$	$\mathcal{M}_{\text{test}}$	AUC
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-3B-Instruct	95.9
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-7B-Instruct	<b>99.5</b>
Qwen2.5-VL-3B-Instruct	LLaVA-NeXT-Video	91.7
Qwen2.5-VL-3B-Instruct	LLaVA-NeXT-Video-DPO	91.2
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-3B-Instruct	98.7
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-7B-Instruct	<b>100.0</b>
Qwen2.5-VL-7B-Instruct	LLaVA-NeXT-Video	96.5
Qwen2.5-VL-7B-Instruct	LLaVA-NeXT-Video-DPO	95.7
LLaVA-NeXT-Video	Qwen2.5-VL-3B-Instruct	63.7
LLaVA-NeXT-Video	Qwen2.5-VL-7B-Instruct	71.5
LLaVA-NeXT-Video	LLaVA-NeXT-Video	<b>100.0</b>
LLaVA-NeXT-Video	LLaVA-NeXT-Video-DPO	<b>100.0</b>
LLaVA-NeXT-Video-DPO	Qwen2.5-VL-3B-Instruct	53.6
LLaVA-NeXT-Video-DPO	Qwen2.5-VL-7B-Instruct	56.2
LLaVA-NeXT-Video-DPO	LLaVA-NeXT-Video	<b>100.0</b>
LLaVA-NeXT-Video-DPO	LLaVA-NeXT-Video-DPO	<b>100.0</b>

Table 4: AUC-ROC FiMMIA performance metrics for various evaluated **Video** MLLMs.

Overall, the results of the FiMMIA detection capabilities are presented in Table 6. All models show significant success within their own family; however, the success of the attack may decrease when testing on a model from a different family. Nev-

$\mathcal{M}_{\text{origin}}$	$\mathcal{M}_{\text{test}}$	AUC
Qwen2-Audio-7B-Instruct	Qwen2-Audio-7B-Instruct	<b>87.7</b>
Qwen2-Audio-7B-Instruct	Qwen-Audio-Chat	76.0
Qwen-Audio-Chat	Qwen2-Audio-7B-Instruct	61.3
Qwen-Audio-Chat	Qwen-Audio-Chat	<b>100.0</b>

Table 5: AUC-ROC FiMMIA performance metrics for various evaluated **Audio** MLLMs.

ertheless, the metric score for each experiment exceeds 65.0, which indicates the promising transferability of the proposed method. Moreover, average metrics for each modality are quite high, ranging from 80 to 90% AUC-ROC.

Modality	AUC
Image	88.658
Video	88.388
Audio	81.250

Table 6: Average AUC-ROC of FiMMIA per modality. Averaging over the models used for training and evaluating FiMMIA.

Evaluations on the transferability of the model to a different language inputs are presented in Table 7. The results indicate that our method is almost entirely in agreement with those presented in the paper (Song et al., 2025). If  $\Delta < 0$  the amount of samples predicted by FiMMIA as leaked is more than 0.1 in most cases, which corresponds to at least 10% of the dataset. However, if the task allows, we suggest to train FiMMIA for particular dataset and language from scratch to obtain more accurate and reliable results.

Dataset	Model	FiMMIA	MM-DETECT $\Delta$
COCO	Phi-3-vision-128k-instruct	0.00	0.5
	Qwen-VL-Chat	0.00	-1.9
	LLaVA-1.5-7B	0.58	-0.6
	fuyu-8b	0.22	1.0
MMStar	Phi-3-vision-128k-instruct	0.06	3.2
	Qwen-VL-Chat	0.00	3.3
	LLaVA-1.5-7B	0.13	2.8
	fuyu-8b	0.011	-1.2
ScienceQA	Phi-3-vision-128k-instruct	0.10	0.7
	Qwen-VL-Chat	0.00	0.1
	LLaVA-1.5-7B	0.21	1.3
	fuyu-8b	0.19	-0.5

Table 7: Comparison FiMMIA % leakage samples detected of MLLMs on English datasets with MM-DETECT score for image modality.

## 6 Conclusion

This paper introduces FiMMIA, a novel framework that leverages input semantics and strategic perturbations to train a highly effective neural network

for data leakage detection in MLLMs. Our key contribution is a language-agnostic system capable of training robust leakage detection models for any dataset. Designed for extensibility, the framework natively supports neighbor generation across multiple modalities paving the way for future research.

## Limitations

**Scope of the Method** When training FiMMIA, we only target a fine-tuning scenario for the MLLM using a low-rank adapter. The results for pretraining and full fine-tuning may be different due to the capacity scaling laws (Morris et al., 2025), and other factors. We leave these evaluations for further work.

**Determinism and Reproducibility** Even our fine-tuned models’ losses are subject to stochasticity, as the entire hardware–software stack affects inference: GPU model, drivers/CUDA/cuDNN, PyTorch, vLLM/transformers (and commit hashes), flash-attention kernels, tokenizers/checkpoints, precision/quantization, and batching – some of which are non-deterministic or can vary between environments. However, in general, the variance that these factors contribute to evaluation metrics is not substantial.

**Speed and Computational Complexity** In our experiments the inference process took approx. 10 hours on a single GPU for one dataset. Generally, the time complexity of our algorithm scales as  $\mathcal{O}(|D|N(M + E + G))$ , where  $|D|$  is the number of samples in the dataset,  $N$  is the number of neighbors, and  $M, E, G$  are time complexities of the target, embedding and neighbor generation models.

**Model Assumption Dependencies** The method relies on per-sample loss access (a gray-box assumption) and depends on an external model for generating embeddings. The applicability of the method in a strict black-box setting, where such access is unavailable, is not addressed in this work, despite the existence of relevant prior research.

## Ethical consideration

**Use of Public Data** All experiments and evaluations in this study rely exclusively on openly accessible public datasets. No proprietary, confidential, or otherwise sensitive information was involved. This choice supports transparency, facilitates inde-

pendent verification, and avoids any infringement on data-privacy protections.

**Defensive and Constructive Purpose** Our work reconceptualizes membership-inference analysis as a diagnostic and privacy-protecting tool rather than a privacy-threat vector. The method is designed to:

- By identifying cases in which benchmark samples have been inadvertently memorized during training, the approach helps prevent benchmark saturation and dataset contamination, thereby supporting fair and meaningful model comparison.
- The technique offers researchers a practical mechanism for auditing training pipelines to ensure that performance improvements stem from genuine advances rather than overfitting to widely used evaluation sets.
- As competitive leaderboard dynamics can unintentionally encourage data leakage and undermine the long-term value of public benchmarks, our framework contributes to more resilient evaluation standards that promote steady, reliable scientific progress.

## Acknowledgments

The authors would like to express their sincere gratitude to Dmitry Gorbetsky, Yaroslav Grebnyak, Oleg Yangalichin and Artem Chervyakov for their valuable contributions and support in this work.

## References

Marah I Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahua Xu,

Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *CoRR*, abs/2404.14219.

Hassan Ali, Adnan Qayyum, Ala Al-Fuqaha, and Junaid Qadir. 2023. Membership inference attacks on dnns using adversarial perturbations. *arXiv preprint arXiv:2307.05193*.

Junxing Bai, Shanshan Bai, Shiqi Yang, Shi Wang, Shoujie Tan, Panpan Wang, Jiawei Lin, Chaozhe Zhou, and Junrui Zhou. 2023. [Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond](#). *Preprint*, arXiv:2308.12966.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. [Qwen2.5-vl technical report](#). *Preprint*, arXiv:2502.13923.

Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondrej Dusek. 2024. [Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93, St. Julian’s, Malta. Association for Computational Linguistics.

Guangsheng Bao, Yanbin Zhao, Juncai He, and Yue Zhang. 2025. [Glimpse: Enabling white-box methods to use proprietary models for zero-shot llm-generated text detection](#). *Preprint*, arXiv:2412.11506.

Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. 2022. [Membership inference attacks from first principles](#). *Preprint*, arXiv:2112.03570.

Dingfan Chen, Ning Yu, Yang Zhang, and Mario Fritz. 2020. GAN-Leaks: A taxonomy of membership inference attacks against generative models. In *Proceedings of the 2020 ACM SIGSAC conference on computer and communications security*, pages 343–362.

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and et al. 2024. [Are we on the right way for evaluating large vision-language models?](#) *arXiv preprint arXiv:2403.20330*.

Zequan Chen, Jiannan Wu, Wenqian Wang, Weijiang Su, Guangda Chen, Shen Xing, Mingyang Zhong, Qing Zhang, Xin Zhu, Lei Lu, Bo Li, Peihao Luo, Tong Lu, Yi Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning

- for generic visuo-linguistic tasks. *arXiv preprint arXiv:2312.14238*.
- Artem Chervyakov, Ulyana Isaeva, Anton Emelyanov, Artem Safin, Maria Tikhonova, Alexander Kharitonov, Yulia Lyakh, Petr Surovtsev, Denis Shevelev, Vildan Saburov, Vasily Kononov, Elisei Rykov, Ivan Sviridov, Amina Miftakhova, Ilseyar Alimova, Alexander Panchenko, Alexander Kapitanov, and Alena Fenogenova. 2025. [Multimodal evaluation of Russian-language architectures](#). *Preprint*, arXiv:2511.15552.
- Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. 2021. Label-only membership inference attacks. In *International conference on machine learning*, pages 1964–1974. PMLR.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. 2024. [Qwen2-audio technical report](#). *Preprint*, arXiv:2407.10759.
- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*.
- Debeshee Das, Jie Zhang, and Florian Tramèr. 2024. Blind baselines beat membership inference attacks for foundation models. *arXiv preprint arXiv:2406.16201*.
- Jan Dubiński, Antoni Kowalczyk, Stanisław Pawlak, Przemysław Rokita, Tomasz Trzcifski, and Paweł Morawiecki. 2023. [Towards more realistic membership inference attacks on large diffusion models](#). *Preprint*, arXiv:2306.12983.
- Matthew Finlayson, Xiang Ren, and Swabha Swayamdipta. 2024. Logits of api-protected llms leak proprietary information. *arXiv preprint arXiv:2403.09539*.
- Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. 2023. Practical membership inference attacks against fine-tuned large language models via self-prompt calibration. *arXiv preprint arXiv:2311.06062*.
- Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. 2025. MIA-tuner: Adapting large language models as pre-training text detector. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Philadelphia, Pennsylvania, USA.
- Isha Gupta, Rylan Schaeffer, Joshua Kazdan, Ken Ziyu Liu, and Sanmi Koyejo. 2025. Understanding adversarial transfer: Why representation-space attacks fail where data-space attacks succeed. *arXiv preprint arXiv:2510.01494*.
- Skyler Hallinan, Jaehun Jung, Melanie Sclar, Ximing Lu, Abhilasha Ravichander, Sahana Ramnath, Yejin Choi, Sai Praneeth Karimireddy, Niloofar Mireshghallah, and Xiang Ren. 2025. The surprising effectiveness of membership inference with simple n-gram coverage. *arXiv preprint arXiv:2508.09603*.
- Jamie Hayes, Luca Melis, George Danezis, and Emiliano De Cristofaro. 2017. Logan: Membership inference attacks against generative models. *arXiv preprint arXiv:1705.07663*.
- Yu He, Boheng Li, Yao Wang, Mengda Yang, Juan Wang, Hongxin Hu, and Xingyu Zhao. 2024. Is difficulty calibration all we need? towards more practical membership inference attacks. *arXiv preprint arXiv:2409.00426*.
- Benjamin Hilprecht, Martin Härterich, and Daniel Bernau. 2019. Monte carlo and reconstruction membership inference attacks against generative models. *Proceedings on Privacy Enhancing Technologies*.
- Haoyang Hu, Zoran Salcic, Lingyu Sun, Gillian Dobbie, Philip S. Yu, and Xinhui Zhang. 2022a. Membership inference attacks on machine learning: A survey. *ACM Comput. Surv.*, 54(11s):1–37.
- Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. 2022b. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37.
- Pingyi Hu, Zihan Wang, Ruoxi Sun, Hu Wang, and Minhui Xue. 2022c. M<sup>4</sup>i: Multi-modal models membership inference. *arXiv preprint arXiv:2209.06997*.
- Pingyi Hu, Zihan Wang, Ruoxi Sun, Hu Wang, and Minhui Xue. 2022d. M<sup>4</sup>i: Multi-modal models membership inference. In *Advances in Neural Information Processing Systems*, volume 35, pages 1867–1882.
- Myeongseob Ko, Ming Jin, Chenguang Wang, and Ruoxi Jia. 2023. Practical membership inference attacks against large-scale multi-modal models: A pilot study. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4871–4881.
- Ryuto Koike, Liam Dugan, Masahiro Kaneko, Chris Callison-Burch, and Naoaki Okazaki. 2025. [Machine text detectors are membership inference attacks](#). *Preprint*, arXiv:2510.19492.
- Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024a. [Llava-next: Stronger llms supercharge multimodal capabilities in the wild](#).
- Zhan Li, Yongtao Wu, Yihang Chen, Francesco Tonin, Elias Abad Rocamora, and Volkan Cevher. 2024b. [Membership inference attacks against large vision-language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 98645–98674. Curran Associates, Inc.

- Zhan Li, Yongtao Wu, Yihang Chen, Francesco Tonin, Elias Abad Rocamora, and Volkan Cevher. 2024c. Membership inference attacks against large vision-language models. *arXiv preprint arXiv: 2411.02902*.
- Jiawei Lin, Haoqi Yin, Weiran Ping, Yu Lu, Pavlo Molchanov, Andrew Tao, Huiyu Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. 2023. [Vila: On pre-training for visual language models](#). *Preprint*, arXiv:2312.07533.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2015. [Microsoft COCO: Common Objects in Context](#). *CoRR*, abs/1405.0312.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. [Improved baselines with visual instruction tuning](#). *CoRR*, abs/2310.03744.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. [Improved baselines with visual instruction tuning](#). *Preprint*, arXiv:2310.03744.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. [Llava-next: Improved reasoning, ocr, and world knowledge](#).
- Yang Liu, Hongfan Duan, Yan Zhang, Binbin Li, Shaohan Zhang, Wei Zhao, Yonglong Yuan, Jinmao Wang, Chuxiong He, Zhongying Liu, Kechen Chen, and Dahua Lin. 2023b. [Mmbench: Is your multi-modal model an all-around player?](#) *Preprint*, arXiv:2307.06281.
- Yunhui Long, Vincent Bindschaedler, Lei Wang, Diyue Bu, Xiaofeng Wang, Haixu Tang, Carl A Gunter, and Kai Chen. 2018. Understanding membership inferences on well-generalized learning models. *arXiv preprint arXiv:1802.04889*.
- Pengcheng Lu, Sudip Mishra, Tianyi Xia, Leqi Qiu, Kai-Wei Chang, Scott Cheng-Hsin Zhu, Oyvind Tafjord, Peter Clark, and Anirudh Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Saeed Mahloujifar, Huseyin A Inan, Melissa Chase, Esha Ghosh, and Marcello Hasegawa. 2021. Membership inference on word embedding and beyond. *arXiv preprint arXiv:2106.11384*.
- Anqi Mao, Mehryar Mohri, and Yutao Zhong. 2023. [Cross-entropy loss functions: Theoretical analysis and applications](#). *arXiv preprint arXiv:2304.07288*. Published in ICML 2023.
- Justus Mattern, Fatemehsadat Miresghallah, Zhijing Jin, Bernhard Schoelkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. 2023. Membership inference attacks against language models via neighbourhood comparison. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11330–11343.
- Fatemehsadat Miresghallah, Archit Uniyal, Tianhao Wang, David K Evans, and Taylor Berg-Kirkpatrick. 2022. An empirical analysis of memorization in fine-tuned autoregressive language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1816–1826.
- Ryoto Miyamoto, Xin Fan, Fuyuko Kido, Tsuneo Matsumoto, and Hayato Yamana. 2025. [OpenLlm-mia: A controlled benchmark revealing the limits of membership inference attacks on large vision-language models](#). *arXiv preprint arXiv: 2510.16295*.
- Mazda Moayeri, Kiarash Banihashem, and Soheil Feizi. 2022. [Explicit tradeoffs between adversarial and natural distributional robustness](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- John X. Morris, Chawin Sitawarin, Chuan Guo, Narine Kokhlikyan, G. Edward Suh, Alexander M. Rush, Kamalika Chaudhuri, and Saeed Mahloujifar. 2025. How much do language models memorize? *arXiv preprint arXiv: 2505.24832*.
- Hamid Mozaffari and Virendra Marathe. 2024. [Semantic membership inference attack against large language models](#). In *Neurips Safe Generative AI Workshop 2024*.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.
- Vinod Nair and Geoffrey E. Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 807–814. Omnipress.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. 2023. [Dinov2: Learning robust visual features without supervision](#). *arXiv preprint arXiv: 2304.07193*.
- Maxwell Reid, Nikita Savinov, Dmitry Teplyashin, Danil Lepikhin, Timothy Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *arXiv preprint arXiv:2403.05530*.

- Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages 1135–1144. ACM.
- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787, Singapore. Association for Computational Linguistics.
- Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. 2018. MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models. *arXiv preprint arXiv:1806.01246*.
- Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*, volume 80, pages 4596–4604. PMLR.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024. Detecting pretraining data from large language models. In *The Twelfth International Conference on Learning Representations*.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18.
- Congzheng Song and Ananth Raghunathan. 2020. Information leakage in embedding models. In *Proceedings of the 2020 ACM SIGSAC conference on computer and communications security*, pages 377–390.
- Dandan Song, Shuai Chen, Guanhua Chen, Fan Yu, Xiuyue Wan, and Bo Wang. 2024. Milebench: Benchmarking mllms in long context. *arXiv preprint arXiv:2404.18532*.
- Dingjie Song, Sicheng Lai, Mingxuan Wang, Shunian Chen, Lichao Sun, and Benyou Wang. 2025. Both text and images leaked! a systematic analysis of data contamination in multimodal LLM. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 10527–10542, Suzhou, China. Association for Computational Linguistics.
- Liwei Song and Prateek Mittal. 2021. Systematic evaluation of privacy risks of machine learning models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2615–2632.
- Liwei Song, Reza Shokri, and Prateek Mittal. 2019. Membership inference attacks against adversarially robust deep learning models. In *2019 IEEE Security and Privacy Workshops (SPW)*, pages 50–56. IEEE.
- Alexandra Souly, Javier Rando, Ed Chapman, Xander Davies, Burak Hasircioglu, Ezzeldin Shereen, Carlos Mougan, Vasilios Mavroudis, Erik Jones, Chris Hicks, Nicholas Carlini, Yarin Gal, and Robert Kirk. 2025. Poisoning attacks on llms require a near-constant number of poison samples. *arXiv preprint arXiv: 2510.07192*.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Keanealy, Lucas Beyer, Xiaohai Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman, Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi, Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petriani, Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Paparas, Divyashree Shivakumar Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne, Idan Szepktor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wieting, Jonathan Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jyotinder Singh, Kat Black, Kathy Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine, Marina Coelho, Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael Moynihan, Min Ma, Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Nilay Chauhan, Noveen Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Rubenstein, Phil Culliton, Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh Shivanna, Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu, Ryan Mullins, Sammy Jerome, Sara Smoot, Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti Sheth, Siim Pöder, Sijal Bhatnagar, Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi Liu, Trevor Yacovone, Tyler Liechty, Uday Kalra, Utku Evci, Vedant Misra, Vincent Roseberry, Vlad Feinberg, Vlad Kolesnikov, Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein Zhu, Zichuan Wei, Zoltan Egyed, Victor Cotruta,

- Minh Giang, Phoebe Kirk, Anand Rao, Kat Black, Nabila Babar, Jessica Lo, Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas Gonzalez, Zach Gleicher, Tris Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Jean-Baptiste Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier Bachem, Armand Joulin, Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot. 2025. [Gemma 3 technical report](#). *Preprint*, arXiv:2503.19786.
- Teun van der Weij, Felix Hofstätter, Ollie Jaffe, Samuel F. Brown, and Francis Rhys Ward. 2024. Ai sandbagging: Language models can strategically underperform on evaluations. *arXiv preprint arXiv:2406.07358*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](#). *Preprint*, arXiv:2409.12191.
- Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Benjamin Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh Sandha, Siddhartha Venkat Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. 2025. Livebench: A challenging, contamination-free LLM benchmark. In *The Thirteenth International Conference on Learning Representations*.
- Wikipedia. 2025. Standard score — Wikipedia, The Free Encyclopedia. [https://en.wikipedia.org/wiki/Standard\\_score](https://en.wikipedia.org/wiki/Standard_score). [Online; accessed 17-November-2025].
- Rui Xu, Ze Wang, Ren-Zhang Fan, and Pengfei Liu. 2024. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*.
- Jing Xue, Zhishen Sun, Haishan Ye, Luo Luo, Xiangyu Chang, Ivor Tsang, and Guang Dai. 2025. Privacy leaks by adversaries: Adversarial iterations for membership inference attack. *arXiv preprint arXiv:2506.02711*.
- Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael I Jordan. 2020. Greedy attack and gumbel attack: Generating adversarial examples for discrete data. *Journal of Machine Learning Research*, 21(43):1–36.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. [Privacy risk in machine learning: Analyzing the connection to overfitting](#). In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pages 268–282.
- Jinhua Yin, Peiru Yang, Chen Yang, Huili Wang, Zhiyang Hu, Shangguang Wang, Yongfeng Huang, and Tao Qi. 2025. Black-box membership inference attack for llms via prior knowledge-calibrated memory probing. *arXiv preprint arXiv: 2511.01952*.
- Dmitry Zmitrovich, Aleksandr Abramov, Andrey Kalmykov, Vitaly Kadulin, Maria Tikhonova, Ekaterina Taktasheva, Danil Astafurov, Mark Baushenko, Artem Snegirev, Tatiana Shavrina, Sergei S. Markov, Vladislav Mikhailov, and Alena Fenogenova. 2024. [A family of pretrained transformer language models for Russian](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 507–524, Torino, Italia. ELRA and ICCL.

## A Appendix

### A.1 Attack model neural network architecture

The detailed architecture of the FiMMIA is provided below.

#### 1. Input Data:

- `loss_input`: A tensor fed into the `loss_component`.
- `embedding_input`: A tensor fed into the `embedding_component`.

#### 2. `loss_component`:

- A Linear layer: 1 input feature  $\rightarrow$  `projection_size` output features.
- Dropout(0.2) and ReLU (Nair and Hinton, 2010) activation.

#### 3. `embedding_component`:

- A Linear layer: `embedding_size`  $\rightarrow$  `embedding_size // 2`.
- Dropout(0.2) and ReLU (Nair and Hinton, 2010) activation.
- A Linear layer: `embedding_size // 2`  $\rightarrow$  512.
- Dropout(0.2) and ReLU (Nair and Hinton, 2010) activation.

#### 4. Concatenation (`torch.hstack`):

- The outputs from the `loss_component` (`projection_size`) and the `embedding_component`(512) are concatenated into a single vector of size  $2 * \text{projection\_size}$ .

#### 5. `attack_encoding`:

- A series of 6 fully connected Linear layers with Dropout(0.2) and ReLU (Nair and Hinton, 2010) activations between them:  $2 * \text{projection\_size} \rightarrow 512 \rightarrow 256 \rightarrow 128 \rightarrow 64 \rightarrow 32$ .
- The final Linear layer:  $32 \rightarrow 2$  (output logits for classification).

## 6. Output:

- The model returns the logits (size 2).
- If labels are provided, it also calculates and returns the cross-entropy loss (Mao et al., 2023).

## A.2 Attack model hyperparameters

To construct the neighbor datasets, we generate  $k = 24$  neighbors for each data point. We employ the adafactor optimizer (Shazeer and Stern, 2018) to train the network on our training data over 10 epochs. The batch size is set to 64, meaning each batch contains random triplets. For experiments, we use a learning rate of  $2 \times 10^{-6}$ .

## A.3 Models Details

Table 8 contains information about multimodal LLMs used for the experiments. As the number of MLLMs trained with a focus on russian is limited, we evaluate our method using known open-source models. Although it may contribute to higher ROC-AUC scores we observe in our experiments due to the models being adapted to vastly new domain, it also helps us alleviate possible effects related to the possibility of our evaluation datasets' traces being already present in models' training data.

## A.4 English Models Details

Table 9 contains information about multimodal LLMs used for the language transfer experiments. All models are selected from the following paper (Song et al., 2025).

## A.5 TPR at low FPR (FPR=5%) results

Here we report the True Positive Rate (TPR) at a low False Positive Rate (FPR), which measures the detection rate at a meaningful threshold. The modality of image is presented in Table 12, the video in Table 10 and the audio accordingly in Table 11.

## A.6 Description of the distribution shift detection pipelines

For the information on the features extracted from image and audio data see Table 13.

Model	Parameters	Context length	Hugging Face Hub link	Citation
Qwen2-VL-7B-Instruct	7B	32K	<a href="#">Qwen/Qwen2-VL-7B-Instruct</a>	Wang et al. (2024)
Qwen2.5-VL-3B-Instruct	3B	128K	<a href="#">Qwen/Qwen2.5-VL-3B-Instruct</a>	Bai et al. (2025)
Qwen2.5-VL-7B-Instruct	7B	128K	<a href="#">Qwen/Qwen2.5-VL-7B-Instruct</a>	
gemma-3-4b-it	4B	128K	<a href="#">google/gemma-3-4b-it</a>	Team et al. (2025)
gemma-3-12b-it	12B	128K	<a href="#">google/gemma-3-12b-it</a>	
llama3-llava-next-8b-hf	8B	128K	<a href="#">llava-hf/llama3-llava-next-8b-hf</a>	Li et al. (2024a)
LLaVA-NeXT-Video	7B	4K	<a href="#">llava-hf/LLaVA-NeXT-Video-7B-hf</a>	Liu et al. (2024b)
LLaVA-NeXT-Video-DPO	7B	4K	<a href="#">llava-hf/LLaVA-NeXT-Video-7B-DPO-hf</a>	
Qwen2-Audio-7B-Instruct	7B	32K	<a href="#">Qwen/Qwen2-Audio-7B-Instruct</a>	Chu et al. (2024)
Qwen/Qwen-Audio-Chat	7B	32K	<a href="#">Qwen/Qwen-Audio-Chat</a>	Chu et al. (2023)

Table 8: General information about used multimodal LLMs for experiments.

Model	Parameters	Context length	Hugging Face Hub link	Citation
Phi-3-vision-128k-instruct	8B	128K	<a href="#">microsoft/Phi-3-vision-128k-instruct</a>	(Abdin et al., 2024)
LLaVA-1.5-7B	7B	16K	<a href="#">llava-hf/llava-1.5-7b-hf</a>	(Liu et al., 2024a)
Qwen-VL-Chat	7B	8K	<a href="#">Qwen-VL-Chat</a>	(Bai et al., 2023)
fuyu-8b <sup>9</sup>	8B	16K	<a href="#">adept/fuyu-8b</a>	

Table 9: General information about used multimodal LLMs used for the language transfer experiments.

$\mathcal{M}_{\text{origin}}$	$\mathcal{M}_{\text{test}}$	AUC	TPR
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-3B-Instruct	95.9	85.8
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-7B-Instruct	99.5	98.4
Qwen2.5-VL-3B-Instruct	LLaVA-NeXT-Video	91.7	52.9
Qwen2.5-VL-3B-Instruct	LLaVA-NeXT-Video-DPO	91.2	62.9
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-3B-Instruct	98.7	95.4
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-7B-Instruct	100.0	100.0
Qwen2.5-VL-7B-Instruct	LLaVA-NeXT-Video	96.5	80.8
Qwen2.5-VL-7B-Instruct	LLaVA-NeXT-Video-7B-DPO	95.7	82.1
LLaVA-NeXT-Video	Qwen2.5-VL-3B-Instruct	63.7	6.0
LLaVA-NeXT-Video	Qwen2.5-VL-7B-Instruct	71.5	70.0
LLaVA-NeXT-Video	LLaVA-NeXT-Video-7B	100.0	100.0
LLaVA-NeXT-Video	LLaVA-NeXT-Video-7B-DPO	100.0	100.0
LLaVA-NeXT-Video-7B-DPO	Qwen2.5-VL-3B-Instruct	53.6	60.0
LLaVA-NeXT-Video-7B-DPO	Qwen2.5-VL-7B-Instruct	56.2	43.0
LLaVA-NeXT-Video-7B-DPO	LLaVA-NeXT-Video-7B	100.0	100.0
LLaVA-NeXT-Video-7B-DPO	LLaVA-NeXT-Video-7B-DPO	100.0	100.0

Table 10: AUC-ROC and TPR at low FPR (FPR=5%) FiMMIA performance metrics for various evaluated Video MLLMs.

$\mathcal{M}_{\text{origin}}$	$\mathcal{M}_{\text{test}}$	AUC	TPR
Qwen2-Audio-7B-Instruct	Qwen2-Audio-7B-Instruct	87.7	61.9
Qwen2-Audio-7B-Instruct	Qwen-Audio-Chat	76.0	74.5
Qwen-Audio-Chat	Qwen2-Audio-7B-Instruct	61.3	62.7
Qwen-Audio-Chat	Qwen-Audio-Chat	100.0	100.0

Table 11: AUC-ROC and TPR at low FPR (FPR=5%) FiMMIA performance metrics for various evaluated Audio MLLMs.

$\mathcal{M}_{\text{origin}}$	$\mathcal{M}_{\text{test}}$	AUC	TPR
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-3B-Instruct	96.2	86.1
Qwen2.5-VL-3B-Instruct	Qwen2-VL-7B-Instruct	86.0	39.1
Qwen2.5-VL-3B-Instruct	Qwen2.5-VL-7B-Instruct	88.0	53.0
Qwen2.5-VL-3B-Instruct	llama3-llava-next-8b-hf	90.2	59.9
Qwen2.5-VL-3B-Instruct	gemma-3-4b-it	65.8	6.2
Qwen2.5-VL-3B-Instruct	gemma-3-12b-it	67.9	61.9
Qwen2-VL-7B-Instruct	Qwen2.5-VL-3B-Instruct	78.0	16.5
Qwen2-VL-7B-Instruct	Qwen2-VL-7B-Instruct	96.2	85.1
Qwen2-VL-7B-Instruct	Qwen2.5-VL-7B-Instruct	80.5	35.9
Qwen2-VL-7B-Instruct	llama3-llava-next-8b-hf	78.0	30.6
Qwen2-VL-7B-Instruct	gemma-3-4b-it	77.7	7.2
Qwen2-VL-7B-Instruct	gemma-3-12b-it	73.7	67.8
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-3B-Instruct	92.8	73.8
Qwen2.5-VL-7B-Instruct	Qwen2-VL-7B-Instruct	93.1	77.0
Qwen2.5-VL-7B-Instruct	Qwen2.5-VL-7B-Instruct	98.1	94.0
Qwen2.5-VL-7B-Instruct	llama3-llava-next-8b-hf	95.8	83.1
Qwen2.5-VL-7B-Instruct	gemma-3-4b-it	95.4	71.8
Qwen2.5-VL-7B-Instruct	gemma-3-12b-it	94.5	66.1
llama3-llava-next-8b-hf	Qwen2.5-VL-3B-Instruct	94.6	78.6
llama3-llava-next-8b-hf	Qwen2-VL-7B-Instruct	90.0	65.7
llama3-llava-next-8b-hf	Qwen2.5-VL-7B-Instruct	96.6	90.9
llama3-llava-next-8b-hf	llama3-llava-next-8b-hf	97.7	93.3
llama3-llava-next-8b-hf	gemma-3-4b-it	99.1	98.2
llama3-llava-next-8b-hf	gemma-3-12b-it	99.5	99.6
gemma-3-4b-it	Qwen2.5-VL-3B-Instruct	76.0	20.2
gemma-3-4b-it	Qwen2-VL-7B-Instruct	71.5	19.6
gemma-3-4b-it	Qwen2.5-VL-7B-Instruct	85.2	42.7
gemma-3-4b-it	llama3-llava-next-8b-hf	86.5	41.7
gemma-3-4b-it	gemma-3-4b-it	99.4	98.0
gemma-3-4b-it	gemma-3-12b-it	98.7	92.7
gemma-3-12b-it	Qwen2.5-VL-3B-Instruct	84.1	49.4
gemma-3-12b-it	Qwen2-VL-7B-Instruct	81.3	50.0
gemma-3-12b-it	Qwen2.5-VL-7B-Instruct	91.2	74.2
gemma-3-12b-it	llama3-llava-next-8b-hf	93.3	77.2
gemma-3-12b-it	gemma-3-4b-it	99.4	97.6
gemma-3-12b-it	gemma-3-12b-it	99.7	98.4

Table 12: AUC-ROC and TPR at low FPR (FPR=5%) FiMMIA performance metrics for various evaluated Image MLLMs.

<b>Feature Type</b>	<b>Image Features</b>	<b>Audio Features</b>
<b>Texture/Pattern</b>	<ul style="list-style-type: none"> <li>• Local Binary Patterns (LBP) histogram</li> <li>• SIFT Bag of Visual Words (BoVW)</li> </ul>	<ul style="list-style-type: none"> <li>• MFCCs (mean coefficients)</li> <li>• Chroma features (mean)</li> <li>• Tonnetz features (mean)</li> </ul>
<b>Spectral/Frequency</b>	<ul style="list-style-type: none"> <li>• DCT coefficients (low-frequency)</li> </ul>	<ul style="list-style-type: none"> <li>• Spectral centroid (mean)</li> <li>• Spectral bandwidth (mean)</li> <li>• Spectral rolloff (mean)</li> </ul>
<b>Color/Energy</b>	<ul style="list-style-type: none"> <li>• HSV histograms (H, S, V channels)</li> </ul>	<ul style="list-style-type: none"> <li>• RMS energy (mean)</li> <li>• Zero-crossing rate (mean)</li> </ul>
<b>Temporal/Rhythmic</b>	<ul style="list-style-type: none"> <li>• —</li> </ul>	<ul style="list-style-type: none"> <li>• Tempogram features (mean)</li> </ul>

Table 13: Statistical Features Extracted for Image and Audio Classification

# A Browser-based Open Source Assistant for Multimodal Content Verification

Rosanna Milner<sup>1</sup>, Michael Foster<sup>1</sup>, Olesya Razuvayevskaya<sup>1</sup>, Ian Roberts<sup>1</sup>,  
Valentin Porcellini<sup>2</sup>, Denis Teyssou<sup>2</sup>, Kalina Bontcheva<sup>1</sup>,

<sup>1</sup>University of Sheffield, <sup>2</sup>AFP Medialab,

Correspondence: [rosanna.milner@sheffield.ac.uk](mailto:rosanna.milner@sheffield.ac.uk)

## Abstract

Disinformation and false content produced by generative AI pose a significant challenge for journalists and fact-checkers who must rapidly verify digital media information. While there is an abundance of NLP models for detecting credibility signals such as persuasion techniques, subjectivity, or machine-generated text, such methods often remain inaccessible to non-expert users and are not integrated into their daily workflows as a unified framework. This paper demonstrates the VERIFICATION ASSISTANT, a browser-based tool designed to bridge this gap. The VERIFICATION ASSISTANT, a core component of the widely adopted VERIFICATION PLUGIN (140,000+ users), allows users to submit URLs or media files to a unified interface. It automatically extracts content and routes it to a suite of backend NLP classifiers, delivering actionable credibility signals, estimating AI-generated content, and providing other verification guidance in a clear, easy-to-digest format. This paper showcases the tool’s architecture, its integration of multiple NLP services, and its real-world application to detecting disinformation.

## Acknowledgments

This work has been co-funded by the UK’s innovation agency (Innovate UK) grant 10039055 (approved under the Horizon Europe Programme as vera.ai, EU grant agreement 101070093) under action number 2020-EU-IA-0282.

## 1 Introduction and Related Work

Digital disinformation poses a significant threat to democratic societies. The rapid advancement of generative AI, which can produce plausible text, images, and videos in seconds (Zhou and Zafarani, 2021), exacerbates this problem. This technological shift has dramatically increased the volume and sophistication of “fake news”, making the manual

verification of online content a near-impossible task for journalists and fact-checkers (Guo et al., 2022).

In response, the NLP community has developed a broad range of automated methods for information verification (Sharma et al., 2019; Srba et al., 2026). Beyond high-level fake news detection, recent works have shifted toward more fine-grained, explainable assessments of content credibility levels that mirror the aspects of professional fact-checking (Shu et al., 2017; Srba et al., 2026). Some examples of such credibility signals are the presence of propaganda and persuasion techniques (Piskorski et al., 2023), bias and subjectivity (Maab et al., 2024; Piskorski et al., 2023), and AI-generated text (Gehrmann et al., 2019).

However, a significant gap persists between this state-of-the-art research and its practical application by journalists. Most verification tools require technical expertise, and are published as stand-alone models (Srba et al., 2026). This leaves non-technical users without efficient access to the very tools designed to help them. Among the existing systems to support journalistic work and assist general users, the majority focus on a single functionality, such as news bias detection by AllSides<sup>1</sup> and GroundNews<sup>2</sup>, claim verification (Hassan et al., 2017), or propaganda detection (Da San Martino et al., 2020).

To bridge this “research-to-practice” gap, we introduce the VERIFICATION ASSISTANT, a tool designed to integrate multiple backend NLP and media analysis microservices, making state-of-the-art research accessible within the user’s browser. It represents an NLP-focused component within the VERIFICATION PLUGIN, a Chrome extension with over 140,000 active users<sup>3</sup>. The VERIFICATION ASSISTANT provides a single, unified interface where a user can submit a URL (from a news

<sup>1</sup><https://www.allsides.com>

<sup>2</sup><https://ground.news/extension>

<sup>3</sup><http://u.afp.com/plugin>

article or social media post) or upload their own media. The system first extracts relevant content (text, metadata, images) and then dispatches this content to a configurable set of backend NLP microservices (e.g., machine generated text detection, topic and genre classification, URL domain analysis). The VERIFICATION ASSISTANT’s feature set is the result of continuous, participatory design sessions with a panel of journalists and fact-checkers, ensuring its real-world utility. A key focus of our design lies in supporting the multilingual nature of journalism. Finally, the tool aggregates and displays the results in an informative, user-friendly dashboard, moving beyond simple binary “fake/real” labels and providing nuanced, explainable insights.

## 2 Architecture

The VERIFICATION ASSISTANT is a verification “Swiss army knife” that provides a single entry point to various state-of-the-art verification models. There are three main components: (a) the frontend interface, (b) the assistant backend, and (c) the verification services. The frontend provides the main user interface and is written in JavaScript using React Redux. It is distributed as a Chrome extension with language support for English, French, Spanish, Greek, Italian, Arabic, German, Japanese, Portuguese and Hungarian. The code is open source under an MIT licence and is available at <https://github.com/AFP-Medialab/verification-plugin>. The assistant backend is a Quart server application written in Python, with the source code stored in a private repository hosted by the GATE<sup>4</sup> team at the University of Sheffield. It mediates interaction between the frontend and the verification services. Finally, the verification services host individual NLP models trained to assess specific credibility indicators, such as presence of framing in the input text or the likelihood that the content is machine-generated.

The main workflow is illustrated in Figure 1. To verify an article, post, or item of media, the user can submit a URL through the frontend interface. Alternatively, there is an option to upload an image or video file stored locally. The frontend then passes this through to the assistant backend, which scrapes the contents and extracts the text, images, videos, and links contained in it and returns these to the frontend. When the frontend receives the scraped results, it then sends several requests in

<sup>4</sup><https://gate.ac.uk/>

parallel to the various verification services. The remainder of this section introduces these services and provides details of their results in the order they appear in the frontend interface.

### 2.1 Database of Known Fakes

The Database of Known Fakes (DBKF)<sup>5</sup> stores the claims that have been previously debunked by trusted organisations, such as Snopes<sup>6</sup>. Its multilingual text service passes the extracted text through this database to find a potential match. It considers the first 100 characters from the text to imitate a title and prevent distant matches from occurring. If any matches are found, these are returned to the frontend along with a score for each match which are used for ranking purposes. For example, if one result has a higher score than a second result, then that first result is more relevant for the search. Matches with a score ranking over 40 are presented to the user along with a link to the original debunk in the *Detection of previously fact checked claims* section, as shown in Figure 2.

The Fact Check Semantic (FCSS)<sup>7</sup> is the second part of the component that matches a larger proportion of the text to other pieces of text that already exist in a collection of fact-checking databases, such as Snopes.

### 2.2 URL Domain Analysis

The URL Domain Analysis service<sup>8</sup> collects information about a domain from multiple sources to inform the user about its credibility. For example, the Duke Reporter’s Lab<sup>9</sup> maintains a database of known fact-checking sites. Analysis of social media links is performed on individual accounts rather than site domains. For example, if <https://x.com/BBCNews> was referenced in an article, the service would consider the BBCNews account rather than the entire “x.com” domain. The URL for a site that has been explicitly listed as unreliable is flagged to the user as a *Warning*, as shown in Figure 3a. Where the site has been mentioned as part of a debunk, but is not listed as unreliable, this is flagged to the user as a *Mention* together with the details. Known fact-checking sites are flagged to the user as a *Fact Checker*.

<sup>5</sup><https://www.ontotext.com/knowledgehub/current/weverify-project/>

<sup>6</sup><https://www.snopes.com/>

<sup>7</sup><https://kinit.sk/>

<sup>8</sup><https://cloud.gate.ac.uk/shopfront/displayItem/url-domain-analysis>

<sup>9</sup><https://reporterslab.org/>

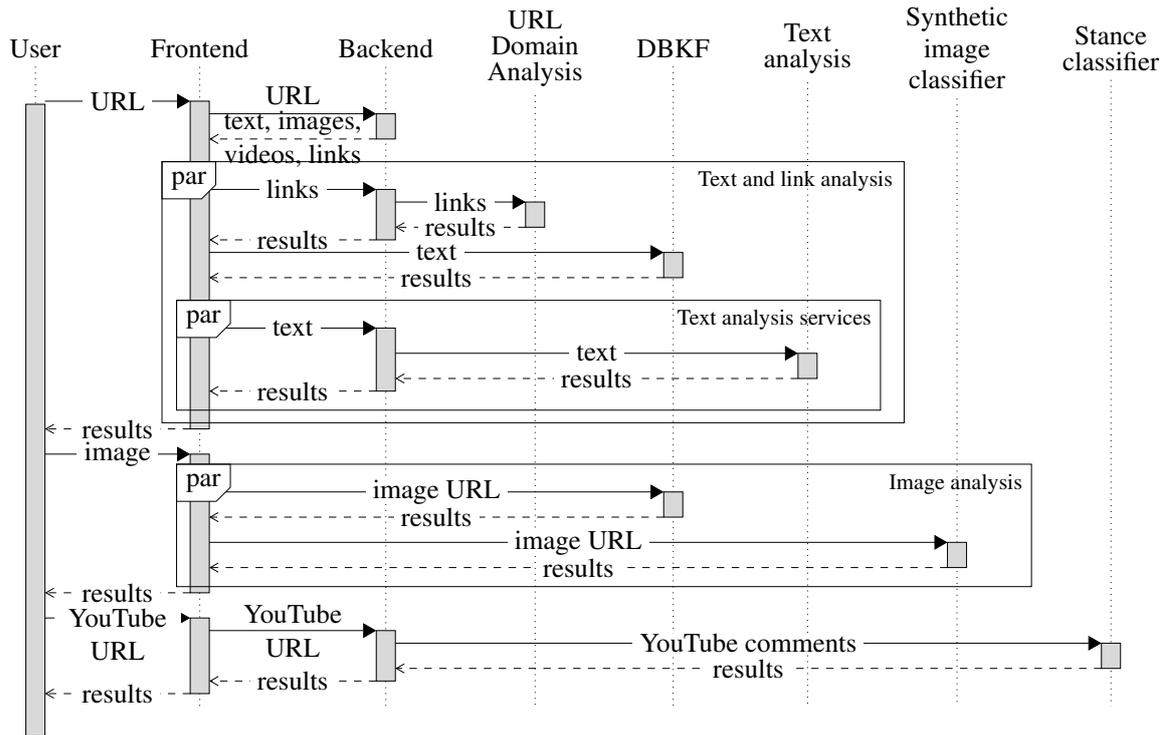
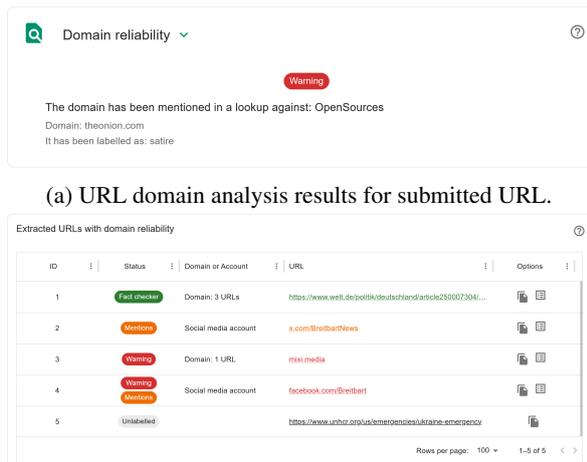


Figure 1: Requests sent by the assistant when checking a typical webpage.



Figure 2: DBKF text service and Fact Check Semantic Search for simulated data.



(b) List of extracted URLs with URL domain analysis results using mock data.

Figure 3: URL domain analysis results.

Extracted links are additionally passed to the source credibility service, with warnings, mentions,

and known fact-checking sites being flagged to the user. The results are then grouped into domains and social media accounts, and are presented in a sortable grid format in the *Extracted URLs with URL Domain Analysis* section, as illustrated in Figure 3b.

### 2.3 Media Analysis

The VERIFICATION ASSISTANT displays extracted image and video thumbnails in the *Extracted media files* section in the order in which they appear on the webpage, to facilitate an easy localisation of media in its original context. illustrated in Figure 4, users may click on an extracted image or video to view more details about it, including a list of tools within the VERIFICATION PLUGIN that are recommended for analysing the media. The possible recommended image analysis tools include image magnifier, metadata retrieval, forensic analysis, OCR, synthetic image detection, geolocalizer and provenance (C2PA). The video analysis tools include video analysis, keyframes, thumbnails, metadata and deepfake. The corresponding image and video analysis tools have been developed by AFP

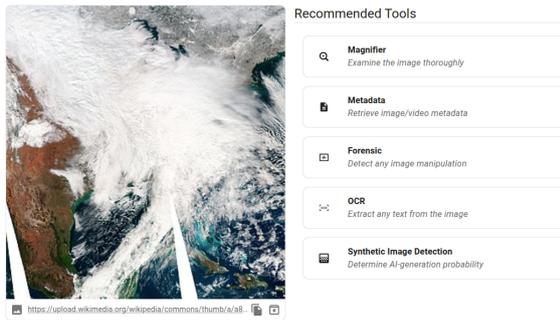


Figure 4: Recommended tools for an image.

Medialab<sup>10</sup>, Borelli Center<sup>11</sup>, ITI CERTH<sup>12</sup>, GRIP UNINA<sup>13</sup> and the University of Sheffield<sup>14</sup>.

Images are additionally sent to the DBKF image search service. Using the image similarity techniques, the service searches for matches in a database of already debunked fakes. If a match is found, this is highlighted to the user as a warning and displayed in the *Detection of previously fact checked claims* section, as described in Section 2.1.

## 2.4 Credibility Signals

The term *credibility signals* refers to a set of context- and content-based indicators that cumulatively contribute to the overall assessment of the credibility of textual information (Srba et al., 2026). The VERIFICATION ASSISTANT incorporates 5 such content-based signals proven to be vital for credibility assessment by prior research (Srba et al., 2026). More specifically, it integrates state-of-the-art classifiers for Framing (Razuvayevskaya et al., 2024; Wu et al., 2023), Genre (Razuvayevskaya et al., 2024; Wu et al., 2023), Persuasion Techniques (Razuvayevskaya et al., 2024; Wu et al., 2023), Subjectivity (Schlicht et al., 2023) and Machine Generated text (Macko et al., 2023) detection. Each signal detector is implemented through a classifier that outputs the associated signal label along with the location in text where the signal is most likely to be present. Together, these credibility signals can be interpreted by the user as an “information nutrition label” (Fuhr et al., 2018). The user can find the results for the credibility signals in the *Extracted text* section separated into different tabs.

<sup>10</sup><https://www.afp.com/en/fact-checking/fact-checking-afp/medialab>

<sup>11</sup><https://ens-paris-saclay.fr/en/research/research-laboratories/centre-borelli>

<sup>12</sup><https://caa.iti.gr/>

<sup>13</sup><https://www.grip.unina.it/>

<sup>14</sup><https://cloud.gate.ac.uk/>

España, aunque no es uno de los principales productores mundiales, ha experimentado un aumento del 23% en el consumo de leche vegetal en los últimos tres años, siendo la de almendras la preferida por el 37% de los consumidores de entre 30 y 40 años, según datos de Nielsen. Este cambio responde a un interés creciente por alternativas más saludables y sostenibles frente a la leche de vaca. Pero las autoridades ahora cuestionan si la leche de almendras es realmente la opción “verde” que muchos creían. Además de la huella hídrica, un nuevo estudio del Instituto Europeo de Nutrición advierte que la leche de almendras industrial contiene apenas un 2% de almendra en promedio, lo que ha suscitado dudas sobre su valor nutricional. “Estamos ante un producto que vende salud, pero ofrece agua con aditivos”, asegura la doctora Ana Moreno, especialista en salud pública. Aunque no se trata aún de una “prohibición formal”, los minoristas ya están siendo advertidos de una posible retirada escalonada de los productos a base de almendras, especialmente aquellos etiquetados como “leche”. Se espera que, bajo la nueva regulación, el término “leche” solo pueda emplearse para productos de origen animal, en línea con recientes directivas de la Unión Europea. La medida ha dividido a consumidores y expertos. Para algunos, se trata de una acción necesaria para frenar prácticas agrícolas insostenibles. Para otros, es una interferencia innecesaria en las

(a) News topic service.

España, aunque no es uno de los principales productores mundiales, ha experimentado un aumento del 23% en el consumo de leche vegetal en los últimos tres años, siendo la de almendras la preferida por el 37% de los consumidores de entre 30 y 40 años, según datos de Nielsen. Este cambio responde a un interés creciente por alternativas más saludables y sostenibles frente a la leche de vaca. Pero las autoridades ahora cuestionan si la leche de almendras es realmente la opción “verde” que muchos creían. Además de la huella hídrica, un nuevo estudio del Instituto Europeo de Nutrición advierte que la leche de almendras industrial contiene apenas un 2% de almendra en promedio, lo que ha suscitado dudas sobre su valor nutricional. “Estamos ante un producto que vende salud, pero ofrece agua con aditivos”, asegura la doctora Ana Moreno, especialista en salud pública. Aunque no se trata aún de una “prohibición formal”, los minoristas ya están siendo advertidos de una posible retirada escalonada de los productos a base de almendras, especialmente aquellos etiquetados como “leche”. Se espera que, bajo la nueva regulación, el término “leche” solo pueda emplearse para productos de origen animal, en línea con recientes directivas de la Unión Europea. La medida ha dividido a consumidores y expertos. Para algunos, se trata de una acción necesaria para frenar prácticas agrícolas insostenibles. Para otros, es una interferencia innecesaria en las

(b) Persuasion techniques service.

España, aunque no es uno de los principales productores mundiales, ha experimentado un aumento del 23% en el consumo de leche vegetal en los últimos tres años, siendo la de almendras la preferida por el 37% de los consumidores de entre 30 y 40 años, según datos de Nielsen. Este cambio responde a un interés creciente por alternativas más saludables y sostenibles frente a la leche de vaca. Pero las autoridades ahora cuestionan si la leche de almendras es realmente la opción “verde” que muchos creían. Además de la huella hídrica, un nuevo estudio del Instituto Europeo de Nutrición advierte que la leche de almendras industrial contiene apenas un 2% de almendra en promedio, lo que ha suscitado dudas sobre su valor nutricional. “Estamos ante un producto que vende salud, pero ofrece agua con aditivos”, asegura la doctora Ana Moreno, especialista en salud pública. Aunque no se trata aún de una “prohibición formal”, los minoristas ya están siendo advertidos de una posible retirada escalonada de los productos a base de almendras, especialmente aquellos etiquetados como “leche”. Se espera que, bajo la nueva regulación, el término “leche” solo pueda emplearse para productos de origen animal, en línea con recientes directivas de la Unión Europea. La medida ha dividido a consumidores y expertos. Para algunos, se trata de una acción necesaria para frenar prácticas agrícolas insostenibles. Para otros, es una interferencia innecesaria en las decisiones de los consumidores adultos. Sin duda, esta decisión parece demasiado severa para un país que valora la libertad de elección.

(c) Subjectivity service.

España, aunque no es uno de los principales productores mundiales, ha experimentado un aumento del 23% en el consumo de leche vegetal en los últimos tres años, siendo la de almendras la preferida por el 37% de los consumidores de entre 30 y 40 años, según datos de Nielsen. Este cambio responde a un interés creciente por alternativas más saludables y sostenibles frente a la leche de vaca. Pero las autoridades ahora cuestionan si la leche de almendras es realmente la opción “verde” que muchos creían. Además de la huella hídrica, un nuevo estudio del Instituto Europeo de Nutrición advierte que la leche de almendras industrial contiene apenas un 2% de almendra en promedio, lo que ha suscitado dudas sobre su valor nutricional. “Estamos ante un producto que vende salud, pero ofrece agua con aditivos”, asegura la doctora Ana Moreno, especialista en salud pública. Aunque no se trata aún de una “prohibición formal”, los minoristas ya están siendo advertidos de una posible retirada escalonada de los productos a base de almendras, especialmente aquellos etiquetados como “leche”. Se espera que, bajo la nueva regulación, el término “leche” solo pueda emplearse para productos de origen animal, en línea con recientes directivas de la Unión Europea. La medida ha dividido a consumidores y expertos. Para algunos, se trata de una acción necesaria para frenar prácticas agrícolas insostenibles. Para otros, es una interferencia innecesaria en las decisiones de los consumidores adultos. Sin duda, esta decisión parece demasiado severa para un país que valora la libertad de elección. Sea como fuere, todo indica que la leche de almendras, tal como la conocemos, tiene los días contados en España.

(d) Machine generated text service.

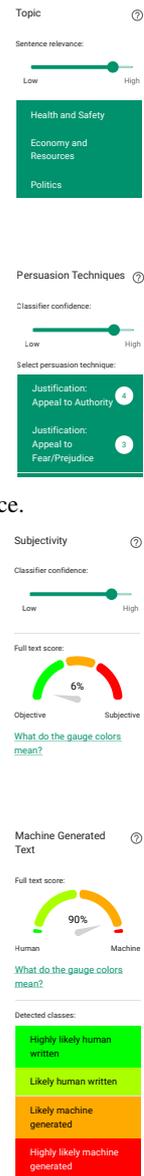


Figure 5: Credibility signals on an AI-generated article.

### 2.4.1 Framing

Framing, referred in VERIFICATION ASSISTANT as *Topic*, is a signal that represents the perspective from which information is presented (Srba et al., 2026). Its purpose is to “frame” information and guide readers toward a particular meaning. The framing classifier (Razuvayevskaya et al., 2024; Wu et al., 2023) has been trained to detect nine main frames: *Economy and Resources, Religious, Ethical and Cultural, Fairness, Equality and Rights, Law and Justice System, Crime and Punishment, Security, Defense and Well-being, Health and Safety, Politics and International Relations*. A single news article frequently incorporates several

overlapping frames simultaneously, and the tool returns all the frames above a certain threshold. For example, Figure 5a shows the top three that have a confidence score  $> 0.8$ . The classifier has been tested on six languages used during fine-tuning and three languages “unseen” during training (Razuvayevskaya et al., 2024). The model demonstrated average performance of  $F_{1_{macro}} = 59.9 \pm 3.1$  and  $F_{1_{micro}} = 61.7 \pm 7.5$  across all 9 languages. The interface highlights sentences that are deemed important by the underlying model in making a decision based on the normalised per-sentence attention scores. The interface allows the user to change the sentence importance threshold by moving a slider.

#### 2.4.2 Genre

Information *genre* is often defined as a way of distinguishing texts based on their writing style (Srba et al., 2026). For the purpose of information credibility detection, the distinction most relevant to the task concerns whether a text maintains an objective tone or incorporates manipulative language. To operationalize this, the incorporated classifier (Razuvayevskaya et al., 2024; Wu et al., 2023) distinguishes between three types of genre: *objective* reporting, *opinionated* pieces, and *satirical* content. Within this taxonomy, objective reporting is characterized by comprehensive coverage of pertinent facts and perspectives. Opinionated news, on the other hand, tends to rely on persuasive or propagandistic techniques, while satirical articles differ from both categories, as they intentionally employ fictionalized material for comedic or critical purposes. Similarly to the framing classifier, the models were tested on both “seen” and “unseen” languages (Razuvayevskaya et al., 2024), with the overall performance of  $F_{1_{macro}} = 49.2 \pm 7.4$  and  $F_{1_{micro}} = 56.7 \pm 6.1$ , averaged across genres and languages. The UI is identical to the framing classifier described in Section 2.4.1.

#### 2.4.3 Persuasion Techniques

*Persuasion techniques*, sometimes referred to as *propaganda techniques* (Piskorski et al., 2023), refer to communication strategies aimed at influencing or manipulating the reader’s opinions. The classifier (Razuvayevskaya et al., 2024; Wu et al., 2023) is trained to identify 23 different techniques which can be organised into six groups: *justification*, *simplification*, *distraction*, *call*, *manipulative wording* and *attack on reputation*. The model achieved an average performance of  $F_{1_{macro}} = 23.7 \pm 5.0$  and

$F_{1_{micro}} = 41.8 \pm 8.6$  across all persuasion techniques and 9 languages, 6 seen and 3 unseen.

Figure 5b presents the user interface. Sentences with a detected persuasion technique(s) are shown as highlighted. The interface also allows the user to hover over a certain highlighted sentence and see the corresponding detected technique(s). A list of the persuasion techniques detected across the complete text is shown on the right hand side. The user can select any technique to only highlight the sentences in which it appears. For this credibility signal, the algorithm provides a confidence score for each sentence-persuasion technique pair. Similarly to the framing classifier, a slider can be moved by the user to change the threshold of these confidence scores for highlighting the associated sentences.

#### 2.4.4 Subjectivity

*Subjectivity* signal refers to the degree to which a news article reflects personal opinions, biases, or emotions rather than strictly objective facts. The integrated classifier (Schlicht et al., 2023), trained based on the CLEF-2023 CheckThat! data challenge (Barrón-Cedeño et al., 2023), returns the score indicating the degree of subjectivity per sentence. The classifier was trained in a multilingual manner, with the test and training sets in English, Turkish and German. The model demonstrated relative robustness across the languages, with the performance of  $F_1 = 0.87$ ,  $F_1 = 0.78$  and  $F_1 = 0.74$  for Turkish, English and German respectively.

As shown in Figure 5c, the interface highlights the subjective sentences identified in text. The overall subjectivity score of the whole text is calculated as a percentage and displayed in a gauge format. The percentage falls into one of three levels: *objective*, *somewhat subjective* and *highly subjective*. Similar to persuasion techniques, a confidence score is provided for sentence subjectivity level and a slider can be moved by the user to change its threshold.

#### 2.4.5 Machine Generated Text

*Machine-generated text* detection as a credibility signal should be distinguished from legitimate automatically generated content, such as machine translation or grammar correction. The classifier (Macko et al., 2023) integrated into the tool, therefore, focuses solely on machine-generated text intended for malicious purposes, such as spreading misinformation. The VERIFICATION ASSISTANT

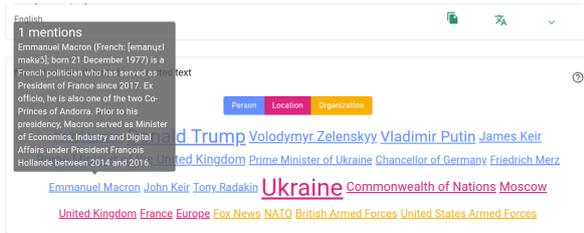


Figure 6: Named entities extracted from a news article on the war in Ukraine.

automatically identifies how likely a piece of text is to be machine generated (e.g. by a large language model), and presents this to the user as a percentage of the full text in a gauge format, as seen in Figure 5d. The text is split into sentences each of which belongs to one of four categories: *highly likely human*, *likely human*, *likely machine generated* and *highly likely machine generated*. The MGT service represents a call to a set of detectors, each trained to identify the generation by a various combinations of language models. Each model was fine-tuned on the data in English, Spanish and Russian languages, and subsequently tested on related languages for each of the training ones: Dutch and German for English, Czech and Ukrainian for Russian, Portuguese and Catalan for Spanish. The evaluation results (Macko et al., 2023) demonstrated high accuracy, with the best-performing model demonstrating  $F_{1_{macro}} = 0.8480$  and weighted  $F_1 = 0.9400$ , averaged across the classes, generator models and languages.

## 2.5 Named Entities

The named entity detector provided by the University of Sheffield identifies the names of people, locations, and organisations in text documents. To extract the named entities, the tool first identifies WikiData concepts in the text. These are then linked to their corresponding DBpedia articles, from which the `rdf:type` field is used for classification. Entities that are not classified as a Person, Location, or Organization are filtered out. The extracted entities are presented to the user as a word cloud, as illustrated in Figure 6. Here, users can select which classes of entities are displayed. The size of each entity in the cloud is proportional to the number of times that entity is mentioned in the text. When the user hovers over an entity, they are shown an abstract from DBpedia along with how many times the entity was mentioned in the text and the link to the corresponding DBpedia article.



Figure 7: Synthesized comments with classifications. Usernames have been blurred to avoid accidental collisions with real (and potentially future) YouTube users.

## 2.6 YouTube Comments with Stance Classifier

When the user submits the URL of a YouTube video, the backend calls the YouTube API to extract the top 10 video comments with their replies. These comments and replies are sent to the stance classifier<sup>15</sup> which identifies comments that *support*, *deny*, or *question* the video, with respect to its title. For comment replies, the classifier identifies whether comments *support*, *deny*, or *query* the original comment. If either are found to be none of these, then the comment is simply labelled as *comment*. This is illustrated in Figure 7 which, to protect the privacy of real users, shows a set of synthesized comments from synthesized users, see Appendix B.

## 3 Evaluation

We collected quantitative feedback from 72 target users—participants of IFCN Global Fact 12<sup>16</sup> conference—in the form of the questionnaires. The survey consisted of 16 questions about the participants, their main occupation, their frequent use of the tool, the features’ usefulness in their workflow, and more qualitative questions about the user interface, their trust in the results or new features they would like to see in updates. Respondents (mainly from AFP, Al Jazeera, dpa, LeadStories, MythDetector, Deutsche Welle and France24) could declare several occupations. Their frequency of use of the tool measured 3.75 out of 5, with 19.44% of respondents using it very often and 41.67% often. Users were asked to rate each new (or enhanced) feature

<sup>15</sup><https://cloud.gate.ac.uk/shopfront/displayItem/stance-classification-multilingual>

<sup>16</sup>Global Fact is the main yearly worldwide gathering of fact-checkers

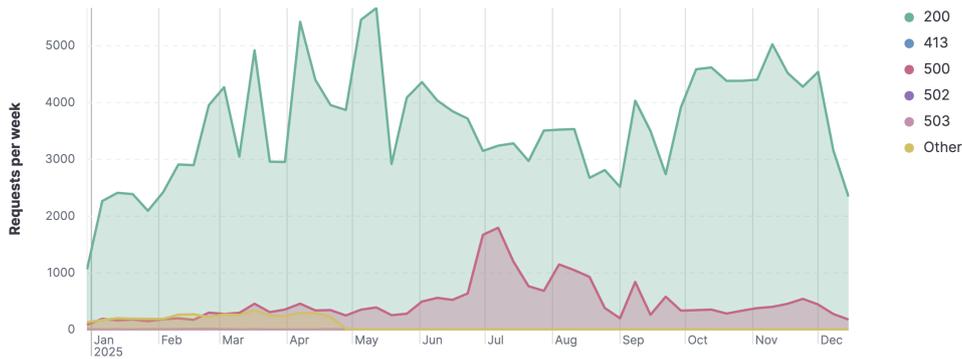


Figure 8: HTTP response code distribution through 2025.

in a scale between Excellent-Good-Average-Poor-Very poor, and this was later transformed into a Likert scale. The overall satisfaction with the VERIFICATION ASSISTANT scored 3.56 out of 5.

Additionally, regular user-centred design (UCD) evaluation sessions were conducted with our target users: professional fact-checkers, journalists, and disinformation researchers. During each session, participants first received a brief overview of the tool and its capabilities. They were then asked to verify a provided piece of information. In addition, participants were encouraged to submit their own content relevant to their daily work—for example, a news article in Greek—and discuss whether they considered the information trustworthy and described how the tool’s outputs influenced their judgment. The primary goal of these sessions was to assess the tool’s practical utility and usability, rather than just raw classifier accuracy. Qualitative feedback was gathered on result agreement, clarity of the output, and the tool’s overall value in their workflow. This feedback was central to our iterative development cycle, leading directly to practical UI refinements. Participants also suggested features or interface improvements that could enhance usability. For example, user feedback prompted the addition of explanatory pop-ups that clarify what each classifier measures and how to interpret its results.

Alongside continuously gathering user feedback, we regularly collect logs of service calls and returned errors to evaluate the robustness of the VERIFICATION ASSISTANT. Figure 8 shows the distribution of HTTP response status codes from the start until the end of 2025. As can be seen, only 14.50% of the requests received a server error response, with the majority of these errors being due to the scraping failures. Overall, 30.82% of user

requests resulted in scraper errors, which could be a result of the users submitting URLs that are either incorrect or cannot be accessed by the tool. Overall, Figure 9 (Appendix A) demonstrates the clear growing trend from the beginning until the end of 2025 in terms of the number of queries submitted by the end-users, with an overall of 18,106 requests submitted from the tool page.

#### 4 Conclusions and Future Work

This paper has presented the VERIFICATION ASSISTANT, a Chrome extension that directly bridges the gap between advanced NLP research and the daily workflow of journalists and fact checkers. By providing a common interface to a plethora of text classifiers, the VERIFICATION ASSISTANT empowers users to analyse content in terms of AI-generation, subjectivity, and other credibility signals. The VERIFICATION ASSISTANT’s integration within a plugin used by over 140,000 users demonstrates its real-world value. Its design is continually refined based on feedback from a panel of fact-checkers, journalists and researchers, ensuring it remains relevant to user needs.

Future work is divided into two main streams. From an engineering perspective, the primary challenge is to develop robust, long-term solutions for content scraping from social media platforms and news sites, which frequently alter their structure and limit access. From a research perspective, we aim to enhance the VERIFICATION ASSISTANT’s utility by integrating more complex AI models. Additionally, increasing the transparency and explainability of their outputs is paramount for building and maintaining user trust. Finally, the recent migration from Webpack to WXT has enabled browser-agnostic functionality, with Firefox support currently under development.

## References

- Alberto Barrón-Cedeño, Firoj Alam, Tommaso Caselli, Giovanni Da San Martino, Tamer Elsayed, Andrea Galassi, Fatima Haouari, Federico Ruggeri, Julia Maria Struß, Rabindra Nath Nandi, Gullal S. Cheema, Dilshod Azizov, and Preslav Nakov. 2023. [The clef-2023 checkthat! lab: Checkworthiness, subjectivity, political bias, factuality, and authority](#). In *Advances in Information Retrieval: 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part III*, page 506–517, Berlin, Heidelberg. Springer-Verlag.
- Giovanni Da San Martino, Shaden Shaar, Yifan Zhang, Seunghak Yu, Alberto Barrón-Cedeño, and Preslav Nakov. 2020. [Prta: A system to support the analysis of propaganda techniques in the news](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 287–293, Online. Association for Computational Linguistics.
- Norbert Fuhr, Anastasia Giachanou, Gregory Grefenstette, Iryna Gurevych, Andreas Hanselowski, Kalervo Jarvelin, Rosie Jones, YiquN Liu, Josiane Mothe, Wolfgang Nejdl, Isabella Peters, and Benno Stein. 2018. [An information nutritional label for online documents](#). *SIGIR Forum*, 51(3):46–66.
- Sebastian Gehrmann, Hendrik Strobelt, and Alexander Rush. 2019. [GLTR: Statistical detection and visualization of generated text](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 111–116, Florence, Italy. Association for Computational Linguistics.
- Zhijiang Guo, Michael Sejr Schlichtkrull, and Andreas Vlachos. 2022. [A survey on automated fact-checking](#). *Trans. Assoc. Comput. Linguistics*, 10:178–206.
- Naeemul Hassan, Fatma Arslan, Chengkai Li, and Mark Tremayne. 2017. [Toward automated fact-checking: Detecting check-worthy factual claims by claim-buster](#). In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '17*, page 1803–1812, New York, NY, USA. Association for Computing Machinery.
- Jonas Maab, Eimara Marrese-Taylor, and Sarah Taylor. 2024. [Media bias detection across families of language models](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4083–4098, Mexico City, Mexico. Association for Computational Linguistics.
- Dominik Macko, Robert Moro, Adaku Uchendu, Jason Lucas, Michiharu Yamashita, Matúš Pikuliak, Ivan Srba, Thai Le, Dongwon Lee, Jakub Simko, and Maria Bielikova. 2023. [MULTITuDE: Large-scale multilingual machine-generated text detection benchmark](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9960–9987, Singapore. Association for Computational Linguistics.
- Jakub Piskorski, Nicolas Stefanovitch, Giovanni Da San Martino, and Preslav Nakov. 2023. [SemEval-2023 task 3: Detecting the category, the framing, and the persuasion techniques in online news in a multilingual setup](#). In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 2343–2361, Toronto, Canada. Association for Computational Linguistics.
- Olesya Razuvayevskaya, Ben Wu, João A. Leite, Freddy Heppell, Ivan Srba, Carolina Scarton, Kalina Bontcheva, and Xingyi Song. 2024. [Comparison between parameter-efficient techniques and full fine-tuning: A case study on multilingual news article classification](#). *PLOS ONE*, 19(5):e0301738.
- Ipek Baris Schlicht, Lynn Khellaf, and Defne Altiok. 2023. [Dwreco at checkthat!-2023: Enhancing subjectivity detection through style-based data sampling](#). 3497:306–317.
- Karishma Sharma, Feng Qian, He Jiang, Natali Ruchansky, Ming Zhang, and Yan Liu. 2019. [Combating fake news: A survey on identification and mitigation techniques](#). *ACM Trans. Intell. Syst. Technol.*, 10(3).
- Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. [Fake news detection on social media: A data mining perspective](#). *CoRR*, arXiv:1708.01967.
- Ivan Srba, Olesya Razuvayevskaya, João A. Leite, Robert Moro, Ipek Baris Schlicht, Sara Tonelli, Francisco Moreno García, Santiago Barrio Lottmann, Denis Teyssou, Valentin Porcellini, Carolina Scarton, Kalina Bontcheva, and Maria Bielikova. 2026. [A survey on automatic credibility assessment using textual credibility signals in the era of large language models](#).
- Ben Wu, Olesya Razuvayevskaya, Freddy Heppell, João A. Leite, Carolina Scarton, Kalina Bontcheva, and Xingyi Song. 2023. [SheffieldVeraAI at SemEval-2023 task 3: Mono and multilingual approaches for news genre, topic and persuasion technique classification](#). pages 1995–2008.
- Xinyi Zhou and Reza Zafarani. 2021. [A survey of fake news: Fundamental theories, detection methods, and opportunities](#). *ACM Comput. Surv.*, 53(5):109:1–109:40.

## A User request trends

Figure 9 shows the number of unique requests or inputs submitted by the users during 2025. A clear growing trend can be observed in terms of the increasing number of events.

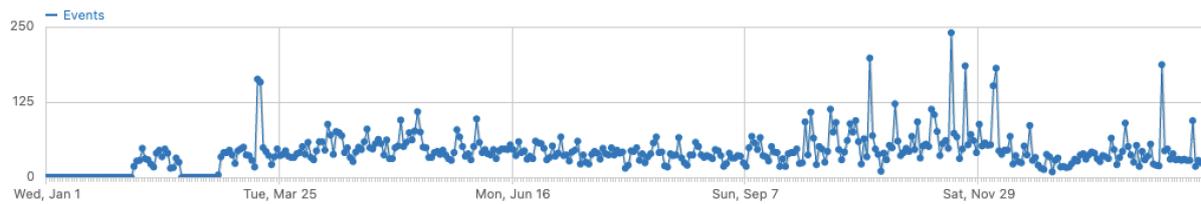


Figure 9: Temporal distribution in the number of events (queries) submitted from the VERIFICATION ASSISTANT page during 2025.

## B Prompt for synthesised YouTube comments

The prompt given to Gemini to generate mock YouTube comments and replies for Figure 7: “I am trying to test a stance classifier. The tool takes a string representing a YouTube comment or twitter post or something and classifies it as either supporting, denying, or questioning a particular statement, or just as a ‘comment’ if it doesn’t support, deny, or question. I need some sample test input for this tool that doesn’t come from anywhere. Please can you generate 100 random inputs with a selection of support, deny, questions and comments. With one input on each line?”

# Inferno: End-to-end Agent-based FHIR Resource Synthesis from Free-form Clinical Notes

Johann Frei<sup>1</sup>   Nils Feldhus<sup>2,3,4</sup>   Lisa Raithel<sup>2,3,4</sup>  
Roland Roller<sup>4</sup>   Alexander Meyer<sup>2,5</sup>   Frank Kramer<sup>1</sup>

<sup>1</sup>IT-Infrastructure for Translational Medical Research, University of Augsburg

<sup>2</sup>BIFOLD – Berlin Institute for the Foundations of Learning and Data   <sup>3</sup>Technische Universität Berlin

<sup>4</sup>German Research Center for Artificial Intelligence (DFKI), Berlin

<sup>5</sup>IKIM, Charité - Universitätsmedizin Berlin

## Abstract

For clinical data integration and healthcare services, the HL7 FHIR standard has established itself as a desirable format for interoperability between complex health data. Previous attempts at automating the translation from free-form clinical notes into structured FHIR resources address narrowly defined tasks and rely on modular approaches or LLMs with instruction tuning and constrained decoding. As those solutions frequently suffer from limited generalizability and structural inconformity, we propose an end-to-end framework powered by LLM agents, code execution, and healthcare terminology database tools to address these issues. Our solution, called Inferno, is designed to adhere to the FHIR document schema and competes well with a human baseline in predicting FHIR resources from unstructured text. The implementation features a front end for custom and synthetic data and both local and proprietary models, supporting clinical data integration processes and interoperability across institutions. Gemini 2.5-Pro excels in our evaluation on synthetic and clinical datasets, yet ambiguity and feasibility of collecting ground-truth data remain open problems.

## 1 Introduction

Large language models (LLMs) have demonstrated strong performance in clinical and biomedical domains, as they have been shown to encode domain-specific knowledge (Singhal et al., 2023; Moor et al., 2023). They are increasingly used to answer clinical questions by processing relevant documents at inference time (Zakka et al., 2024; Chen et al., 2025a; Wang et al., 2024). However, this retrieval-based approach incurs significant latency and computational cost, as documents must be re-processed for every query. This limits usability for tasks such as retrospective analysis or study planning with multiple queries over the same data (Coromilas et al., 2021; Leibig et al., 2022).

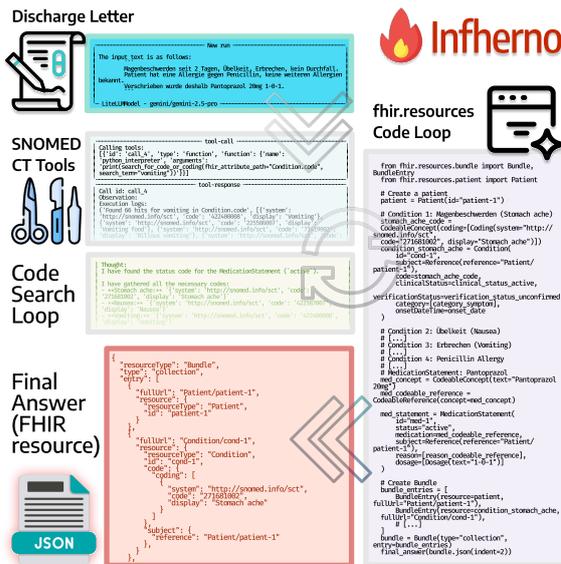


Figure 1: Illustrative example of how Inferno, an agentic approach for FHIR resource synthesis, processes a discharge letter (top left, cyan) using SNOMED CT tools (light blue) and terminology search (green) and fhir.resources code loops (purple, right). After a few iterations including tool calls and observations from a Python executor, the LLM agent proceeds to produce a final answer (red) in a FHIR/JSON format, representing the clinical information on patients and medications.

A more scalable solution is to extract structured representations from clinical text in advance. If the extracted structure preserves the relevant information, it can be queried and reused instantly across multiple applications. This is particularly important for healthcare service providers and clinical data integration efforts (Leroux et al., 2017; Hong et al., 2019; Pimenta et al., 2023). Here, the FHIR<sup>1</sup> standard provides a flexible and interoperable format for representing healthcare data and is increasingly adopted to support standardized access to complex medical information.

Conventional information extraction (IE) meth-

<sup>1</sup><https://fhir.org/>

ods, such as classical named entity recognition (NER), relation extraction, and entity linking, are typically designed for narrowly defined tasks and fixed schemata. As such, they lack the flexibility to adapt to complex clinical contexts and often fail to produce complete, structured clinical representations. In contrast, LLMs have shown promise for IE methods when they are framed as structured prediction tasks (Dagdelen et al., 2024). A recurring challenge in such tasks is ensuring that the output adheres to a specified schema, particularly when downstream components require well-structured inputs (Tavanaei et al., 2024). This is especially true in domains like healthcare, where semantic correctness and schema compliance are essential. Various approaches have been proposed to guide LLM outputs toward structural conformity, including fine-tuning, pre-training, instruction tuning, or constrained decoding (Shin et al., 2021; Geng et al., 2023). Some of them have dealt with text-to-FHIR translation (Sharma et al., 2023; Li et al., 2024; Tabari et al., 2025; Pope and Patooghy, 2025), but often encountered inconsistencies with the desired schema. Agentic LLM approaches that “reason” through intermediate steps using external tools have emerged as a promising solution. Inspired by frameworks like Toolformer (Schick et al., 2023) and ReAct (Yao et al., 2023), models perform multiple tool-augmented reasoning steps, with validation and retry mechanisms to ensure correct output.

As our key contribution, we propose an end-to-end framework that transforms unstructured clinical text into rich, semantically accurate FHIR representations using an agentic LLM-approach. This enables holistic information extraction (Zhang et al., 2025; Shao et al., 2025), supports integration of both legacy and new data, and fosters interoperability across institutions. Our contribution involves:

- (1) An end-to-end implementation for text-to-FHIR translation using LLM agents, SNOMED CT terminology integration and FHIR schema validation;
- (2) Evaluation on real-life and synthetic data, with quantitative and qualitative error analyses across both proprietary and open-source LLMs to characterize failures and their severity;
- (3) A lightweight demonstrator with front-end functionality, supporting both locally run and API-based state-of-the-art LLMs.

## 2 Background

**FHIR and SNOMED CT** FHIR (Fast Healthcare Interoperability Resources) is a widely adopted standard for exchanging healthcare-related data, developed by the HL7 organization. FHIR defines resources as nested documents, often encoded in JSON, with well-defined types, required fields, enumerations, and references to other resources. A single FHIR resource can represent a broad range of entities, from patients and conditions to administrative structures like coverage or questionnaires.<sup>2</sup> FHIR facilitates the structured encoding of complex medical information in an interoperable fashion.

A key feature of FHIR is the integration of internal and external code systems, composed as ValueSets to reference specific entities and concepts. Certain data elements may be constrained to a fixed, FHIR-internal code system to define the set of valid data values.<sup>3</sup> For certain fields, concepts can be referenced from external coding systems such as SNOMED CT<sup>4</sup> or LOINC<sup>5</sup>, and the set of valid data values can be further constrained by individual ValueSets, e.g., to limit data entries for body site to the subset of SNOMED CT concepts that only refer to body structures. To search for codes and terms in a specific ValueSet, FHIR terminology servers provide a standardized interface for querying valid concepts. These servers commonly support multiple external code systems in addition to the FHIR-internal code systems.

While the FHIR schema is capable to accurately and verbosely capture complex clinical situations, it is also subject to structural and semantic ambiguity. Practitioners often use only relevant subsets of data elements depending on their specific use cases. In addition, the standard does not always enforce the encoding of certain information into an unambiguous representation. For instance, a bone fracture of the left limb may be expressed as a *Fracture of bone* SNOMED CT concept along with the *bodySite* element referring to the *Structure of left hand* concept, or purely by referring to the *Fracture of bone of left hand* concept. Dosage information could be phrased only by a free-form

<sup>2</sup>See an example for a Patient resource object at: <https://hl7.org/fhir/R4/patient-example.json.html>

<sup>3</sup>For instance, Condition.clinicalStatus only allows the values active, recurrence, relapse, inactive, remission, and resolved.

<sup>4</sup><https://www.snomed.org/what-is-snomed-ct>

<sup>5</sup><https://loinc.org/get-started/what-loinc-is/>

text element, or by fully utilizing all relevant structured elements, rendering both approaches valid. Clinical notes may also be rather imprecise or ambiguous and require additional and subjective interpretation to fully infer the intended meaning, yet this issue also affects other, non-FHIR-based IE systems. Therefore, comparing predicted and ground-truth FHIR data for semantic equivalence and correctness remains a non-trivial task.

**Related Work** Sharma et al. (2023) presented a pipeline for digitizing prescription images into FHIR using separate components for extraction, normalization, and entity linking, limited to this particular task and mostly small-scale models. Li et al. (2024) first applied LLMs to clinical text-to-FHIR transformation with human-annotated data<sup>6</sup>, but were limited to MedicationStatement resources and faced JSON parsing issues, whereas our validation loop ensures format conformity. Tabari et al. (2025) integrated a syntactic validator and zero- and few-shot strategies into their text-to-FHIR pipeline. Their setup is constrained to sentence-level conversion and exhibits less transparency due to the separation between the OpenAI model and the validator. In contrast, *Infherno*'s tool-calling approach offers a higher degree of transparency and a larger variety of model choices. Pope and Patooghy (2025) explored a variety of FHIR-related tasks as a benchmark, but simplified them to short QA-style problems and also did not consider any elaborate pipeline with tools. Lee et al. (2025) presented FHIR-AgentBench, a comprehensive benchmark for evaluating LLM agents on clinical question answering over FHIR-structured EHR data. Unlike Pope and Patooghy's simplified tasks, they assess complex multi-step retrieval and reasoning over realistic FHIR resources, though their focus remains on querying existing data rather than generation. Idrissi-Yaghir et al. (2025) presented FHIR Workbench for evaluating text-to-FHIR generation, though models struggled without tool augmentation or validation mechanisms. Riquelme Tornel et al. (2025) used GPT-4o and Llama-3.2 alongside clustering and retrieval generation approaches to perform automated FHIR mappings on MIMIC-IV (instead of free text), but missed out on evaluating the results manually. Finally, Schmiedmayer et al. (2025) aimed for an inverse perspective on the translation task by developing a mobile application that

<sup>6</sup>The human-annotated FHIR-GPT data has not been open-sourced to the best of our knowledge.

allows users to interact with FHIR resources via an LLM, while Ehtesham et al. (2025) presented an MCP-based agent for summarization and interpretation. Both represent a FHIR-to-text scenario which is focused on patient understanding.

### 3 *Infherno*, an Agentic Approach

Building on recent work on LLM agents in the medical domain (Liao et al., 2025; Rose et al., 2025; Chen et al., 2025b; Wang et al., 2025), we propose an agentic framework that incorporates tool calls and coding to generate structured FHIR output from unstructured clinical text.

The core task is to transform an unstructured clinical text into semantically corresponding FHIR representations. Our approach follows the Thought-Code-Observation structure proposed as the ReAct framework by Yao et al. (2023), and is implemented using the Smolagents (Roucher et al., 2025) library which supports multi-step LLM agents with Python-code execution. Figure 1 presents a simplified example of the *Infherno* pipeline<sup>7</sup>: Given a discharge letter, *Infherno* which is equipped with tools accessing SNOMED CT, performing Terminology Search, and executing Python code, is tasked to extract information pertinent to patients and medications. In the following, we describe each component:

**Prompt Structure** To guide the agent's behavior, we include relevant contextual information into the prompt (Figure 4, top left). This includes the unstructured input text, a list of target FHIR resource types, supported ValueSets, example code snippets demonstrating FHIR object creation, and a set of instructions on desired behaviors and constraints.

**Terminology Search** To integrate FHIR-specific codes that conform to its specification, we provide our agentic system with an external, retrieval-augmented generation-based function to query particular terms in a set of supported FHIR ValueSets. This enables the agent to rely on external code systems, in particular the SNOMED CT ontology, to retrieve potential search results and include them into its context window. The external function call binds to an external FHIR terminology server to obtain a valid query response.

**Structured Data as Code** Within the agent code execution stage, the agent is incentivized to use the

<sup>7</sup>Figure 4 in Appendix B shows the extended version.

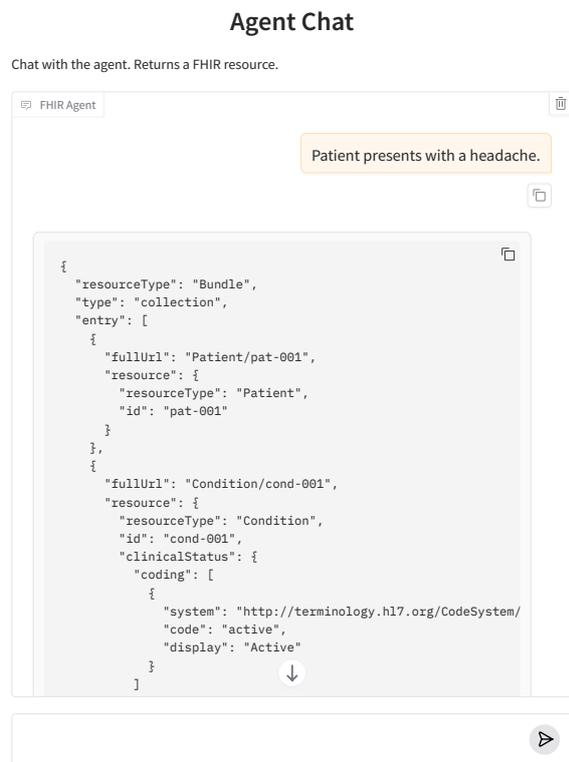


Figure 2: Front end of Infherno showing a short input text and the final answer as given by Gemini-2.5-Pro in the *Agent Chat* function.

`fhir.resources`<sup>8</sup> Python module to create FHIR-conform data instances in a object-oriented fashion. This approach is crucial as it is able to catch morphological and syntactic errors early within the life cycle of the agent loop, and avoids cumbersome data validation that may arise from a purely JSON-centric FHIR document generation by the LLM. Since the library can directly provide error feedback, it can also facilitate the recovery from erroneous code predicted initially by the agent.

**Output Formatting** As part of the Smolagents framework, the code agent can stop the agent loop by the `final_answer` function call. Hereby, the agent is instructed to use the JSON-based object serialization of the `fhir.resources` module. This ensures that the response provides a structurally valid, FHIR-compliant JSON output. To deliver all generated FHIR resources to the user, the agent is instructed to aggregate them into a FHIR *Bundle* that encapsulates the complete set.

<sup>8</sup><https://github.com/nazrulworld/fhir.resources/tree/8.0.0>

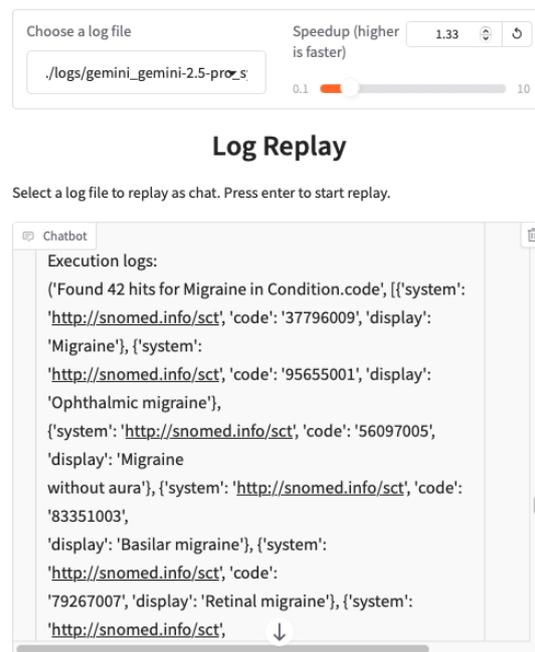


Figure 3: Front end of Infherno showing an intermediate step (Terminology Search) during the text-to-FHIR translation with the *Log Replay* function.

**Front End** The visual interface of Infherno is built on top of Gradio<sup>9</sup> and allows the user to enter arbitrary clinical notes or select pre-defined examples from our synthetic dataset (Figure 2). Intermediate steps including the tool calls and tool responses as well as the reasoning processes (“thoughts” and “observations”) of the ReAct framework (Yao et al., 2023) are shown at inference time. A *Log Replay* tab (Figure 3) also enables to simulate the execution of already conducted experiments at custom speed without the need of an API key. It supports the default Gemini API and OpenAI API via LiteLLM<sup>10</sup> and both API and local Hugging Face models – a feature inherited from Smolagents. The front-end app is available on <https://github.com/j-frei/Infherno>.

## 4 Experiments

To validate our agentic approach, we apply the agent to a set of medical documents to transform the unstructured text into a set of FHIR resources. Such individual comparison is highly complex due to the depth and richness of the FHIR schema and

<sup>9</sup><https://www.gradio.app/>

<sup>10</sup><https://docs.litellm.ai/docs/>

the complexity of clinical language. Therefore, we conduct two experiments, involving both an automatic and a manual inspection and interpretation to assess the prediction quality.

**Data** For the experiments, we use two data sources. We source ten anonymized English EHR documents from the n2c2 2018 challenge (Henry et al., 2020) with drug-related NER annotations. The documents are used with focus on MedicationStatement, given the prior ground truth data. To craft ground-truth FHIR data, we first remove any drug-related NER spans if they do not refer to a medication instruction (e.g., drug allergies are not considered as such), normalize the drug-related NER span annotation into their common concepts to remove redundant mentions and create their corresponding MedicationStatement objects as ground truth. As the n2c2 texts are anonymized, common Protected Health Information (PHI) elements like patient names do not occur and cannot be extracted.

We also rely on ten ChatGPT-synthesized documents resembling medical discharge letters in German which include name and birth date mentions as well as described conditions and medication instructions, allowing us to share PHI-containing documents with remote LLM endpoints. The raw texts required manual editing. Obvious placeholder names, such as “Max Mustermann”, were replaced in the synthetic texts to ensure realistic, non-repetitive patient names and addresses to improve the authenticity of the data.<sup>11</sup> To obtain suitable reference FHIR data, we annotate the documents from the corpus by manually extracting the relevant corresponding FHIR resources, referred to as human baseline (HB). For these documents, we support Patient, Condition, and MedicationStatement as FHIR resources, since we consider these resource types to fit best to key clinical entities.

In general, the FHIR R4 release is targeted as it currently is the latest *normative* release version.

**Experimental Setup** The first experiment aims to implement an evaluation requiring only minimal manual intervention or semantic interpretation. To measure the overlap between the predicted and ground-truth MedicationStatement concepts, we list and compare the .medication fields and manually categorize predicted concepts into TP, FP,

<sup>11</sup>The documents are publicly available on GitHub at <https://github.com/j-frei/Infherno>. One sample is included in the Appendix Figure 5.

FN to eventually calculate precision, recall and  $F_1$  score. As the ground truth involves subjective interpretation of what drug mention is an actual medication instruction, we also add an *relaxed* score setting that ignores certain ground truth entries that were jointly ignored or skipped by all LLMs. We also track the use of SNOMED CT codes and their correctness, as certain predictions only encode the concepts as text field. For the first experiment, both the n2c2 and synthetic data is used.

For the second experiment, we extensively compare the manual annotation and the generated annotation by verifying individual *items* of each FHIR object across multiple FHIR resources. We define an item as a single unit of information, that may refer to, for instance, a single birthDate field but could also refer to a nested object item that describes a reference to a concept from an external coding system. Since the internal structure of certain objects is only meaningful in its entirety, we consider them as monolithic items in the evaluation, rather than decomposing their components. For this experiment, we only use synthetic data as it also covers synthetic PHI elements. Items are stratified into primary (e.g., Condition.code) and secondary (e.g., Condition.verificatiOnStatus) items to account for the information importance differences. Evaluation decision details are highlighted in Appendix Section A.

**Models** For our agentic approach, we use both commercial and open-weight models for the first setup. This includes Claude Sonnet 4.5, Gemini-2.5 Pro<sup>12</sup>, GPT-5 as well as DeepSeek V3.1, Qwen3-235B-A22B-2507, and Qwen3-8B. As the second experiment incurs substantial manual effort, we select Gemini 2.5 Pro as target model as it performed best in the first setup on n2c2 data.

## 5 Results & Discussion

According to Table 1, all LLMs demonstrate high precision with minimal fabrication<sup>13</sup>, but recall varies substantially on n2c2 data, where Gemini 2.5 outperforms all other models. The eased scores reveal that most failures represent cases systematically difficult for all models. For instance, drug information mentioned prior to or outside of the DISCHARGE MEDICATION section is often skipped, as it remains unclear whether such prior mentions

<sup>12</sup><https://ai.google.dev/gemini-api/docs/models>

<sup>13</sup>We found the FP in the synthetic data to be an error in the ground truth.

Data	Model	Scores						Eased Scores			Concepts
		TP	FP	FN	Pr	Re	$F_1$	FN	Re	$F_1$	total/with codes/with correct codes
n2c2 27436 tokens	Claude Sonnet 4.5	78	0	91	<b>1.0</b>	0.462	0.632	31	0.716	0.834	79 / 78 / 78
	DeepSeek V3.1 Chat	82	0	85	<b>1.0</b>	0.491	0.659	25	0.766	0.868	84 / 84 / 82
	Gemini-2.5 Pro	101	0	66	<b>1.0</b>	<b>0.605</b>	<b>0.754</b>	6	<b>0.944</b>	<b>0.971</b>	104 / 104 / 104
	GPT-5	68	0	100	<b>1.0</b>	0.405	0.576	40	0.630	0.773	69 / 63 / 42
	Qwen3-235B-A22B-2507	76	0	92	<b>1.0</b>	0.452	0.623	32	0.704	0.826	77 / 76 / 76
	Qwen3-8B	32	0	134	<b>1.0</b>	0.193	0.323	74	0.302	0.464	36 / 36 / 35
synthetic 4065 tokens	Claude Sonnet 4.5	13	1	3	0.929	0.813	0.867	2	0.867	0.897	14 / 14 / 14
	DeepSeek V3.1 Chat	15	1	1	<b>0.938</b>	<b>0.938</b>	<b>0.938</b>	0	<b>1.0</b>	<b>0.968</b>	16 / 16 / 16
	Gemini-2.5 Pro	14	1	1	0.933	0.933	0.933	0	<b>1.0</b>	0.966	15 / 15 / 15
	GPT-5	15	1	1	<b>0.938</b>	<b>0.938</b>	<b>0.938</b>	0	<b>1.0</b>	<b>0.968</b>	16 / 12 / 9
	Qwen3-235B-A22B-2507	15	1	1	<b>0.938</b>	<b>0.938</b>	<b>0.938</b>	0	<b>1.0</b>	<b>0.968</b>	16 / 15 / 15
	Qwen3-8B	15	1	1	<b>0.938</b>	<b>0.938</b>	<b>0.938</b>	0	<b>1.0</b>	<b>0.968</b>	16 / 15 / 15

Table 1: Evaluation scores from the first experiment, including Precision (Pr), Recall (Re),  $F_1$  scores of *Inferno* with various LLMs evaluated on n2c2 and synthetic data. Best scores are in **bold**.

Category	Worse than HB		Neutral		Better than HB	
	prim	sec	prim	sec	prim	sec
importance						
semantically related	0	4	0	4	0	0
completely identical	0	0	121	83	0	0
lacking in HB	0	10	0	23	13	67
lacking in PD	6	15	0	67	0	0
value difference	0	10	0	12	5	1
semantic halluc. / invalid	1	9	0	0	0	0
total	46		314		86	

Table 2: Second experiment: Results from the manual analysis between predicted (PD) and human baseline (HB) indicating the success and failure cases of *Inferno* for **primary** and **secondary** items. Examples for different categories are shown in Appendix Table 3.

should be considered superseded. Performance on synthetic data confirms that all models extract reliably when ambiguity is low.

Table 2 presents detailed manual analysis comparing predictions against human annotations. Agreement is highest for primary items carrying essential clinical information. Notably, in nearly twice as many cases of quality differences, the model performed better than the human baseline rather than worse. Confabulations remain rare at the semantic level. Divergence occurs primarily on secondary items, which typically involve vague descriptions or ambiguous phrasing rather than explicitly stated primary elements.

**Key Findings** The validation highlights several important observations. First, the phrasing of the input text plays a critical role in annotation consistency. Vague or ambiguous expressions frequently lead to disagreements between the predicted and reference annotations, particularly for secondary items. In contrast, plainly stated and well-structured information is more reliably and consistently captured.

Second, many divergences can be attributed to

the partially subjective nature of FHIR in fringe cases. Minor or nonspecific health issues often fall into a gray area. These may either be excluded or encoded in different ways, such as as a Condition or an Observation. Since the experimental setup allowed only the use of Patient, MedicationStatement, and Condition resource types, the agent was not permitted to use the Observation resource, which limited some of its encoding options.

Furthermore, the *Inferno* agent appears to be more cautious when deciding whether to encode uncertain symptoms. At the same time, it demonstrates stronger recall for clearly stated information that human annotators sometimes overlook, especially with a state-of-the-art LLM like Gemini 2.5. For example, the agent successfully included an address field that was missing in the human annotation. It also inferred an onsetDateTime by subtracting six weeks from the encounter date, which is a detail the human annotator did not encode.

These findings indicate that while human annotations are prone to fatigue and inconsistency, especially in repetitive and detail-oriented tasks, automated agents benefit from their ability to process dense text data using their large context as receptive field. As a result, they can achieve more reliable and comprehensive structured data extraction from our clinical text samples.

## 6 Conclusion

In conclusion, *Inferno* presents a robust and effective framework and interface for transforming unstructured clinical data into standardized FHIR resources. Its agentic design, integrating external knowledge and validation, addresses critical chal-

allenges in clinical information extraction, paving the way for improved data interoperability in healthcare. Future work includes the fine-tuning of smaller language models on the text-to-FHIR task and the integration of more FHIR resource types while further strengthening the robustness, and evaluate the approach on more diverse datasets.

## Limitations

**Dataset and Annotations** Manual evaluation of a larger, more diverse dataset was infeasible given the labor intensity and expertise required. We rely on a single annotator for the human baseline, which inherently introduces a degree of subjectivity.

**Resource Types** Furthermore, our scope was intentionally limited to a subset of FHIR resource types (Patient, Condition, MedicationStatement). Expanding to a broader range of FHIR resources would likely necessitate more verbose guidance in the system prompt, potentially increasing computational cost and latency.

**Legal Remark** Finally, from a legal perspective, it is important to note that Infherno interacts with a FHIR terminology server that includes an initialized SNOMED CT ontology. Therefore, a SNOMED CT license may be required if self-hosting a FHIR terminology server is desired.

**CFG Baseline** While we experimented with context-free grammar-based (CFG) approaches, we found that there are several reasons to object to this design choice:

- Creating a fully conformant FHIR grammar is a major engineering challenge.
- Applying a schema-based decoding may lead to a constrained decoding misalignment issue, which may result in divergences between constrained and unconstrained distributions, and may lead to generation instabilities and poor semantic outputs.
- While considering schema-based generation as an alternative to (fhir.resources-based) code-based FHIR Bundle assembly, other tasks like SNOMED CT code search must also be integrated into the pipeline process. There is no clear way of integrating all components in a non-agentic way. Orchestrating the pipeline in a non-agentic, multi-step pipeline flow rely on an inflexible, rigid process, which could negatively affect the final quality especially in complex clinical

situations that do not fit well into a rigid data transformation process.

- An agentic-based flow allows for certain semantic cross-checks within a FHIR Bundle, that cannot be verified through a CFG-based approach.
- In general, an agentic-based flow simplifies the addition of more (custom/user-defined) tools and validation checks.

**Multilingual Support** While one part of our evaluation focuses on a proprietary model (Gemini 2.5 Pro), our key aim of our work is to remain agnostic of specific LLMs. Consequently, support of other languages depends on the ability of the used LLM to process individual languages correctly rather than our system implementation. The system prompt is written in English and code switching is not used apart from the language in the input text. Since SNOMED CT is mostly an English system, English must be used to query for SNOMED CT codes.

## Ethics Statement

Depending on the selection of the LLM, we want to emphasize that users should be careful in selecting what data they enter. Most of the real-world medical datasets have licences and usage restrictions, so we recommend to use synthetic data only. Users should acknowledge the risk of leaking private data and de-identification.

## Acknowledgments

We thank the reviewers of the EACL 2026 and EMNLP 2025 System Demonstrations tracks for their valuable feedback. This research is funded by the Berlin Institute for the Foundations of Learning and Data (BIFOLD, ref. 01IS18037A) and the German Federal Ministry of Research, Technology and Space (MoMoTuBo, ref. FKZ01ZZ2008).

## References

- Hanjie Chen, Zhouxiang Fang, Yash Singla, and Mark Dredze. 2025a. [Benchmarking large language models on answering and explaining challenging medical questions](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3563–3599, Albuquerque, New Mexico. Association for Computational Linguistics.
- Shan Chen, Pedro José Ferreira Moreira, Yuxin Xiao, Samuel Schmidgall, Jeremy L. Warner, Hugo Aerts,

- Thomas Hartvigsen, Jack Gallifant, and Danielle Bitterman. 2025b. [Medbrowsecomp: Benchmarking medical deep research and computer use](#). In *The Second Workshop on GenAI for Health: Potential, Trust, and Policy Compliance*.
- Ellie J. Coromilas, Stephanie Kochav, Isaac Goldenthal, Angelo Biviano, Hasan Garan, Seth Goldberg, Joon-Hyuk Kim, Ilhwan Yeo, Cynthia Tracy, Shant Ayanian, Joseph Akar, Avinainder Singh, Shashank Jain, Leandro Zimmerman, Maurício Pimentel, Stefan Osswald, Raphael Twerenbold, Nicolas Schaerli, Lia Crotti, and 58 others. 2021. [Worldwide survey of covid-19-associated arrhythmias](#). *Circulation: Arrhythmia and Electrophysiology*, 14(3):e009458.
- John Dagdelen, Alexander Dunn, Sanghoon Lee, Nicholas Walker, Andrew S Rosen, Gerbrand Ceder, Kristin A Persson, and Anubhav Jain. 2024. [Structured information extraction from scientific text with large language models](#). *Nature Communications*, 15(1):1418.
- Abul Ehtesham, Aditi Singh, and Saket Kumar. 2025. [Enhancing clinical decision support and ehr insights through llms and the model context protocol: An open-source mcp-fhir framework](#). *arXiv*, abs/2506.13800.
- Saibo Geng, Martin Josifoski, Maxime Peyrard, and Robert West. 2023. [Grammar-constrained decoding for structured NLP tasks without finetuning](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10932–10952, Singapore. Association for Computational Linguistics.
- Sam Henry, Kevin Buchan, Michele Filannino, Amber Stubbs, and Ozlem Uzuner. 2020. [2018 n2c2 shared task on adverse drug events and medication extraction in electronic health records](#). 27(1):3–12.
- Na Hong, Andrew Wen, Feichen Shen, Sunghwan Sohn, Chen Wang, Hongfang Liu, and Guoqian Jiang. 2019. [Developing a scalable fhir-based clinical data normalization pipeline for standardizing and integrating unstructured and structured electronic health record data](#). *JAMIA Open*, 2(4):570–579.
- Ahmad Idrissi-Yaghir, Kamyar Arzideh, Henning Schäfer, Bahadır Eryilmaz, Mikel Bahn, Yutong Wen, Katarzyna Borys, Eva Hartmann, Cynthia Schmidt, Obioma Pelka, Johannes Haubold, Christoph M Friedrich, Felix Nensa, and René Hirsch. 2025. [Using a diverse test suite to assess large language models on fast health care interoperability resources knowledge: Comparative analysis](#). *J Med Internet Res*, 27:e73540.
- Gyubok Lee, Elea Bach, Eric Yang, Tom Pollard, Alistair Johnson, Edward Choi, Yugang jia, and Jong Ha Lee. 2025. [Fhir-agentbench: Benchmarking llm agents for realistic interoperable ehr question answering](#). *arXiv*, abs/2509.19319.
- Christian Leibig, Moritz Brehmer, Stefan Bunk, Dana-Lyn Byng, Katja Pinker, and Lale Umutlu. 2022. [Combining the strengths of radiologists and ai for breast cancer screening: a retrospective analysis](#). *The Lancet Digital Health*, 4(7):e507–e519.
- Hugo Leroux, Alejandro Metke-Jimenez, and Michael J Lawley. 2017. [Towards achieving semantic interoperability of clinical study data with fhir](#). *Journal of biomedical semantics*, 8:1–14.
- Yikuan Li, Hanyin Wang, Halid Z. Yerebakan, Yoshitsuna Shinagawa, and Yuan Luo. 2024. [FHIR-GPT enhances health interoperability with large language models](#). *NEJM AI*, 1(8):A1cs2300301.
- Yusheng Liao, Shuyang Jiang, Yanfeng Wang, and Yu Wang. 2025. [ReflecTool: Towards reflection-aware tool-augmented clinical agents](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13507–13531, Vienna, Austria. Association for Computational Linguistics.
- Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. 2023. [Foundation models for generalist medical artificial intelligence](#). *Nature*, 616(7956):259–265.
- Nuno Pimenta, António Chaves, Regina Sousa, António Abelha, and Hugo Peixoto. 2023. [Interoperability of clinical data through fhir: A review](#). *Procedia Computer Science*, 220:856–861. The 14th International Conference on Ambient Systems, Networks and Technologies Networks (ANT) and The 6th International Conference on Emerging Data and Industry 4.0 (EDI40).
- Tia Pope and Ahmad Patooghy. 2025. [Comparative evaluation of gpt models in fhir proficiency](#). *ACM Trans. Intell. Syst. Technol.* Just Accepted.
- Álvaro Riquelme Tornel, Pedro Costa del Amo, and Catalina Costa Martínez. 2025. [Large language models for automating clinical data standardization: H17 fhir use case](#). *arXiv*, abs/2507.03067.
- Daniel Philip Rose, Chia-Chien Hung, Marco Lepri, Israa Alqassem, Kiril Gashteovski, and Carolin Lawrence. 2025. [MEDDxAgent: A unified modular agent framework for explainable automatic differential diagnosis](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13803–13826, Vienna, Austria. Association for Computational Linguistics.
- Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunismäki. 2025. [‘smolagents’: a smol library to build great agentic systems](#). <https://github.com/huggingface/smolagents>.

- 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 *Thirty-seventh Conference on Neural Information Processing Systems*.
- Paul Schmiedmayer, Adrit Rao, Philipp Zagar, Lauren Aalami, Vishnu Ravi, Aydin Zahedivash, Dong han Yao, Arash Fereydooni, and Oliver Aalami. 2025. [Llmonfhir](#). *JACC: Advances*, 4(6\_Part\_1):101780.
- Chong Shao, Douglas Snyder, Chiran Li, Bowen Gu, Kerry Ngan, Chun-Ting Yang, Jiageng Wu, Richard Wyss, Kueiyu Joshua Lin, and Jie Yang. 2025. [Scalable medication extraction and discontinuation identification from electronic health records using large language models](#). *Journal of Clinical Epidemiology*.
- Megha Sharma, Tushar Vatsal, Srujana Merugu, and Aruna Rajan. 2023. [Automated digitization of unstructured medical prescriptions](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 794–805, Toronto, Canada. Association for Computational Linguistics.
- Richard Shin, Christopher Lin, Sam Thomson, Charles Chen, Subhro Roy, Emmanouil Antonios Platanios, Adam Pauls, Dan Klein, Jason Eisner, and Benjamin Van Durme. 2021. [Constrained language models yield few-shot semantic parsers](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7699–7715, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, and 1 others. 2023. [Large language models encode clinical knowledge](#). *Nature*, 620(7972):172–180.
- Parinaz Tabari, Alfonso Piscitelli, Gennaro Costagliola, and Mattia de Rosa. 2025. [Assessing the potential of an llm-powered system for enhancing fhir resource validation](#). In *Intelligent Health Systems—From Technology to Data and Knowledge*, pages 803–807. IOS Press.
- Amir Tavanaei, Kee Kiat Koo, Hayreddin Ceker, Shaobai Jiang, Qi Li, Julien Han, and Karim Bouyarmane. 2024. [Structured object language modeling \(SO-LM\): Native structured objects generation conforming to complex schemas with self-supervised denoising](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 821–828, Miami, Florida, US. Association for Computational Linguistics.
- Bowen Wang, Jiuyang Chang, Yiming Qian, Guoxin Chen, Junhao Chen, Zhouqiang Jiang, Jiahao Zhang, Yuta Nakashima, and Hajime Nagahara. 2024. [DiReCT: Diagnostic reasoning for clinical notes via large language models](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Wenxuan Wang, Zizhan Ma, Zheng Wang, Chenghan Wu, Jiaming Ji, Wenting Chen, Xiang Li, and Yixuan Yuan. 2025. [A survey of LLM-based agents in medicine: How far are we from baymax?](#) In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10345–10359, Vienna, Austria. Association for Computational Linguistics.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023. [React: Synergizing reasoning and acting in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Cyril Zakka, Rohan Shad, Akash Chaurasia, Alex R. Dalal, Jennifer L. Kim, Michael Moor, Robyn Fong, Curran Phillips, Kevin Alexander, Euan Ashley, Jack Boyd, Kathleen Boyd, Karen Hirsch, Curt Langlotz, Rita Lee, Joanna Melia, Joanna Nelson, Karim Salim, Stacey Tullis, and 3 others. 2024. [Almanac—retrieval-augmented language models for clinical medicine](#). *NEJM AI*, 1(2):AIoa2300068.
- Xiao Yu Cindy Zhang, Carlos R. Ferreira, Francis Rossignol, Raymond T. Ng, Wyeth Wasserman, and Jian Zhu. 2025. [Casereportbench: An llm benchmark dataset for dense information extraction in clinical case reports](#). In *Proceedings of the sixth Conference on Health, Inference, and Learning*, volume 287 of *Proceedings of Machine Learning Research*, pages 527–542. PMLR.

## A Quantitative Analysis

We manually compared the items of each FHIR object in the human baseline annotation with their corresponding agent-generated equivalents. Missing items were tracked using + for those absent in the human baseline and - for those missing in the prediction. Potentially equivalent values were categorized as exact matches (==), semantically equivalent (=), or different (+-).

Items were tagged with / if the prediction value was preferable to the baseline, and with | if the baseline value was preferable. Items were left untagged when no clear preference could be determined or justified.

To distinguish essential from less relevant features during evaluation, we annotated core FHIR objects of major importance, such as patient information or the main diagnosis extracted from the input text, with ! as primary items, and less critical elements, like vaguely described symptoms, with ? as secondary items. The importance level

Level of equivalence		Infherno (w/ Gemini-2.5-Pro)	Human Baseline
<b>Worse than HB</b> (field not referenced) (hallucination)	[ +-?] [X+?]	Condition.bodySite Condition.category	"code": "52795006", "display": "Forehead structure" N/A
<b>Neutral</b> (optional field missing) (total equivalence)	[-?] [==]	MS.dosage .doseAndRate.type Condition.severity	N/A N/A "code": "255604002", "display": "Mild"
<b>Better than HB</b> (inaccurate reference for <i>throw up</i> )	[/+-!]	Condition.code	"code": "422400008", "display": "Vomiting" "code": "422587007", "display": "Nausea (finding)"

Table 3: Examples for level of equivalence and the manual validation between system output and human baseline.

of a FHIR object determined the default importance of its internal items, unless overridden by manually applied, item-specific tags. These overrides were primarily used to demote non-essential items such as `Condition.verificationStatus` or `Patient.name.use`. Conversely, items from otherwise crucial FHIR resources, such as `Condition.subject` or `Condition.code`, were generally considered primary, as they carry core informational content. Some examples are shown in Table 3.

## B Examples

Figure 4 illustrates a complete example of an text-to-FHIR translation flow.

## C Example of a Synthetic Clinical Document

The Figure 5 shows the first document from our synthesized text corpus. All documents are accessible on GitHub at the following url: <https://github.com/j-frei/Infherno>.

...system

You are a code agent with expertise in Information Extraction of medical information from free text. Your task is to translate clinical information into valid FHIR R4 resources, using step-by-step reasoning and supported tooling.

You work by **thinking step by step** in cycles of:

- Thought:** - where you explain your current reasoning and intended approach.
- Code:** - where you write Python code using 'fhir.resources' or supported tools, always ending with 'end\_code'.
- Observation:** - where you receive the printed outputs or results of your code.

You repeat this cycle until you're ready to give the final output using the 'final\_answer()' tool.

At each cycle step, you must start next line plainly with "Thought:", "Code:", or "Observation:" (without quotes), otherwise it will fail.

Tools Available:

You ONLY have access to the following tools:

- Bundle:** (the top-level container)
- Patient**
- Condition**
- MedicationStatement**

All output must be a valid 'bundle' of these resources using the 'fhir.resources' Python package.

Attribute Coding Rules:

- You must only use codings returned from 'search\_for\_code\_or\_coding', and only for these specific FHIR attribute paths:
- Patient.name.use
- Patient.contact.system
- Patient.address.use
- Patient.address.type
- Patient.maritalStatus
- Patient.contact.relationship
- Condition.clinicalStatus
- Condition.verificationStatus
- Condition.category
- Condition.severity
- Condition.code
- Condition.bodySite
- Condition.stage.summary
- Condition.stage
- Condition.evidence
- MedicationStatement.status
- MedicationStatement.statusReason
- MedicationStatement.category
- MedicationStatement.medication
- MedicationStatement.codableConcept
- MedicationStatement.reasonCode
- MedicationStatement.dosage.additionalInstruction
- MedicationStatement.dosage.timing.repeat.dayOfWeek
- MedicationStatement.dosage.timing.repeat.when
- MedicationStatement.dosage.timing.code
- MedicationStatement.dosage.asNeeded
- MedicationStatement.dosage.asNeededCodableConcept
- MedicationStatement.dosage.site
- MedicationStatement.dosage.route
- MedicationStatement.dosage.method
- MedicationStatement.dosage.doseAndRate.type

Use this pattern to search for codes:

For example, to search for a code for "douleurs abdominales" in the Condition.code attribute, you would call:

```
Code: """
search_results =
search_for_code_or_coding(fhir_attribute_path="Condition.code",
search_term="abdominal pain")
print(search_results)
"""
end_code
Observation: ('Found results for 'abdominal pain' in 'Condition.code' truncated to max. of 10 results.', [{'code': '21522001', 'system': 'http://snomed.info/sct', 'display': 'Abdominal pain (finding)', 'code': '162942009', 'system': 'http://snomed.info/sct', 'display': 'Abdominal wall pain (finding)', 'code': '4597903', 'system': 'http://snomed.info/sct', 'display': 'Abdominal wind pain (finding)', 'code': '9991608', 'system': 'http://snomed.info/sct', 'display': 'Abdominal colic (finding)', 'code': '5458694', 'system': 'http://snomed.info/sct', 'display': 'Lower abdominal pain (finding)', 'code': '8332093', 'system': 'http://snomed.info/sct', 'display': 'Upper abdominal pain (finding)', 'code': '116290004', 'system': 'http://snomed.info/sct', 'display': 'Acute abdominal pain (finding)', 'code': '282210013193', 'system': 'http://snomed.info/sct', 'display': 'Abdominal muscle pain (finding)', 'code': '11985807', 'system': 'http://snomed.info/sct', 'display': 'Chronic abdominal pain (finding)', 'code': '162946062', 'system': 'http://snomed.info/sct', 'display': 'Central abdominal pain (finding)'}])
"""
end_code
Example (search for Patient.gender code for a man)
Code: """
search_results =
search_for_code_or_coding(fhir_attribute_path="Patient.gender",
search_term="male")
print(search_results)
"""
end_code
Observation: ('Found results for 'man' in 'Patient.gender' truncated to max. of 10 results.', [{'code': 'male', 'code': 'female', 'code': 'other', 'code': 'unknown'}])
"""
end_code
You need to search and print the results before the creating the final FHIR data in the final step, where you pick the most fitting items from the observed print outputs.


Process Strategy:



- Start with a Thought: describing your plan to extract relevant medical info and convert to FHIR resources.
- Use Code: blocks to either:
  - call 'search_for_code_or_coding' to get valid SNOMED/HL7 codings and print it
  - construct FHIR resource objects ('Patient', 'Condition', etc.) (when you have all the needed info from previous steps)
  - build the final Bundle
  - call 'final_answer(bundle.json(indent=2))' to output the final result
- At each step, explain and print intermediate info you'll need in later steps.
- Never reuse tool parameters unnecessarily; only call tools when needed.
- Stick to valid FHIR attributes and use 'fhir.resources' models.



Final Code Block



When you're ready, build the final bundle.



Let's consider an example where the input text is "Herr Meyer klagt über Bauchschmerzen".



We need to extract the Patient and Condition information from this text and create a FHIR Bundle.



So you will first search for the coding for "abdominal pain" in the Condition.code attribute to find the appropriate code for abdominal pain, as well as other relevant codings (e.g. for Patient.gender).



This will take a few steps, and you will print the results of each search.



Then, you will create a Patient resource with the name "Meyer" and a Condition resource with the found code. Finally, you will bundle them together in a FHIR Bundle, using the codings and codes that you have already queried in the previous steps, like this:


```

## System Prompt

```
Code: """
from fhir.resources.bundle import Bundle, BundleEntry
from fhir.resources.patient import Patient
from fhir.resources.condition import Condition
from fhir.resources.humanname import HumanName
from fhir.resources.codableconcept import CodableConcept
from fhir.resources.reference import Reference
from fhir.resources.reference import Reference

# Build Patient
name = HumanName(use="official", family="Meyer")
patient = Patient(id="pat-001", gender="male", name=name)

# Build Condition
condition_code = Coding(system="http://snomed.info/sct", code="21522001", display="Abdominal pain")
clinical_status = Coding(system="http://terminology.hl7.org/CodeSystem/condition-clinical", code="active")
condition = Condition(subject=Reference(reference="Patient(patient-id)"), condition_code = CodableConcept(coding=condition_code), condition.clinicalStatus = CodableConcept(coding=condition_code))

# Create Bundle
bundle = Bundle(type="collection", entry=[
    BundleEntry(resource=patient),
    BundleEntry(resource=condition)
])

final_answer(bundle.json(indent=2))
"""
end_code
Now, this code example is just for illustration purposes for an imaginary male Patient with the name 'Meyer'. Do not use any information from it directly.
You will need to adapt the code to your specific task and FHIR resources matching the upcoming input text. Yet, the basic task structure remains similar.
Important: The codes/codings should be determined before this step.
Always:
- Use valid 'fhir.resources' models and their documented fields
- Use codings only from 'search_for_code_or_coding'
- Output a single Bundle using 'final_answer' formatted in JSON (and indent=2)
- Verify the codes/codings before using them in the final step.
- Hard-code the picked coding/code data which was observed by the search tool in prior steps
- Don't put any code or text after the final_answer() call. Otherwise, it will fail.
Never:
- NEVER use unverified FHIR resources
- NEVER use unverified codes
- NEVER call tools with dict-style args (e.g. 'X' tool({'arg': val}))
- NEVER chain tool calls in the same block with unpredictable results
- NEVER create FHIR resources to information that is not present in the input text.
- NEVER use the search_for_code_or_coding in the final step.
Now begin your task!
If you solve the task correctly, you will receive a reward of $1,000,000.
```

## Input Document/EHR

The input text is as follows:

Herrnbeschwerden seit 2 Tagen, Übelkeit, Erbrechen, kein Durchfall.  
Patient hat eine Allergie gegen Penicillin, keine weiteren Allergien bekannt.  
Verschrieben wurde deshalb Pantoprazol 20mg 1-0-1.

- LiteLJModel - gemini/gemini-2.5-pro

## FHIR Resource

```
{
  "resourceType": "Bundle",
  "type": "collection",
  "entry": [
    {
      "fullUrl": "Patient/patient-1",
      "resource": {
        "resourceType": "Patient",
        "id": "patient-1"
      }
    },
    {
      "fullUrl": "Condition/cond-1",
      "resource": {
        "resourceType": "Condition",
        "id": "cond-1",
        "clinicalStatus": {
          "coding": [
            {
              "system": "http://terminology.hl7.org/CodeSystem/condition-clinical",
              "code": "active",
              "display": "Active"
            }
          ]
        },
        "verificationStatus": {
          "coding": [
            {
              "system": "http://terminology.hl7.org/CodeSystem/condition-verification-status",
              "code": "unconfirmed",
              "display": "Unconfirmed"
            }
          ]
        },
        "category": [
          {
            "coding": [
              {
                "system": "http://terminology.hl7.org/CodeSystem/condition-category",
                "code": "encounter-diagnosis",
                "display": "Encounter Diagnosis"
              }
            ]
          }
        ],
        "code": {
          "coding": [
            {
              "system": "http://snomed.info/sct",
              "code": "271881002",
              "display": "Stomach ache"
            }
          ]
        },
        "subject": {
          "reference": "Patient/patient-1",
          "onsetDateTime": "2025-06-19"
        }
      }
    }
  ]
}
```

```
Call id: call_4
Observation:
('Found 66 hits for vomiting in Condition.code', [{'system': 'http://snomed.info/sct', 'code': '422400008', 'display': 'Vomiting', 'code': '225566007', 'system': 'http://snomed.info/sct', 'display': 'Vomiting food', 'code': '71419002', 'system': 'http://snomed.info/sct', 'display': 'Vomiting', 'code': '8765009', 'system': 'http://snomed.info/sct', 'display': 'Hematemesis', 'code': '91863907', 'system': 'http://snomed.info/sct', 'display': 'Viral gastroenteritis', 'code': '23971007', 'system': 'http://snomed.info/sct', 'display': 'Acute vomiting', 'code': '49206006', 'system': 'http://snomed.info/sct', 'display': 'Habit vomiting', 'code': '45994004', 'system': 'http://snomed.info/sct', 'display': 'Erotic vomiting', 'code': '6372006', 'system': 'http://snomed.info/sct', 'display': 'Chronic vomiting', 'code': '20647008', 'system': 'http://snomed.info/sct', 'display': 'Vomiting symptom'}])
Last output from code snippet:
None

tool-call
Calling tools:
{'id': 'call_4', 'type': 'function', 'function': {'name': 'python_interpreter', 'arguments': {'print(search_for_code_or_coding(fhir_attribute_path="Condition.code", search_term="vomiting"))'}}}

Thought:
I have found the status code for the MedicationStatement ('active').
I have gathered all the necessary codes:
- Stomach ache: 'http://snomed.info/sct', 'code': '271881002', 'display': 'Stomach ache'
- Nausea: 'http://snomed.info/sct', 'code': '422587007', 'display': 'Nausea'
- Vomiting: 'http://snomed.info/sct', 'code': '422400008', 'display': 'Vomiting'
- Penicillin Allergy: 'http://snomed.info/sct', 'code': '91926005', 'display': 'Allergy to penicillin'
- Encounter Clinical Status: 'http://terminology.hl7.org/CodeSystem/condition-clinical', 'code': 'active', 'display': 'Active'
- Condition Verification Status: 'http://terminology.hl7.org/CodeSystem/condition-verification-status', 'code': 'unconfirmed', 'display': 'Unconfirmed'
- Condition Category (Symptom): 'http://terminology.hl7.org/CodeSystem/condition-category', 'code': 'encounter-diagnosis', 'display': 'Encounter Diagnosis'
- Condition Category (Allergy): 'http://terminology.hl7.org/CodeSystem/condition-category', 'code': 'problem-list-item', 'display': 'Problem List Item'
- MedicationStatement Status: 'http://hl7.org/fhir/CodeSystem/medication-statement-status', 'code': 'active', 'display': 'Active'

Now I will construct the FHIR resources and the final bundle.
1. A 'Patient' resource (anonymous, as no name is given).
2. Three 'Condition' resources for the symptoms (stomach ache, nausea, vomiting), linked to the patient. They will have an 'onsetDateTime' approximation based on 'seit 2 Tagen'.
3. One 'Condition' resource for the penicillin allergy, linked to the patient.
4. One 'MedicationStatement' for Pantoprazol, linked to the patient. It will include the dosage '1-0-1'. The medication will be represented as text. The reasonCode will point to the stomach complaint condition.

I will now write the Python code to create these resources and bundle them. I will hardcode the date for today, assuming the consultation happened
```

## Intermediate Steps (Tool-Calling, Reasoning, Python Code)

Figure 4: Extended example of clinical note synthesis with Infferno including the System Prompt and a longer snippet from the tool calls and generated Python code which yields the FHIR Resource.

Betreff: Arztberichtsbrief - Patienteninformationen

Sehr geehrter Dr. Peters,

hiermit möchte ich Ihnen einen aktuellen Bericht über den Gesundheitszustand von Herrn Uwe Jaeger, geboren am 10. Februar 1975, vorlegen. Herr Jaeger wurde am 20. Juni 2023 in unserer Klinik, dem St. Ursula Krankenhaus, zur weiteren Untersuchung und Behandlung aufgenommen.

Anamnese:

Herr Jaeger suchte unsere Notaufnahme mit anhaltenden Beschwerden im Magen-Darm-Bereich auf. Er berichtete über starke Bauchschmerzen, Übelkeit, Erbrechen und Gewichtsverlust in den letzten vier Wochen. Er verneinte jegliche vorherige Operationen oder relevante Vorerkrankungen. Herr Jaeger ist Nichtraucher und konsumiert keinen Alkohol.

Klinischer Befund:

Bei der körperlichen Untersuchung zeigten sich eine allgemeine Schwäche und ein mäßig abgeschwächter Allgemeinzustand. Der Bauch war diffus druckempfindlich, ohne spürbare Vergrößerungen der Organe. Keine Zeichen einer Peritonitis waren erkennbar. Die übrige körperliche Untersuchung ergab keine auffälligen Befunde.

Diagnostische Maßnahmen:

Um die Ursache der Beschwerden zu ermitteln, wurden bei Herrn Jaeger verschiedene diagnostische Tests durchgeführt. Eine Blutuntersuchung ergab eine erhöhte Anzahl weißer Blutkörperchen und eine leichte Anämie. Der Leberfunktionstest zeigte normale Werte. Ein abdominales Ultraschall wurde durchgeführt, das keine strukturellen Abnormalitäten zeigte. Eine Endoskopie des oberen Verdauungstrakts wurde ebenfalls durchgeführt, bei der eine erosive Gastritis festgestellt wurde.

Diagnose:

Basierend auf den klinischen Symptomen, den Laborergebnissen und der Endoskopie wurde bei Herrn Jaeger die Diagnose einer erosiven Gastritis gestellt.

Therapie:

Um die Symptome zu lindern und die Schleimhaut im Magen zu heilen, wurde Herr Jaeger eine Kombinationstherapie verschrieben. Er erhält eine Protonenpumpenhemmer (PPI) für acht Wochen, um die Magensäureproduktion zu reduzieren. Zusätzlich wurde ihm ein Antazidum verschrieben, um den sofortigen Effekt einer schnellen Symptomlinderung zu erzielen. Er erhielt auch Anweisungen zur Vermeidung von auslösenden Nahrungsmitteln, wie scharfe und säurehaltige Lebensmittel.

Verlauf und Prognose:

Herr Jaeger hat die empfohlene Therapie begonnen und wurde über mögliche Nebenwirkungen und Maßnahmen zur Verbesserung seines Gesundheitszustands aufgeklärt. Wir werden ihn in regelmäßigen Abständen zu Follow-up-Terminen einladen, um den Verlauf seiner Symptome zu überwachen und gegebenenfalls weitere Untersuchungen durchzuführen.

Abschließend möchte ich Ihnen versichern, dass wir die bestmögliche Versorgung für Herrn Jaeger sicherstellen und eng mit ihm zusammenarbeiten werden, um eine schnelle Genesung zu erreichen.

Bei weiteren Fragen stehe ich Ihnen gerne zur Verfügung.

Mit freundlichen Grüßen,

Dr. Anna Karolin Vogel  
Fachärztin für Innere Medizin  
St. Ursula Krankenhaus

Figure 5: The full text from the first document from the synthetic corpus.

# **BOOM**: *Beyond Only One Modality* KIT's Multimodal Multilingual Lecture Companion

Sai Koneru, Fabian Retkowski, Christian Huber, Lukas Hilgert,  
Seymanur Akti, Enes Yavuz Ugan, Alexander Waibel, Jan Niehues  
Karlsruhe Institute of Technology

[firstname.lastname@kit.edu](mailto:firstname.lastname@kit.edu)

## Abstract

The globalization of education and rapid growth of online learning have made localizing educational content a critical challenge. Lecture materials are inherently multimodal, combining spoken audio with visual slides, which requires systems capable of processing multiple input modalities. To provide an accessible and complete learning experience, translations must preserve all modalities: text for reading, slides for visual understanding, and speech for auditory learning. We present **BOOM**, a multimodal multilingual lecture companion that jointly translates lecture audio and slides to produce synchronized outputs across three modalities: translated text, localized slides with preserved visual elements, and synthesized speech. This end-to-end approach enables students to access lectures in their native language while aiming to preserve the original content in its entirety. Our experiments demonstrate that slide-aware transcripts also yield cascading benefits for downstream tasks such as summarization and question answering. The demo video and code can be found at <https://ai4lt.github.io/boom/><sup>1</sup>.

## 1 Introduction

Access to educational content in a learner's native language greatly enhances the learning experience for university students. Localizing lecture material reduces communication barriers, improves accessibility, and enables learners to engage more deeply with complex concepts. As higher education becomes increasingly global, the ability to provide multilingual lecture content both in-person and online has become essential to increase accessibility to educational resources (Muthuswamy and Varshika, 2023; Gambier, 2023).

With the ongoing digitalization of teaching, lecture content itself is inherently multimodal. The

primary modality is the lecture audio, which can be converted into transcripts via Automatic Speech Recognition (ASR) (Pham et al., 2019; Radford et al., 2022). Instructional material is presented through slides, and additional outputs, such as summaries, chapters, and question–answer interactions, can be generated based on the transcript in a cascaded setup using modern Large Language Model (LLM)-based systems to enhance the learning experience (Waibel and Fuegen, 2012; Waibel, 2014; Anderer et al., 2025; Retkowski et al., 2025). To ensure accessibility for all students, including non-native speakers, these outputs should also be available in multiple languages. Effective localization must therefore handle this diversity of content, spanning audio, text, and visual materials, making lecture translation a truly multimodal challenge.

This multimodality introduces complexity but also offers valuable contextual signals. Images often contain additional cues ranging from scene information in natural images and definitions, formulas, diagrams, and domain-specific terminology in slides that help disambiguate spoken content (Nguyen et al., 2025) and support downstream tasks such as Summarization (SUM) and Question Answering (QA). Leveraging these visual cues enables translation systems to move beyond audio-only processing and incorporate richer semantic information throughout the lecture translation pipeline (Waibel, 2018; Chen et al., 2024; Sinhamahapatra and Niehues, 2025).

Machine Translation (MT) forms the foundation of localization, evolving from rule-based systems (Hutchins, 2004) to Neural MT (NMT; Vaswani et al. 2017; Koehn and Knowles 2017; Johnson et al. 2017) and then Speech Translation (ST), which directly translates spoken content. Modern ST handles many languages (Barrault et al., 2023) but often processes short segments, limiting context and potential to benefit from multimodality.

In this work, we address multimodality on both

<sup>1</sup>All released code and models are licensed under the MIT License



(a) Original English Slide



(b) Translated German Slide

Figure 1: Comparison of the English (original) and German (translated) slides. Text outside the images is translated with a unimodal system for efficiency, while text inside the images is translated using a multimodal system.

the input and output sides of lecture localization. On the input side, we incorporate slide screenshots into the ST pipeline to provide contextual grounding that improves translation accuracy and downstream LLM performance. On the output side, we tackle the challenge of localizing lecture slides themselves. Slides often contain text embedded within images, such as diagram labels, equations, or annotations, that existing ST tools typically ignore. Localizing such material requires detecting, recognizing, translating, and re-rendering text while preserving layout, alignment, font style, and visual coherence (illustrated in Figure 1).

To overcome these limitations, we extend the Lecture Translator (LT) software (Huber et al., 2023) with OmniFusion (Koneru et al., 2025), a multilingual multimodal ST model that uses slide images to enrich translation. We further introduce a fully open-source slide translation system capable of translating text inside slide images and rendering it back into its original layout, enabling complete slide localization. Together, these components form a unified multimodal lecture localization pipeline that combines improved ST with synchronized slide translation, significantly enhancing accessibility for learners across languages.

Our main contributions include:

- Adapt and integrate OmniFusion to leverage lecture slide screenshots during live translation, by extracting relevant slides from segmented audio.
- Introduce an open-source image-to-image translation pipeline with modular components, enabling future research on full-image/slide translation and rendering.
- Demonstrate the impact of including images

on ST for downstream NLP tasks across different LLMs, showing performance improvements in different language pairs. We also evaluate several optical character recognition (OCR) models and the translation quality of unimodal and multimodal NMT models for image translation.

## 2 System Description

To fully localize lecture content, including audio and slides, across multiple modalities and languages, and to support accessibility tasks such as SUM and QA, we develop multimodal translation systems. Our approach performs multimodal ST and leverages the resulting transcripts for downstream LLM tasks. We also translate slides by converting text and images into the target language while preserving their layout and visual coherence.

To map visual context to each audio segment and improve usability, we built a PDF viewer that displays slides with overlaid captions synchronized to the presenter’s selected slide (Figure 8). This interface enables participants to follow translations while viewing slides and allows the system to automatically identify which slide corresponds to each audio segment, providing essential context for multimodal ST.

In this section, we first describe the multimodal ST pipeline, including how slide images are extracted and associated with audio segments. Next, we outline how the resulting translations are used for downstream SUM and QA. Finally, we present the slide-translation process, detailing how text embedded within slide images is detected, translated, and re-rendered. Additionally, details about Text-to-Speech (TTS) are described in Appendix A.1.

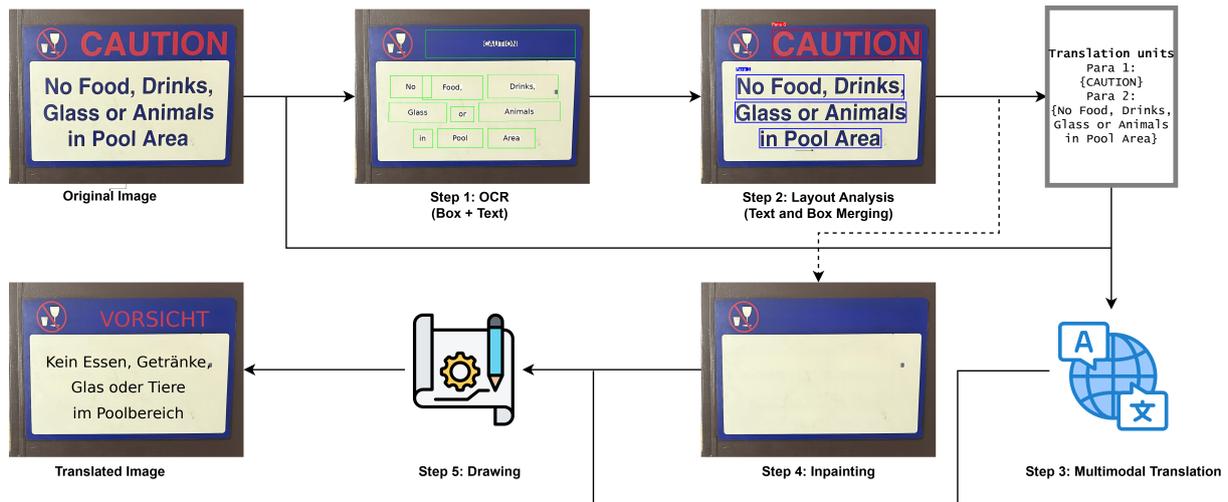


Figure 2: Overview of the image translator pipeline. Arrows indicate the inputs to each step. All steps are model-based except for drawing, which uses heuristic rules.

## 2.1 Multimodal Speech Translation

Several ST systems support translation across multiple languages, but they are not directly suitable for live lectures. Most are trained for offline tasks with fixed segmentation, which is incompatible with streaming audio, and require simultaneous translation policies (Niehues et al., 2016, 2018; Polák et al., 2023) to determine when enough audio has been received. Existing systems are either unimodal, ignoring slides, or multimodal but lack multilingual support. Lecture scenarios demand both multimodality and multilinguality.

To address these challenges, we adopt the OmniFusion model for multimodal ST, which supports multiple languages and has been shown to improve quality when integrating slides. Since it is trained primarily on clean speech, we fine-tune it on noisy data<sup>2</sup>. For streaming translation, we follow the LT policy (Huber et al., 2023), combining voice-activity detection with Local-Agreement to produce low-latency outputs.

Accurate visual context is crucial for effective translation. The PDF viewer tracks the slide displayed during each audio segment, allowing us to extract a screenshot from the middle of the segment and feed it to the ST model. This provides relevant visual cues, improving translation quality, especially for technical content, while enabling participants to follow translations in real time.

<sup>2</sup>[https://huggingface.co/skoneru/OmniFusion\\_v2](https://huggingface.co/skoneru/OmniFusion_v2)

### 2.1.1 Summarization & Question-Answering

Beyond translating spoken content, lecture material should also be chaptered (Zechner and Waibel, 2000a,b; Schneider et al., 2025; Retkowski and Waibel, 2024), meaning split into coherent functional and semantic sections, and then summarized in multiple languages and made available for interactive QA. To support these tasks, we use the transcribed multimodal ST output as context. Although modern LLMs can handle long contexts efficiently, their context window is still limited, so we adopt the following strategy.

For summarization, lectures are first translated into multiple languages. Each lecture is then divided into chapters, which prevents context-window overflow and also produces conceptually cleaner summaries, since chapters contain locally coherent content and avoid the topic drift that often appears in global summaries. For each chapter, we generate several forms of compressed representations. These include transcript compressions at multiple ratios such as 50 percent, 70 percent, and 90 percent, as well as length-controlled summaries whose size is determined by the length of the source section (Retkowski and Waibel, 2025). All summaries are first produced in English to benefit from the stronger performance of LLMs on English text and are then translated into the target languages.

For QA, we follow a similar approach: the English transcript, organized by chapters, is used with Retrieval-Augmented Generation (RAG) to query an LLM (Anderer et al., 2025), and the resulting answers are translated into the target languages.

Model	CER ( $\downarrow$ )	TER ( $\downarrow$ )	Sub.	Del.	Ins.	Average Time ( $\downarrow$ ) (Seconds)
EasyOCR	56.44	57.44	1488	29337	553	0.22
Paddle-OCR-v4	11.31	16.53	880	2791	2435	0.06
Paddle-OCR-v5	13.48	16.91	1717	2639	3014	0.10
Qwen-2.5-VL 7B	13.54	12.77	413	2348	3144	5.10

Table 1: Performance of OCR models on the VISTRA benchmark. Evaluations are restricted to English text in signboards and similar visual contexts, and therefore do not reflect performance across broader OCR domains.

## 2.2 Slide Translation

Another challenge for making lectures accessible is translating slides into multiple languages. Slides contain both editable text and images with embedded text. For editable text, we use a Python-based PowerPoint parser<sup>3</sup> to extract text blocks and translate them with standard unimodal MT, avoiding multimodal models due to computational cost.

Text inside images cannot be directly extracted, often lacks surrounding linguistic context, and relies on visual elements for interpretation, making multimodal translation necessary. After translation, text must be reinserted into the original image to preserve layout and visual meaning. To address this, we propose an **image-translation pipeline** that detects, recognizes, translates, and re-renders text within slide images (Figure 2).

### 2.2.1 Optical Character Recognition

The system begins with extracting text from slide images using PaddleOCR v5 (Cui et al., 2025), which supports multiple languages and outputs both recognized text and bounding boxes, typically at the word or character level. While sufficient for translation, these detections do not form coherent segments or preserve semantic structure, requiring layout analysis.

### 2.2.2 Layout Analysis

We then apply layout analysis using the Hi-SAM model<sup>4</sup> (Ye et al., 2025b), which predicts block-level regions and their constituent lines. OCR boxes are grouped into block-level and line-level segments, producing sentence-like units suitable for translation. Layout analysis also preserves structural cues, such as grouping, font size, and color, that aid re-rendering. For instance, bullet list items or diagram labels are grouped to maintain consistent formatting.

<sup>3</sup><https://pypi.org/project/python-pptx/>

<sup>4</sup>[sam\\_vit\\_l\\_0b3195.pth](https://github.com/inesmsahin/simple-lama-inpainting/)

### 2.2.3 Multimodal Translation

Text from each block is concatenated and translated using OmniFusion adapted from Qwen Omni 2.5 7B (Ye et al., 2025a) and SeedX PPO 7B (Cheng et al., 2025), which leverages the slide image as visual context. This multimodal approach is particularly helpful for short, ambiguous, or visually grounded text.

### 2.2.4 Inpainting

Before inserting the translated text, the original text regions are removed using Simple-LaMa<sup>5</sup> (Suvorov et al., 2021), a lightweight inpainting model that reconstructs the background with minimal artifacts, preserving slide quality.

### 2.2.5 Drawing

Translated text is then rendered back onto the slide. Fully automatic diffusion-based methods proved unsuitable because repeated edits gradually degraded clarity. Instead, a heuristic drawing module estimates original text styling and positions the translated text within the same layout and line structure. This preserves alignment, spatial organization, and overall visual coherence, ensuring the localized slide matches the structure and intent of the original.

## 3 Experiments

### 3.1 Evaluation Data & Metrics

Since no dataset directly provides lecture slides with ground-truth translations, summaries, and QA pairs, we evaluate our approach on established benchmarks that approximate these tasks. For image translation, we use the VISTRA benchmark (Salesky et al., 2024), which contains real-world images such as street signs with ground-truth OCR and translations for English  $\rightarrow$  German, Chinese,

<sup>5</sup><https://github.com/enesmsahin/simple-lama-inpainting/>

Model	de			es			ru			zh		
	BLEU (↑)	ChrF (↑)	COMET (↑)	BLEU (↑)	ChrF (↑)	COMET (↑)	BLEU (↑)	ChrF (↑)	COMET (↑)	BLEU (↑)	ChrF (↑)	COMET (↑)
<i>OCR Predicted + Line-level</i>												
SeedX 7B PPO	6.7	21.3	50.9	18.3	48.8	68.9	10.8	37.8*	65.6*	0.6	7.4	62.8
Tower-Instruct 7B	4.5	23.3	50.5	11.6	40.0	63.2	7.3	28.9	59.1	3.5*	17.2	63.1
OmniFusion	9.2*	25.3*	53.5*	19.8*	50.7*	70.4*	11.0*	34.8	64.6	1.3	22.1*	67.6*
<i>OCR Predicted + Layout-level</i>												
SeedX 7B PPO	10.3	23.7	53.1	28.4*	56.8*	74.0	17.3*	43.4*	71.2*	2.0	14.8	67.8
Tower-Instruct 7B	11.2	27.4	53.7	19.1	46.4	68.2	10.6	30.7	63.3	8.4*	22.9	68.3
OmniFusion	13.6*	30.1*	56.9*	28.1	56.2	74.5*	15.2	36.7	68.5	5.4	27.9*	71.4*
<i>Ground-Truth (OCR + Segmentation)</i>												
SeedX 7B PPO	14.5	27.4	57.8	35.6	<b>63.1*</b>	81.8	<b>23.5*</b>	<b>49.2*</b>	<b>78.9*</b>	13.9	34.5	83.4
Tower-Instruct 7B	11.0	31.6	59.2	28.1	53.2	75.4	15.1	34.4	69.2	<b>23.3*</b>	37.5	83.1
OmniFusion	<b>18.4*</b>	<b>35.0*</b>	<b>62.2*</b>	<b>36.9*</b>	62.5	<b>81.9*</b>	20.4	38.8	74.0	16.5	<b>43.5*</b>	<b>84.6*</b>

Table 2: Comparison of translation quality across models on the VISTRA benchmark. OCR-predicted results rely on PaddleOCR-v5. The best score within each evaluation setting is marked with \*, and the best overall is **bold**.

Russian, Spanish. OCR performance is measured using Character Error Rate (CER), Term Error Rate (TER; Snover et al. 2006), and latency. Translation quality is evaluated with BLEU, ChrF using SacreBLEU<sup>6</sup> (Post, 2018), and COMET<sup>7</sup> (Rei et al., 2022). For downstream tasks, we use the MCIF dataset of ACL talks (Papi et al., 2025b) and report normalized BERTScore to evaluate generated summaries and answers.

### 3.2 Image Translation

We evaluate our complete image-translation pipeline along three dimensions: OCR accuracy, translation quality, and component runtime.

**OCR Evaluation.** Table 1 summarizes OCR performance of several open-source systems and the vision LLM Qwen-2.5-VL (7B; Bai et al. 2025). EasyOCR<sup>8</sup> performs the worst due to its lightweight and less robust design. PaddleOCR v4 and v5 achieve similar and much higher accuracy, while Qwen-2.5-VL matches PaddleOCR but suffers from very high latency (0.1s → 5s per image). Considering accuracy, latency, and language coverage, PaddleOCR v5 provides the best trade-off and is used for all subsequent experiments.

**Translation Quality** Table 2 presents translation results for both unimodal LLMs, Tower 7B (Alves et al., 2024) and SeedX, and the multimodal OmniFusion model. To evaluate the impact of input segmentation, we compare line-level segmentation (where each OCR line is treated independently), block-level segmentation (where lines are grouped within layout regions), and ground-truth OCR plus segmentation as an upper bound.

Overall, OmniFusion consistently outperforms unimodal translation in most languages, showing that visual context from images helps disambiguate short or visually grounded text, such as diagram labels or signs. Ground-truth OCR and segmentation yield the best performance, highlighting the importance of accurate text extraction and layout grouping. Block-level segmentation improves translation over line-level segmentation, confirming that coherent sentence-like units are critical for high-quality output. Unimodal translation performs better in Russian, indicating potentially less reliance on visual context in this direction.

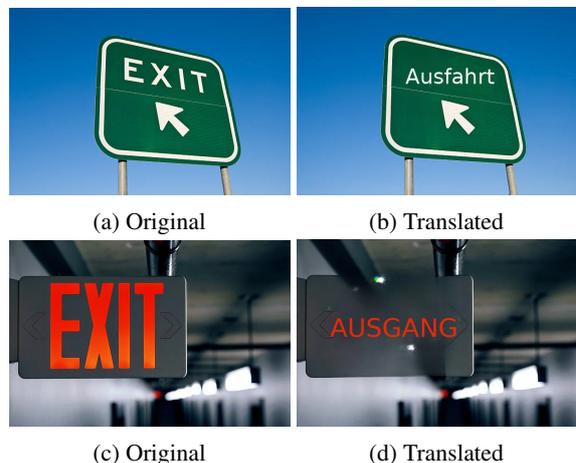


Figure 3: Example illustrating that our Image Translator uses context for disambiguation. The word “Exit” can mean “Ausgang” in the context of a pedestrian exit and “Ausfahrt” in the context of a car exit. Our translator correctly leverages the visual context to produce different translations, even when the source text is identical in both scenarios.

Table 4 in Appendix A shows inference times for different components. Layout analysis and translation are slowest, whereas OCR and image ren-

<sup>6</sup>nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1

<sup>7</sup>Unbabel/wmt22-comet-da

<sup>8</sup>https://github.com/JaidedAI/EasyOCR

Language	ST Input	LLaMA 3.1 8B	GPT OSS 20B	Mistral Small 3.2 24B
<b>Summarization</b>				
English	🔊	18.4	12.1	18.1
	🔊🖼️	20.5	12.7	19.7
German	🔊	20.6	18.0	21.7
	🔊🖼️	23.4	18.9	24.1
Italian	🔊	22.5	18.9	25.4
	🔊🖼️	24.4	19.7	26.3
Chinese	🔊	35.7	31.9	35.9
	🔊🖼️	35.3	31.7	35.8
<b>Question Answering</b>				
English	🔊	31.5	23.0	34.5
	🔊🖼️	34.5	22.0	35.4
German	🔊	32.0	21.5	37.2
	🔊🖼️	33.6	22.5	37.6
Italian	🔊	33.7	19.4	36.2
	🔊🖼️	34.7	20.5	34.7
Chinese	🔊	35.8	30.5	32.4
	🔊🖼️	35.4	30.0	32.7

Table 3: Summarization and Question Answering performance of different LLMs on the MCIF test dataset based on translations of the presentations with Omni-Fusion. Reported is BERTScore ( $\uparrow$ ), rescaled with the baseline. 🔊: Audio only, 🔊🖼️: Audio + Image.

dering add relatively minor overhead, suggesting that optimizing efficiency for these would provide the largest latency gains. Figure 3 illustrates an example in which multimodal translation disambiguates text using visual context, demonstrating the practical benefit of incorporating images.

### 3.3 Downstream Tasks

We analyze how downstream performance on the MCIF benchmark (Papi et al., 2025b), specifically for Summarization and Question Answering, is affected when the transcript used as context is generated by the multimodal speech-translation system. Using the task instructions provided by MCIF, we prompt each evaluated model directly with the translated talk transcript produced by our pipeline. We evaluate three LLMs: LLaMA 3.1 8B (Grattafiori et al., 2024), GPT-OSS 20B (OpenAI, 2025), and Mistral-Small 3.2 24B (Jiang et al., 2023)<sup>9 10</sup>. This setup allows us to measure how using audio-only transcripts compared to multimodal transcripts that also incorporate slide information influences downstream task performance.

**Summarization.** As shown in Table 3, summaries generated from audio+image input (🔊🖼️) consistently outperform those based on audio-only (🔊) across most languages and models, even

though the summarization models are text-only. The gains are most pronounced in English, German and Italian, while results for Chinese slightly degraded. We presume this is because English domain terminology appears in references for Latin-alphabet languages, while the lexical distance between English and Chinese prevents the models from consistently benefiting from additional context provided in English language.

**Question Answering.** In most settings, the results for QA are slightly better when incorporating visual context, though the gains are much less pronounced compared to summarization. In most cases, we observe small gains but also performance regression in four out of twelve language-model combinations. We assume that we do not see higher improvements and regression because the LLM does not receive the image data itself but just the (through multimodality improved) textual context which is not enough for the model to answer the questions more reliably.

## 4 Related Work

Streaming ST has been extensively studied in the last decade (Macháček et al., 2023; Guo et al., 2025; Papi et al., 2025a). Several lecture translation tools have also leveraged ST (Cho et al., 2013; Niehues et al., 2016; Son Nguyen et al., 2020; Müller et al., 2016; Dessloch et al., 2018; Huber et al., 2023), but these systems primarily rely on audio input. In contrast, our work extends lecture translation to multimodal input, incorporating visual cues from slides, and multimodal output, producing translated audio and slides in multiple languages.

Image-to-image translation remains relatively under-explored. Several research works focus on road sign translation (Gao et al., 2001; Yang et al., 2001; Zhang et al., 2002; Chen et al., 2002, 2004) facing many similar challenges to translating images in academic slides. Interest in this area is growing with the availability of larger datasets (Zuo et al., 2025; Li et al., 2025; Zhuang et al., 2025), but most existing work focuses solely on text translation within images, without addressing the aligned re-rendering of the visual content. An initial step in this direction is (Tian et al., 2025), which explicitly models the rendering process. Our image translation pipeline provides a modular foundation, enabling researchers to integrate models at any stage from OCR to translation and rendering, without

<sup>9</sup><https://mistral.ai/news/mistral-small-3-1>

<sup>10</sup><https://huggingface.co/mistralai/Mistral-Small-3.2-24B-Instruct-2506>

needing to implement additional components.

## 5 Conclusion

This paper presents a multimodal, multilingual lecture translation system that leverages multiple input modalities to generate translations across different output modalities. Future work includes conducting human evaluations to assess the quality of translated slides and audio, enabling targeted improvements to the system.

## Limitations

To assess the effectiveness of our slide translation, we use the VISTRA benchmark as a proxy. However, this benchmark does not fully reflect translation quality in the lecture domain, nor does it allow us to evaluate the quality of rendered slides. Human evaluation is therefore needed to assess the rendering quality of translated slides, including layout preservation and visual coherence. For SUM and QA, we conduct evaluation only after the entire talk has been translated, which does not accurately simulate a live lecture scenario. Benchmarks with questions aligned to the lecture timeline would provide more realistic and informative evaluations for our use-case.

## Acknowledgments

The research leading to these results was supported by European Union’s Horizon Europe programme grant agreement No. 101213369 (DVPS) and No. 101135798 (Meetween), The German Federal Ministry of Education, Research (BMBF) under the Robotics Institute Germany (RIG) and "How is AI Changing Science? Research in the Era of Learning Algorithms" (HiAICS) project.

## References

- Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G C de Souza, and André F T Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv [cs.CL]*.
- Katharina Anderer, Karin Müller, Lukas Strobel, Matthias Wölfel, Jan Niehues, and Kathrin Gerling. 2025. Making lecture videos accessible for students who are blind or have low vision through ai-assisted navigation and visual question answering. In *Proceedings of the 27th International ACM SIGACCESS*

*Conference on Computers and Accessibility, ASSETS '25*, New York, NY, USA. Association for Computing Machinery.

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, and 1 others. 2023. Seamless: Multilingual expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*.
- Xilin Chen, Jie Yang, Jing Zhang, and Alex Waibel. 2002. Automatic detection of signs with affine transformation. In *Sixth IEEE Workshop on Applications of Computer Vision, 2002.(WACV 2002). Proceedings.*, pages 32–36. IEEE.
- Xilin Chen, Jie Yang, Jing Zhang, and Alex Waibel. 2004. Automatic detection and recognition of signs from natural scenes. *IEEE Transactions on image processing*, 13(1):87–99.
- Zhe Chen, Heyang Liu, Wenyi Yu, Guangzhi Sun, Hongcheng Liu, Ji Wu, Chao Zhang, Yu Wang, and Yanfeng Wang. 2024. M<sup>3</sup>AV: A multimodal, multi-genre, and multipurpose audio-visual academic lecture dataset. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9041–9060, Bangkok, Thailand. Association for Computational Linguistics.
- Shanbo Cheng, Yu Bao, Qian Cao, Luyang Huang, Liyan Kang, Zhicheng Liu, Yu Lu, Wenhao Zhu, Jingwen Chen, Zhichao Huang, Tao Li, Yifu Li, Huiying Lin, Sitong Liu, Ningxin Peng, Shuaijie She, Lu Xu, Nuo Xu, Sen Yang, and 7 others. 2025. Seed-X: Building strong multilingual translation LLM with 7B parameters. *arXiv [cs.CL]*.
- Eunah Cho, Christian Fügen, Teresa Herrmann, Kevin Kilgour, Mohammed Mediani, Christian Mohr, Jan Niehues, Kay Rottmann, Christian Saam, Sebastian Stüker, and 1 others. 2013. A real-world system for simultaneous translation of german lectures. In *INTERSPEECH*, pages 3473–3477.
- Cheng Cui, Ting Sun, Manhui Lin, Tingquan Gao, Yubo Zhang, Jiakuan Liu, Xueqing Wang, Zelun Zhang, Changda Zhou, Hongen Liu, Yue Zhang, Wenyu Lv, Kui Huang, Yichao Zhang, Jing Zhang, Jun Zhang, Yi Liu, Dianhai Yu, and Yanjun Ma. 2025. Paddleocr 3.0 technical report. *Preprint*, arXiv:2507.05595.
- Florian Desseloch, Thanh-Le Ha, Markus Müller, Jan Niehues, Thai-Son Nguyen, Ngoc-Quan Pham, Elizabeth Salesky, Matthias Sperber, Sebastian Stüker, Thomas Zenkel, and 1 others. 2018. Kit lecture translator: Multilingual speech translation with one-shot learning. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 89–93.

- Yves Gambier. 2023. Audiovisual translation and multimodality: What future? *Media and intercultural communication: a multidisciplinary journal.*, 1(1):1–16.
- Jiang Gao, Jie Yang, Ying Zhang, and Alex Waibel. 2001. Text detection and translation from natural scenes. Technical report.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models. *arXiv [cs.AI]*.
- Shoutao Guo, Xiang Li, Mengge Liu, Wei Chen, and Yang Feng. 2025. Streamuni: Achieving streaming speech translation with a unified large speech-language model. *arXiv preprint arXiv:2507.07803*.
- Christian Huber, Tu Anh Dinh, Carlos Mullov, Ngoc-Quan Pham, Thai-Binh Nguyen, Fabian Retkowski, Stefan Constantin, Enes Ugan, Danni Liu, Zhaolin Li, and 1 others. 2023. End-to-end evaluation for low-latency simultaneous speech translation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 12–20.
- W John Hutchins. 2004. The georgetown-ibm experiment demonstrated in january 1954. In *Conference of the Association for Machine Translation in the Americas*, pages 102–114. Springer.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. **Mistral 7b**. *CoRR*, abs/2310.06825.
- Melvin Johnson, Mike Schuster, Quoc Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Vi  gas, Martin Wattenberg, Greg Corrado, and 1 others. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.
- Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *International Conference on Machine Learning*, pages 5530–5540. PMLR.
- Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*. Association for Computational Linguistics.
- Sai Koneru, Matthias Huck, and Jan Niehues. 2025. **Omnifusion: Simultaneous multilingual multimodal translations via modular fusion**. *Preprint*, arXiv:2512.00234.
- Jungil Kong, Jihoon Park, Beomjeong Kim, Jeongmin Kim, Dohee Kong, and Sangjin Kim. 2023. Vits2: Improving quality and efficiency of single-stage text-to-speech with adversarial learning and architecture design. *arXiv preprint arXiv:2307.16430*.
- Bo Li, Shaolin Zhu, and Lijie Wen. 2025. **MIT-10M: A large scale parallel corpus of multilingual image translation**. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5154–5167, Abu Dhabi, UAE. Association for Computational Linguistics.
- Dominik Mach  cek, Raj Dabre, and Ondr  j Bojar. 2023. Turning whisper into real-time transcription system. 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: System Demonstrations*. Association for Computational Linguistics (ACL).
- Markus M  ller, Sarah F  nfer, Sebastian St  ker, and Alex Waibel. 2016. Evaluation of the kit lecture translation system. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 1856–1861.
- Vimala Venugopal Muthuswamy and G Varshika. 2023. Analysing the influence of cultural distance and language barriers on academic performance among international students in higher education institutions. *Journal of International Students*, 13(3):415–440.
- Thai-Binh Nguyen, Ngoc-Quan Pham, and Alexander Waibel. 2025. Cocktail-party audio-visual speech recognition. In *Proc. Interspeech 2025*, pages 1828–1832.
- Jan Niehues, Thai Son Nguyen, Eunah Cho, Thanh-Le Ha, Kevin Kilgour, Markus M  ller, Matthias Sperber, Sebastian St  ker, and Alex Waibel. 2016. Dynamic transcription for low-latency speech translation. In *Interspeech*, pages 2513–2517.
- Jan Niehues, Ngoc-Quan Pham, Thanh-Le Ha, Matthias Sperber, and Alex Waibel. 2018. Low-latency neural speech translation. *arXiv preprint arXiv:1808.00491*.
- OpenAI. 2025. **gpt-oss-120b & gpt-oss-20b model card**. *CoRR*, abs/2508.10925.
- Sara Papi, Peter Pol  k, Dominik Mach  cek, and Ondr  j Bojar. 2025a. How “real” is your real-time simultaneous speech-to-text translation system? *Trans. Assoc. Comput. Linguist.*, 13:281–313.
- Sara Papi, Maike Z  fle, Marco Gaido, Beatrice Savoldi, Danni Liu, Ioannis Douros, Luisa Bentivogli, and Jan Niehues. 2025b. **Mcif: Multimodal crosslingual instruction-following benchmark from scientific talks**. *Preprint*, arXiv:2507.19634.

- Ngoc-Quan Pham, Thai-Son Nguyen, Jan Niehues, Markus Müller, Sebastian Stüker, and Alexander Waibel. 2019. Very deep self-attention networks for end-to-end speech recognition. *arXiv preprint arXiv:1904.13377*.
- Peter Polák, Brian Yan, Shinji Watanabe, Alex Waibel, and Ondřej Bojar. 2023. Incremental blockwise beam search for simultaneous speech translation with controllable quality-latency tradeoff. *arXiv preprint arXiv:2309.11379*.
- 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.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. [Robust speech recognition via large-scale weak supervision](#). *arXiv preprint*.
- Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. [COMET-22: Unbabel-IST 2022 submission for the metrics shared task](#). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Fabian Retkowsky and Alexander Waibel. 2024. [From text segmentation to smart chaptering: A novel benchmark for structuring video transcriptions](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 406–419, St. Julian’s, Malta. Association for Computational Linguistics.
- Fabian Retkowsky and Alexander Waibel. 2025. [Zero-shot strategies for length-controllable summarization](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 551–572, Albuquerque, New Mexico. Association for Computational Linguistics.
- Fabian Retkowsky, Maike Züfle, Andreas Sudmann, Dinah Pfau, Shinji Watanabe, Jan Niehues, and Alexander Waibel. 2025. [Summarizing speech: A comprehensive survey](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 27263–27294, Suzhou, China. Association for Computational Linguistics.
- Elizabeth Salesky, Philipp Koehn, and Matt Post. 2024. Benchmarking visually-situated translation of text in natural images. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1167–1182.
- Felix Schneider, Marco Turchi, and Alex Waibel. 2025. Policies and evaluation for online meeting summarization. *arXiv preprint arXiv:2502.03111*.
- Supriti Sinhamahapatra and Jan Niehues. 2025. Do slides help? multi-modal context for automatic transcription of conference talks. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 16111–16121, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231.
- Thai Son Nguyen, Jan Niehues, Eunah Cho, Thanh-Le Ha, Kevin Kilgour, Markus Muller, Matthias Sperber, Sebastian Stueker, and Alex Waibel. 2020. Low latency asr for simultaneous speech translation. *arXiv e-prints*, pages arXiv–2003.
- Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. 2021. Resolution-robust large mask inpainting with fourier convolutions. *arXiv preprint arXiv:2109.07161*.
- Yanzhi Tian, Zeming Liu, Zhengyang Liu, Chong Feng, Xin Li, He-Yan Huang, and Yuhang Guo. 2025. Prim: Towards practical in-image multilingual machine translation. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 13693–13708.
- 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.
- Alexander Waibel. 2014. Translation and integration of presentation materials in cross-lingual lecture support. US Patent App. 14/302,149.
- Alexander Waibel. 2018. Translation training with cross-lingual multi-media support. US Patent 9,892,115.
- Alexander Waibel and Christian Fuegen. 2012. Simultaneous translation of open domain lectures and speeches. US Patent 8,090,570.
- Jie Yang, Jiang Gao, Ying Zhang, Xilin Chen, and Alex Waibel. 2001. An automatic sign recognition and translation system. In *Proceedings of the 2001 workshop on Perceptive user interfaces*, pages 1–8.
- Hanrong Ye, Chao-Han Huck Yang, Arushi Goel, Wei Huang, Ligeng Zhu, Yuanhang Su, Sean Lin, An-Chieh Cheng, Zhen Wan, Jinchuan Tian, and 1 others. 2025a. Omnivinci: Enhancing architecture and data for omni-modal understanding llm. *arXiv preprint arXiv:2510.15870*.
- Maoyuan Ye, Jing Zhang, Juhua Liu, Chenyu Liu, Bao-cai Yin, Cong Liu, Bo Du, and Dacheng Tao. 2025b.

Hi-sam: Marrying segment anything model for hierarchical text segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(03):1431–1447.

Klaus Zechner and Alex Waibel. 2000a. Diasumm: Flexible summarization of spontaneous dialogues in unrestricted domains. In *COLING 2000 Volume 2: The 18th International Conference on Computational Linguistics*.

Klaus Zechner and Alex Waibel. 2000b. Minimizing word error rate in textual summaries of spoken language. In *1st Meeting of the North American Chapter of the Association for Computational Linguistics*.

Ying Zhang, Bing Zhao, Jie Yang, and Alex Waibel. 2002. Automatic sign translation. In *INTER-SPEECH*, pages 645–648.

Wanru Zhuang, Wenbo Li, Zhibin Lan, Xu Han, Peng Li, and Jinsong Su. 2025. Patimt-bench: A multi-scenario benchmark for position-aware text image machine translation in large vision-language models. *arXiv preprint arXiv:2509.12278*.

Fei Zuo, Kehai Chen, Yu Zhang, Zhengshan Xue, and Min Zhang. 2025. [InImageTrans: Multimodal LLM-based text image machine translation](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 20256–20277, Vienna, Austria. Association for Computational Linguistics.

## A.2 User Interface Screenshots

## A Appendix

Step	Time (seconds)
OCR	0.46
Layout Analysis	2.93
Multimodal Translation	3.10
Inpainting	0.42
Drawing	0.18

Table 4: Inference time for each step in the pipeline for translating the image shown in Figure 2.

### A.1 Text-to-Speech

In Figure 5, the interface of the TTS output can be seen. It is possible to select between the simultaneous and consecutive modes. The simultaneous mode can be used when listening to the TTS output via headphones during the talk. The consecutive mode is suitable in dialog scenarios where the TTS output is paused as long as the system recognizes speaker input. We use the VITS/VITS2 (Kim et al., 2021; Kong et al., 2023) and Kokoro-82M to generate audio together with a rule-based streaming algorithm to segment input text into segments.

## (a) English translation with segmentation into multiple chapters.

The screenshot shows a video player interface for the English translation. On the left, there is a 'Table of Contents' sidebar with a list of chapters: Intro, Intelligence, Mathematics, Comparing Computers to Human Beings, Recognition, Heuristics, and Assigning Credit. The main content area is titled 'English' and displays the first chapter, '1 Intro', which is currently expanded. The text of the 'Intro' chapter is visible, starting with 'This is a Q & A excerpt on the topic of artificial intelligence from a lecture by Richard Feynman from September 26, 1985. I re-recorded the audience questions because they're barely audible in the original. Question: Do you think there will ever be a machine that will think like human beings and be more intelligent than human beings?'. Below the text, there are seven buttons labeled 'Show' corresponding to the other chapters in the table of contents. The interface includes a play button on the left and a close button on the right.

## (b) German translation with segmentation into multiple chapters.

The screenshot shows a video player interface for the German translation. On the left, there is a 'Table of Contents' sidebar with a list of chapters: Einleitung, Intelligenz, Mathematik, Vergleich von Computern und Menschen, Anerkennung, Heuristik, and Zuweisung von Gutschriften. The main content area is titled 'German' and displays the first chapter, '1 Einleitung', which is currently expanded. The text of the 'Einleitung' chapter is visible, starting with 'Dies ist ein Frage-Antwort-Ausschnitt zum Thema künstliche Intelligenz aus einem Vortrag von Richard Feynman vom 26. September 1985. Ich habe die Fragen des Publikums neu aufgenommen, weil sie in der Originalaufnahme kaum zu hören sind. Frage: Glauben Sie, dass es jemals eine Maschine geben wird, die wie Menschen denkt und intelligenter ist als Menschen?'. Below the text, there are seven buttons labeled 'Show' corresponding to the other chapters in the table of contents. The interface includes a play button on the left and a close button on the right.

Figure 4: Translations of the YouTube video “Richard Feynman: Can Machines Think?” (<https://www.youtube.com/watch?v=ipRvjs7q1DI>). Subfigure (a) shows the English version; subfigure (b) shows the German version.

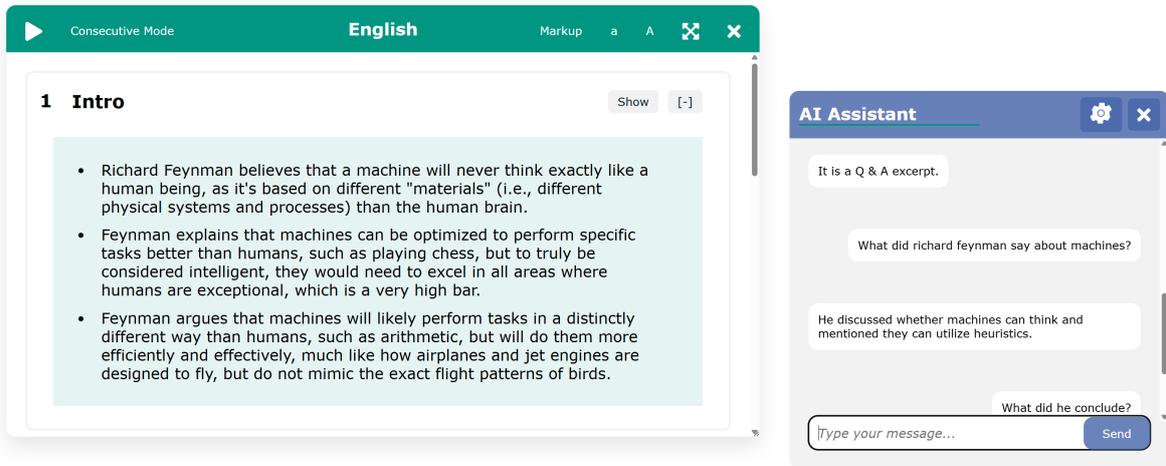


Figure 5: Summarization and Question Answering user interface. The summaries are shown for each chapter in all languages.



Figure 6: Slide viewer interface with multilingual navigation options. Users can switch between languages, browse slides independently of the presenter through an out-of-sync mode, and subsequently use the sync toggle to realign with the live presentation.

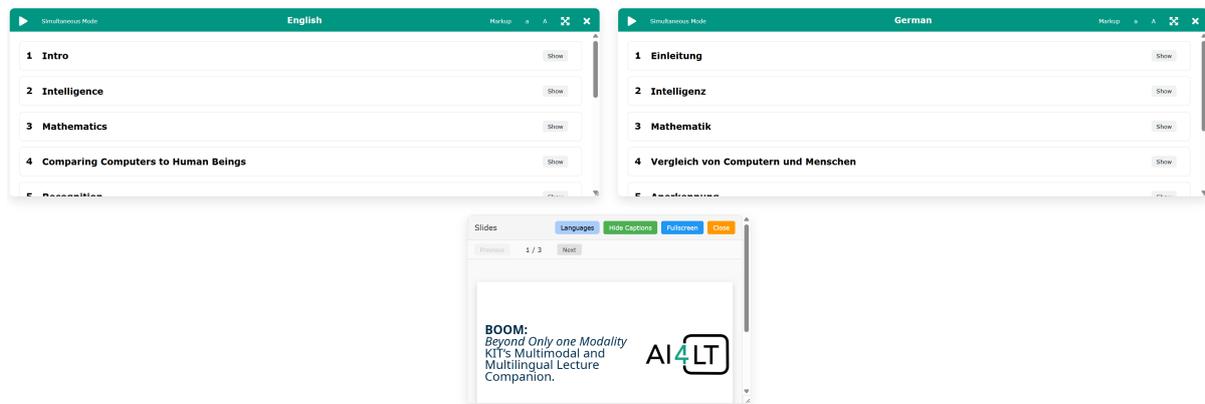


Figure 7: The interface also allows to see the translations in multiple languages along with the current slide.



Figure 8: Participant full screen view of the slide interface showing slides with caption overlay in German.

# PEFT-FACTORY: Unified Parameter-Efficient Fine-Tuning of Autoregressive Large Language Models

Robert Belanec<sup>♣†</sup>, Ivan Srba<sup>†</sup>, Maria Bielikova<sup>†</sup>

<sup>♣</sup> Faculty of Information Technology, Brno University of Technology, Brno, Czechia

<sup>†</sup> Kempelen Institute of Intelligent Technologies, Bratislava, Slovakia

{name.surname}@kinit.sk

Demonstration video: <https://youtu.be/Q3kxvly0-XY>

Installation package: <https://pypi.org/project/peftfactory>

Live demo: <https://peftfactory.kinit.sk>

## Abstract

Parameter-Efficient Fine-Tuning (PEFT) methods address the increasing size of Large Language Models (LLMs). Currently, many newly introduced PEFT methods are challenging to replicate, deploy, or compare with one another. To address this, we introduce PEFT-FACTORY, a unified framework for efficient fine-tuning LLMs using both off-the-shelf and custom PEFT methods. While its modular design supports extensibility, it natively provides a representative set of 19 PEFT methods, 27 classification and text generation datasets addressing 12 tasks, and both standard and PEFT-specific evaluation metrics. As a result, PEFT-FACTORY provides a ready-to-use, controlled, and stable environment, improving replicability and benchmarking of PEFT methods. PEFT-FACTORY is a downstream framework that originates from the popular LLaMA-Factory, and is publicly available at <https://github.com/kinit-sk/PEFT-Factory>.

## 1 Introduction

Large Language Models (LLMs) (Minaee et al., 2024; Radford et al., 2019; Dubey et al., 2024; Raffel et al., 2020) achieved remarkable results in many different Natural Language Processing (NLP) tasks mainly after the introduction of the transformer architecture (Vaswani, 2017). However, for its great scaling capabilities (Kaplan et al., 2020), the size of the model in terms of trainable parameters is continuously increasing accordingly. The growing number of LLM parameters rendered fine-tuning computationally expensive, data-hungry, and hardly accessible for many researchers and practitioners.

Parameter-Efficient Fine-Tuning (PEFT) methods (Xu et al., 2023; Ding et al., 2023; Lialin et al., 2023; Han et al., 2024) aim to address these issues by training only a small percentage of the full model’s parameters while achieving performance

comparable to that of full model fine-tuning. Such a decrease in trainable parameters can be achieved by adding new parameters (Houlsby et al., 2019), selecting specific parameters (Ben Zaken et al., 2022) for training, or by reparameterizing the model with a smaller number of parameters (Hu et al., 2022).

Due to their effectiveness, PEFT methods have gained popularity and become an attractive research area, with many new contributions being introduced each year (Han et al., 2024). However, the large number of newly introduced PEFT methods makes it harder to compare their contributions, resulting in only a few established PEFT methods (LoRA variants in most cases) being used in practice, while others, which may be more effective, remain largely unused. Moreover, many new PEFT methods lack a fully functional open-source implementation, and essential details on the experimental setup often prevent fellow researchers from replicating their results. Therefore, many researchers (Asai et al., 2022; Shi and Lipani, 2024; Tang et al., 2025) have to rely on reported performance metrics, while there is a risk of not replicating exactly the same experimental setup for their own methods (making a comparison potentially unfair), as well as it is not feasible to rerun the existing solutions on additional datasets/tasks. Lastly, many established PEFT methods lack proper evaluation on autoregressive LLMs.

To tackle these accumulating problems, we **introduce PEFT-FACTORY, an easy-to-use and modular framework for efficient fine-tuning and evaluation of LLMs using different PEFT methods**. PEFT-FACTORY is based on the popular and open-source fine-tuning framework *LLaMA-Factory* (Zheng et al., 2024b). It is built using PyTorch (Paszke et al., 2019) and utilizes open-source Python modules for training LLMs, including Transformers (Wolf et al., 2020), PEFT (Mangrulkar et al., 2022), TRL (von Werra et al., 2020), and Adapters (Poth et al., 2023).

Our main contributions are as follows:

- PEFT-FACTORY provides a support for off-the-shelf methods from popular PEFT provider frameworks like *HuggingFace PEFT* (Mangrulkar et al., 2022) or *Adapters* (Poth et al., 2023) as well as dynamic loading of *custom user-created PEFT methods*. In contrast to the existing solutions, it provides so-far-missing support for *soft prompt-based*, *adapter-based* and *selective* PEFT methods; as well as for *classification* tasks.
- PEFT-FACTORY natively provides a representative set of *19 PEFT methods*, *27 classification and text generation datasets* addressing *12 unique tasks*, and standard as well as PEFT-specific *evaluation metrics*. This ready-to-use setup enables quick adoption and experimentation by researchers and practitioners, significantly improving the currently limited replicability and benchmarking of PEFT methods.
- PEFT-FACTORY is designed with future *extensibility* in mind and provides a fully open-source codebase for anyone to use. It implements a standardized PEFT interface to enable modular addition of newly created PEFT methods. Similarly, it allows easy extension for additional datasets.

## 2 Related Work

There has been a significant rise in the number of frameworks used for training of (not only) LLMs. LLaMA-Factory (Zheng et al., 2024b) is a recent addition to such frameworks and offers end-to-end and easy-to-use training of LLMs ranging across all of the stages (from pre-training to alignment via reinforcement learning). LLaMA-Factory also provides a graphical user interface called LLaMABoard, implemented in Gradio (Abid et al., 2019), which enhances the ease of use of LLaMA-Factory. Despite being a really popular and useful tool for LLM training, LLaMA-Factory still provides fine-tuning only with LoRA (Hu et al., 2022) and its variants, namely QLoRA (Detters et al., 2023), DoRA (Liu et al., 2024), LoRA+ (Hayou et al., 2024), PiSSA (Meng et al., 2024), and Galore (Zhao et al., 2024b). With the recent update, LLaMA-Factory also allows Orthogonal Fine-Tuning (OFT) (Qiu et al., 2023), which utilizes

the Cayley transformation (Cayley, 1846) to fine-tune only orthogonal vectors. Nevertheless, the selection of PEFT methods in LLaMA-Factory still remains limited. Lastly, LLaMA-Factory primarily focuses on text-generation problems and does not incorporate the possibility of casting text-generation problems as classification tasks. Our framework PEFT-FACTORY addresses both the limited number of available PEFT methods and the potential for fine-tuning LLMs for classification.

There are also other LLM training frameworks that are less easy to run (compared to LLaMA-Factory) and have their specific benefits. FastChat (Zheng et al., 2023) is a specialized framework for training LLMs for chat-completion. LitGPT (AI, 2023) and LMFlow (Diao et al., 2023) are extensible and convenient general training frameworks that support various generative models and training methods. Axolotl (Axolotl maintainers and contributors, 2023) is a terminal-based tool for efficient post-training of LLMs without sacrificing functionality or scale. Open-Instruct (Wang et al., 2023a) focuses on instruction fine-tuning for LLMs and provides multiple models and recipes for this purpose. H2O LLM Studio<sup>1</sup> is a more enterprise-oriented, all-in-one tool that also provides a graphical interface for developing and deploying LLM models. GPT4All (Anand et al., 2023) creates a user-friendly interface around llamacpp. ColossalAI (Li et al., 2023) focuses on delivering a framework for distributed fine-tuning.

In addition, LLaMA-Adapter (Zhang et al., 2023) and LLaMA-Accessory (Gao et al., 2023) are more lightweight frameworks, where LLaMA-Adapter adds trainable adapters to (not only) LLaMA models and LLaMA-Accessory provides a full toolkit for LLM development. LLaMA-Adapter is often implemented in previously-mentioned frameworks, such as LitGPT. Table 1 provides a summary of unique PEFT-FACTORY features when compared with popular fine-tuning frameworks as well as our upstream framework LLaMA-Factory. Based on our analysis of related frameworks and to the best of our knowledge, we have identified 3 key features that are currently missing or limited, and are novel in our work: 1) training of LLMs with other than reparametrization-based PEFT methods, 2) modular and dynamic addition of new PEFT methods, and 3) support for training and evaluation of LLMs for classification.

<sup>1</sup><https://github.com/h2oai>

	Reparameterized	Soft Prompt-Based	Adapter-Based	Selective	Classification Datasets	Classification Metrics	Extensibility
<b>PEFT-FACTORY</b>	8	5	4	2	✓	✓	datasets, models, PEFT methods
LLaMA-Factory	7	0	0	0	✗	✗	datasets, models
FastChat	3	0	0	0	✗	✗	datasets, models
LitGPT	2	0	0	0	✓	✓	datasets, models
LMFlow	3	0	0	0	✗	✗	datasets, models
Axolotl	2	0	0	0	✗	✗	datasets, models
Open-Instruct	3	0	0	0	✗	✗	✗
H2O LLM Studio	3	0	0	0	✓	✓	datasets, models
GPT4All	2	0	0	0	✗	✗	✗

Table 1: Comparison of PEFT-FACTORY to popular fine-tuning frameworks. Only PEFT-FACTORY allows for out-of-the-box non-reparameterization efficient fine-tuning with the extensibility of additional and custom fine-tuning methods. Comparison at the level of individual PEFT methods can be found in Table 3 of Appendix B.

### 3 PEFT-FACTORY

The PEFT-FACTORY consists of four main components: 1) PEFT Methods, 2) Datasets, 3) Models, and 4) Metrics, as also depicted in Figure 1.

In the *PEFT methods* component, we design and implement support for reparameterized, soft prompt-based, adapter-based, and selective PEFT methods, from HuggingFace PEFT (Mangrulkar et al., 2022) and Adapters (Poht et al., 2023) PEFT provider frameworks. We also provide a custom PEFT interface for more advanced users to provide and dynamically load their custom PEFT methods into PEFT-FACTORY. Currently, we include 19 different PEFT methods (out of them, 7 are natively provided by the LLaMA-Factory). Full listing of PEFT methods covered by PEFT-FACTORY can be found in Table 3 of Appendix B.

The core of the *Datasets* component is the dataset loader supporting datasets from classification tasks, with the possibility of adding separate instructions for instruction fine-tuned models (a missing feature of LLaMA-Factory). Additionally, we include and adapt multiple well-known classification benchmarks, as well as text-generation tasks, totalling 27 datasets.

Regarding *Models*, PEFT-FACTORY leverages the existing support provided by LLaMA-Factory. It enables the utilization of a wide range of models from different model families, spanning from 0.5 (e.g., Qwen 2.5 (Yang et al., 2024)) to 671 (e.g., DeepSeek R1 (Guo et al., 2025)) billion parameters. For demonstration purposes, we selected Llama-3.2-1B-Instruct (Dubey et al., 2024) as it is a popular representative of a reasonable size that allows fast training to demonstrate PEFT-FACTORY.

Within the *Metrics* component, we add classification and performance-based metrics into the evaluation of LLMs trained using PEFT methods. This includes the addition of standard classification metrics, such as accuracy and F1, as well as the

PSCP metric (Belanec et al., 2025), which incorporates various efficiency factors into the results.

#### 3.1 Off-The-Shelf PEFT Support

There are many different PEFT methods included off-the-shelf within the PEFT provider libraries, such as HuggingFace PEFT and Adapters. In the current state, we include 10 different off-the-shelf PEFT methods that we have tested with different state-of-the-art LLMs, namely, from the **Adapters library** – *Parallel Adapter* (He et al., 2022), *Bottleneck Adapter* (Houlsby et al., 2019), and *Sequential Bottleneck Adapter* (Pfeiffer et al., 2020); and from the **HuggingFace PEFT** – *Prompt Tuning* (Lester et al., 2021), *Prefix Tuning* (Li and Liang, 2021), *P-Tuning* (Liu et al., 2023), *P-Tuningv2* (Liu et al., 2022b), *MTP* (Wang et al., 2023b), *LNTuning* (Zhao et al., 2024a), and *IA<sup>3</sup>* (Liu et al., 2022a). Moreover, PEFT-FACTORY enables to easily *add more of such off-the-shelf PEFT methods simply by updating two constants in the code*.

From the implementation perspective, to include support for these libraries in PEFT-FACTORY, we created a unified `PeftArguments` class that inherits both the `PEFTConfig` and `AdapterConfig` classes for typing purposes. This joint class is then used for parsing the parameters from configurations via `HFArgumentParser`. We store all supported off-the-shelf PEFT methods in several constants, specifically lists `HF_PEFT_METHODS` and `ADAPTERS_METHODS`, along with their counterpart mapping dictionary constants `PEFT_CONFIG_MAPPING` and `ADAPTERS_CONFIG_MAPPING`. If the PEFT method in the configuration is contained within the mapping dictionaries, a specific config is used. Otherwise, it will default to `PeftArguments` (in file `hparams/parser.py`, function `_parse_train_args`). Importantly, every PEFT method, whether added from the Adapters

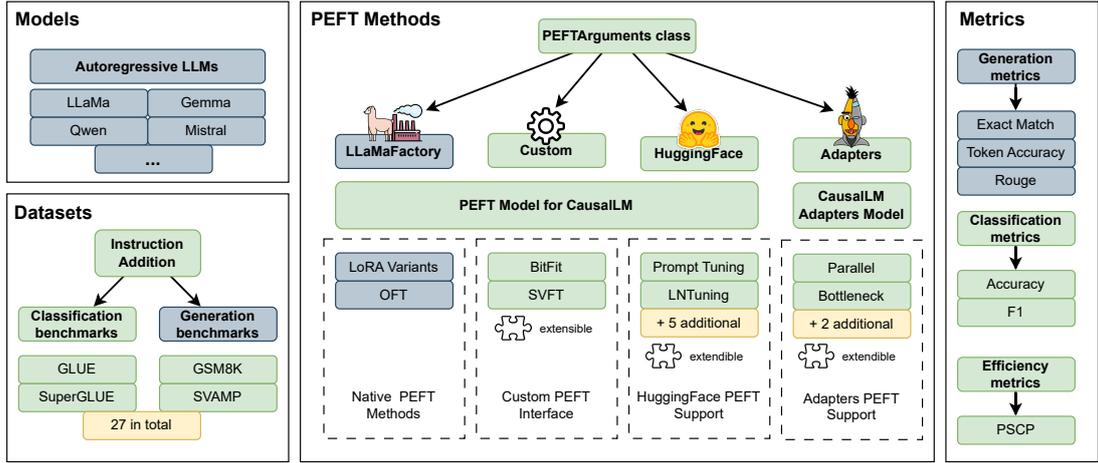


Figure 1: Diagram representing the components of PEFT-FACTORY. The four main overarching components of PEFT-FACTORY are PEFT Methods, Datasets, Models, and Metrics, which are further defined by their subcomponents. Components represented by green color are implemented in PEFT-FACTORY, components in blue color are native to LLaMA-Factory (Zheng et al., 2024a). Additionally, the Adapters library requires a different model class than the rest of the PEFT provider frameworks.

library or the HuggingFace PEFT library, comes with its own set of parameters or hyperparameters that can be tuned. PEFT-FACTORY automatically detects the parameters required by a specific PEFT method and uses the `HFAArgumentParser` to parse them from config files (if parameters are not specified in the config, default values are used from the original implementation). This allows easy configuration and hyperparameter tuning via a YAML config file or the Gradio user interface.

The newly created and parsed `PeftArguments` are then forwarded to the model loader, where the `init_adapter` method creates the PEFT model (in file `model/model_utils/adapter.py`, functions `_setup_custom_peft`, `_setup_adapters_peft` and `_setup_hf_peft`). From this part, we leave the loading, training, and model saving on PEFT libraries and the LLaMA-Factory framework.

### 3.2 Custom PEFT Interface

In contrast to LLaMA-Factory, PEFT-FACTORY implements a dynamic loading mechanism for custom PEFT methods, ultimately enabling its extensibility and modularity. This design allows researchers and practitioners to seamlessly integrate custom PEFT implementations without modifying the core codebase of PEFT-FACTORY. To demonstrate our Custom PEFT Interface, we replicated 2 PEFT methods that are not a part of any off-the-shelf PEFT framework, namely BitFit (Ben Zaken et al., 2022) and SVFT (Lingam et al., 2024) (located in the `peft` directory).

During the process of dynamic loading, the `peft_loader` module automatically discovers and loads PEFT methods from a structured directory hierarchy (in file `extras/peft_loader.py`, function `discover_custom_peft_methods`). Each custom PEFT method is organized in its own subdirectory containing two essential components: a `config.py` file defining a `PeftConfig` subclass, and a `model.py` file implementing a `BaseTuner` subclass. The configuration and model subclasses need to inherit the `PeftConfig` dataclass and `BaseTuner` abstract class from the HuggingFace PEFT library. The methods required to be implemented are then specified by the `BaseTuner` description in `tuner_utils.py`.

The loader validates each implementation by checking for required attributes (`peft_type` for configurations and `prefix` for model classes) before registration. The loader dynamically loads the config and model subclasses, registering them via the `register_peft_method` function (in file `peft_loader.py`), which adds the config and model to the constants of the Hugging Face PEFT library. Additionally, this process is also defined by the Algorithm 1, which explains how dynamic loading is implemented.

After the dynamic loading, the corresponding custom PEFT method is handled similarly to off-the-shelf methods as described in Section 3.1. As a result, the ease of configuration and hyperparameter tuning for the newly added PEFT methods also remains unchanged.

---

**Algorithm 1** Dynamic PEFT Method Discovery

---

```
1: Input: PEFT directory path  $D$ 
2: Output: Dictionary  $M$  mapping method
   names to (config, model) tuples
3:  $M \leftarrow \emptyset$ 
4: for each subdirectory  $d$  in  $D$  do
5:   if  $d$  contains config.py and model.py
     then
6:     Load config class  $C$  from config.py
7:     Load model class  $T$  from model.py
8:     if  $C$  validates and  $T$  validates then
9:        $M[\text{name}(d)] \leftarrow (C, T)$ 
10:    end if
11:  end if
12: end for
```

---

To add a new method, it is required to match the directory structure inside the *PEFT methods directory* (the directory can be specified by the environment variable `PEFT_DIR` with `./peft` as the default directory) to match the organizational structure and class inheritance. We provide information about example method templates in Appendix C as well as in the PEFT-FACTORY documentation<sup>2</sup>.

This plugin-style architecture promotes code reusability and enables fast prototyping of novel PEFT methods. Researchers can develop and test new methods separately, with the framework automatically integrating them at runtime.

The dynamic loading approach has proven particularly valuable for comparative studies, allowing researchers to evaluate multiple PEFT variants under identical experimental conditions without code duplication or version control conflicts.

### 3.3 Improved Dataset Loader

Besides a native support of text generation tasks (inherited from the LLaMA-Factory), we add support for classification tasks (in case of autoregressive models, the classification task  $Pr_{\theta}(y|X)$  is cast as a generation  $Pr_{\theta}(Y|X)$  task). To this end, we adapt and include multiple classification benchmarks, including GLUE (Wang et al., 2018) and SuperGLUE (Wang et al., 2019). This was largely possible due to the dataset loading feature included in LLaMA-Factory. However, we needed to enhance the dataset loader with an additional parameter called *instruction*. This optional attribute was added to the *Dataset-*

---

<sup>2</sup>For detailed information, please visit the [Adding PEFT Methods](#) section of the PEFT-FACTORY documentation.

*tAttr* class (in file `data/processor/parser.py`), and during data preprocessing, the instruction is prepended to the input text for the LLM (in file `data/processor/converter.py`, class `AlpacaDatasetConverter`). This allows adding instructions that have dataset-specific formatting (e.g., ones recommended by the dataset authors)<sup>3</sup> or are designed for instruction tuning tasks.

**Addition of classification datasets.** From the GLUE benchmark, we include 8 classification datasets separated into 6 tasks, namely **natural language inference (NLI)** – *MNLI* (Williams et al., 2018), *QNLI* (Rajpurkar et al., 2016), *RTE* (Dagan et al., 2005; Bar Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009); **paraphrase classification** – *QQP*<sup>4</sup>, *MRPC* (Dolan and Brockett, 2005); **sentiment classification** – *SST-2* (Socher et al., 2013); **sentence similarity** – *STS-B* (Cer et al., 2017) and **acceptability classification** – *CoLA* (Warstadt et al., 2019).

From SuperGLUE, we include 7 datasets separated into 4 tasks, namely **natural language inference (NLI)** – *CB* (De Marneffe et al., 2019); **question answering** – *MultiRC* (Khashabi et al., 2018), *ReCoRD* (Zhang et al., 2018), *BoolQ* (Clark et al., 2019), *COPA* (Roemmele et al., 2011); **word sense disambiguation** – *WiC* (Pilehvar and Camacho-Collados, 2019) and **coreference resolution** – *WSC* (Levesque et al., 2011).

**Addition of generation datasets.** We also include generation datasets that are commonly used to benchmark generative LLMs. We cover 6 datasets for reasoning and natural language understanding separated into 3 tasks, namely **question answering** – *MMLU* (Hendrycks et al., 2021), *PIQA* (Bisk et al., 2020), *SIQA* (Sap et al., 2019), *OBQA* (Khot et al., 2019); **natural language inference (NLI)** – *HellaSwag* (Zellers et al., 2019); **commonsense reasoning** – *WinoGrande* (Sakaguchi et al., 2021); 3 datasets for mathematical problem solving separated into 3 tasks, namely **question answering** – *MathQA* (Amini et al., 2019); **math word problems** – *GSM8K* (Cobbe et al., 2021) and **simple math problems** – *SVAMP* (Patel et al., 2021); and 3 datasets for code generation, namely *Conala* (Yin et al., 2018), *CodeAlpacaPy* (Chaudhary, 2023), and *APPS* (Hendrycks et al., 2021).

---

<sup>3</sup>For detailed information, please visit the [Adding Datasets](#) section of the PEFT-FACTORY documentation.

<sup>4</sup><https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs>

**Adapting and preprocessing datasets.** Some datasets may require further (mostly minor) format changes to be compatible with the input formatting of PEFT-FACTORY. We further describe this preprocessing in the Appendix B.1.

### 3.4 Classification and Efficiency Metrics

PEFT-FACTORY also calculates classification and efficiency metrics during the prediction phase in addition to already existing token accuracy and semantic similarity metrics (i.e., Rouge (Lin, 2004) and Bleu (Papineni et al., 2002)). From the classification metrics, PEFT-FACTORY implements standard Accuracy and F1 metrics. To measure efficiency during the evaluation of PEFT methods, PEFT-FACTORY implements the PSCP metric (Belanec et al., 2025), a highly adjustable metric that considers various efficiency parameters (e.g., number of parameters, memory usage, inference time).

We implement these metrics within the `train/sft/metric.py` file for supervised fine-tuning, following the pattern from LLaMA-Factory and utilizing separate data classes for each metric. Specifically, we name the classes `ComputeAccuracy`, `ComputeF1`, and `ComputePSCP`. For classification, we include a binary flag attribute within the training arguments, called `compute_classification_metrics`, which enables or disables the computation of classification metrics. For the efficiency metrics, we include a binary flag `compute_pscp`. Additional information on the usage of efficiency metrics can be found in Appendix B.2.

## 4 PEFT-FACTORY-enabled Use Cases

The extensibility of PEFT methods and datasets, together with a ready-to-use, controlled, and stable environment, is a key factor of PEFT-FACTORY that aims to promote further research on PEFT methods. To demonstrate how PEFT-FACTORY improves reproducibility and benchmarking of PEFT methods, we present two specific use cases.

### 4.1 PEFT Methods Reproducibility

How the modular design of PEFT-FACTORY promotes reproducibility and transparency of newly created PEFT methods can be seen in the use case, when fellow researchers and practitioners develop new PEFT methods.

Currently, when a new PEFT method is developed, the published source code is often not fully functional or difficult to reproduce. In addition, the

authors often have to implement code for training and evaluation of the PEFT method from scratch, which is often repetitive, increases the probability of mistakes in the code, and is prone to inconsistencies in the final results.

In our scenario, authors only need to create a minimum number of files that are directly and solely connected to the design of the PEFT method itself. If the authors maintain the structure compatible with the PEFT-FACTORY custom PEFT interface, they can simply share it within the PEFT methods directory, create a configuration for training and evaluation, and run experiments on vast amounts of datasets and autoregressive models. Additionally, if the authors choose to implement their method inside any of the supported PEFT provider frameworks (i.e., Hugging Face PEFT or Adapters), only a small change is needed to contribute it to the next version of the PEFT-FACTORY<sup>5</sup>.

### 4.2 PEFT Methods Benchmarks

Another possible use case of PEFT-FACTORY is to benchmark PEFT methods. To this end, PEFT-Factory provides a standardized and reproducible environment that eliminates inconsistencies in experimental setups (e.g., different seeds, hyperparameters or dataset splits), allowing researchers to reliably compare PEFT methods under identical conditions. To illustrate the benchmarking capability, Table 2 provides results from fine-tuning three PEFT methods on four different datasets using the LLaMA-3.2-1B-Instruct (Dubey et al., 2024) autoregressive model. Even such a small demonstrative comparison would require significant codebase preparation to execute the experiments, which PEFT-FACTORY eliminates to a minimum. From this evaluation, we can see that BitFit achieves the highest results in most of the datasets.

As a more complex benchmarking use case, we refer to our parallel work (Belanec et al., 2025), in which we introduce the PEFT-Bench – a benchmark of the efficiency of PEFT methods fully conducted within PEFT-FACTORY. PEFT-Bench provides the first unified, end-to-end benchmarking suite for evaluating PEFT methods on modern autoregressive LLMs, covering 27 datasets, 12 task types, and 7 diverse PEFT techniques. This benchmark was only possible due to PEFT-FACTORY, which serves as the underlying engine.

---

<sup>5</sup>We provide information on how to request the addition of a new PEFT provider method in the [Contributing page](#) of PEFT-FACTORY documentation

Method	SST-2	CoLA	WSC	SVAMP
BitFit	<b>97.5</b>	86.9	<b>55.2</b>	<b>92.3</b>
IA <sup>3</sup>	95.3	85.3	3.6	84.1
Prefix Tuning	96.3	<b>88.8</b>	0.8	91.4

Table 2: Macro F1 results to demonstrate the benchmarking use case of PEFT-FACTORY on different datasets for different PEFT methods.

PEFT-FACTORY thus allows the community to easily extend the PEFT-Bench with new PEFT methods or even design new benchmarks with minimal effort. By ensuring experiment equivalency, replicability, and ease of extensibility, PEFT-FACTORY empowers researchers and practitioners to rigorously evaluate existing PEFT approaches and accelerate the development of new ones.

## 5 Conclusion and Future Work

We introduce PEFT-FACTORY, a modular and extensible framework for fine-tuning modern autoregressive models using recent and diverse PEFT methods. PEFT-FACTORY not only provides a way to utilize PEFT methods but also implements support for various PEFT providers and a custom PEFT interface to promote replicability and transparency when designing new PEFT methods. When comparing PEFT-FACTORY to various popular fine-tuning frameworks, as well as to our upstream framework, LLaMA-Factory, its novelty lies in supporting different PEFT methods, classifying tasks with custom instructions, and providing PEFT- and dataset-level extensibility.

**Sustainability and Maintenance.** To keep up with the updates included in LLaMA-Factory (which often include support of new LLMs or improvements in the training pipeline), we will regularly release a new version of PEFT-FACTORY (this includes merging the upstream changes into our repository). Additionally, to include new features in PEFT-FACTORY itself, we will regularly release a separate version increment. Each change will be documented in the changelog of the specific release.

As the next steps, we would like to increase support for additional PEFT off-the-shelf methods, as well as reproduce some popular PEFT methods that are not currently supported by any of the PEFT provider frameworks. We believe that PEFT-FACTORY is an important and enabling tool that will promote the research of PEFT methods and allow their fair and consistent evaluation.

## Acknowledgments

This work was partially funded by the European Union, under the project LorAI - Low Resource Artificial Intelligence, GA No. [101136646](#); and by the European Union NextGenerationEU through the Recovery and Resilience Plan for Slovakia under the project No. 09I01-03-V04-00006.

Part of the research results was obtained using the computational resources procured in the national project, National Competence Centre for High Performance Computing (project code: 311070AKF2), funded by ERDF, EU Structural Funds Informatization of Society, Operational Program Integrated Infrastructure.

## References

- Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. 2019. Gradio: Hassle-free sharing and testing of ml models in the wild. *arXiv preprint arXiv:1906.02569*.
- Lightning AI. 2023. Litgpt. <https://github.com/Lightning-AI/litgpt>.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). 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 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. 2023. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. <https://github.com/nomic-ai/gpt4all>.
- Akari Asai, Mohammadreza Salehi, Matthew Peters, and Hannaneh Hajishirzi. 2022. [ATTEMPT: Parameter-efficient multi-task tuning via attentional mixtures of soft prompts](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6655–6672, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Axolotl maintainers and contributors. 2023. [Axolotl: Open source llm post-training](#).
- Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second PASCAL recognising textual entailment challenge.

- Robert Belanec, Branislav Pecher, Ivan Srba, and Maria Bielikova. 2025. *Peft-bench: A parameter-efficient fine-tuning methods benchmark*. *arXiv preprint*.
- Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. 2022. *BitFit: Simple parameter-efficient fine-tuning for transformer-based masked language-models*. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1–9, Dublin, Ireland. Association for Computational Linguistics.
- Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. 2009. The fifth PASCAL recognizing textual entailment challenge.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Arthur Cayley. 1846. Sur quelques propriétés des déterminants gauches. *Journal für die reine und angewandte Mathematik*.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. *SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation*. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. ACL.
- Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation. <https://github.com/sahil280114/codealpaca>.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. *BoolQ: Exploring the surprising difficulty of natural yes/no questions*. 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 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. *The pascal recognising textual entailment challenge*. In *Proceedings of the First International Conference on Machine Learning Challenges: Evaluating Predictive Uncertainty Visual Object Classification, and Recognizing Textual Entailment, MLCW’05*, page 177–190, Berlin, Heidelberg. Springer-Verlag.
- Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*, volume 23, pages 107–124.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: efficient finetuning of quantized llms. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.
- Shizhe Diao, Rui Pan, Hanze Dong, Ka Shun Shum, Jipeng Zhang, Wei Xiong, and Tong Zhang. 2023. Lmflow: An extensible toolkit for finetuning and inference of large foundation models. *arXiv preprint arXiv:2306.12420*.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, and 1 others. 2023. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235.
- William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the International Workshop on Paraphrasing*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. 2023. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third PASCAL recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, pages 1–9. Association for Computational Linguistics.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. 2024. *Parameter-efficient fine-tuning for large models: A comprehensive survey*. *Transactions on Machine Learning Research*.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. Lora+: efficient low rank adaptation of large models. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org.

- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. [Towards a unified view of parameter-efficient transfer learning](#). In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Daniel Khoshabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. 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 252–262.
- Tushar Khot, Ashish Sabharwal, and Peter Clark. 2019. [What’s missing: A knowledge gap guided approach for multi-hop question answering](#). 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 2814–2828, Hong Kong, China. Association for Computational Linguistics.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hector J Levesque, Ernest Davis, and Leora Morgenstern. 2011. The Winograd schema challenge. In *AAAI Spring Symposium: Logical Formalizations of Commonsense Reasoning*, volume 46, page 47.
- Shenggui Li, Hongxin Liu, Zhengda Bian, Jiarui Fang, Haichen Huang, Yuliang Liu, Boxiang Wang, and Yang You. 2023. [Colossal-ai: A unified deep learning system for large-scale parallel training](#). In *Proceedings of the 52nd International Conference on Parallel Processing*, ICPP ’23, page 766–775, New York, NY, USA. Association for Computing Machinery.
- 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.
- Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. 2023. Scaling down to scale up: A guide to parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.15647*.
- 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.
- Vijay Lingam, Atula Tejaswi Neerkaje, Aditya Vavre, Aneesh Shetty, Gautham Krishna Gudur, Joydeep Ghosh, Eunsol Choi, Alex Dimakis, Aleksandar Bojchevski, and sujay sanghavi. 2024. [SVFT: Parameter-efficient fine-tuning with singular vectors](#). In *2nd Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ICML 2024)*.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohata, Tenghao Huang, Mohit Bansal, and Colin Raffel. 2022a. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. Dora: weight-decomposed low-rank adaptation. In *Proceedings of the 41st International Conference on Machine Learning*, ICML’24. JMLR.org.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022b. [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.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023. Gpt understands, too. *AI Open*.
- 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>.
- Fanxu Meng, Zhaohui Wang, and Muhan Zhang. 2024. Pissa: principal singular values and singular vectors

- adaptation of large language models. In *Proceedings of the 38th International Conference on Neural Information Processing Systems, NIPS '24*, Red Hook, NY, USA. Curran Associates Inc.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. **Bleu: a method for automatic evaluation of machine translation**. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. Pytorch: an imperative style, high-performance deep learning library.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. **Are NLP models really able to solve simple math word problems?** In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020. **MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7654–7673, Online. Association for Computational Linguistics.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. **WiC: the word-in-context dataset for evaluating context-sensitive meaning representations**. 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 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.
- Clifton Poth, Hannah Sterz, Indraneil Paul, Sukannya Purkayastha, Leon Engländer, Timo Imhof, Ivan Vulić, Sebastian Ruder, Iryna Gurevych, and Jonas Pfeiffer. 2023. **Adapters: A unified library for parameter-efficient and modular transfer learning**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 149–160, Singapore. Association for Computational Linguistics.
- Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. 2023. Controlling text-to-image diffusion by orthogonal finetuning. *Advances in Neural Information Processing Systems*, 36:79320–79362.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of EMNLP*, pages 2383–2392. Association for Computational Linguistics.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *AAAI spring symposium: logical formalizations of commonsense reasoning*, pages 90–95.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavathula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. **Social IQa: Commonsense reasoning about social interactions**. 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 4463–4473, Hong Kong, China. Association for Computational Linguistics.
- Zhengxiang Shi and Aldo Lipani. 2024. **DePT: Decomposed prompt tuning for parameter-efficient fine-tuning**. In *The Twelfth International Conference on Learning Representations*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on EMNLP*, pages 1631–1642.
- Pengwei Tang, Xiaolin Hu, and Yong Liu. 2025. **ADePT: Adaptive decomposed prompt tuning for parameter-efficient fine-tuning**. In *The Thirteenth International Conference on Learning Representations*.
- A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.
- Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert,

- Shengyi Huang, Kashif Rasul, and Quentin Galouédec. 2020. Trl: Transformer reinforcement learning. <https://github.com/huggingface/trl>.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A sticker 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. *arXiv preprint arXiv:1804.07461*.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023a. How far can camels go? exploring the state of instruction tuning on open resources. *Preprint*, arXiv:2306.04751.
- Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogério Feris, Huan Sun, and Yoon Kim. 2023b. Multitask prompt tuning enables parameter-efficient transfer learning. In *The Eleventh International Conference on Learning Representations*.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments. *Transactions of the ACL*, 7:625–641.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the ACL: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. ACL.
- 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, and 3 others. 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.
- Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. 2023. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. *arXiv preprint arXiv:2312.12148*.
- Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, and 25 others. 2024. Qwen2.5 technical report. *ArXiv*, abs/2412.15115.
- Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. 2018. Learning to mine aligned code and natural language pairs from stack overflow. In *Proceedings of the 15th international conference on mining software repositories*, pages 476–486.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. 2023. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*.
- Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018. Record: Bridging the gap between human and machine commonsense reading comprehension. *arXiv preprint arXiv:1810.12885*.
- Bingchen Zhao, Haoqin Tu, Chen Wei, Jieru Mei, and Cihang Xie. 2024a. Tuning layernorm in attention: Towards efficient multi-modal LLM finetuning. In *The Twelfth International Conference on Learning Representations*.
- Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. 2024b. Galore: memory-efficient llm training by gradient low-rank projection. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- 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. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyang Luo. 2024a. LlamaFactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 400–410, Bangkok, Thailand. Association for Computational Linguistics.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024b. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

## A Ethical Considerations

The experiments in this paper were conducted using publicly available datasets, including SST-2, CoLA, WSC, and SVAMP, as cited by the original authors. As we were unable to determine the licenses for all used datasets, we have opted to use them in the limited form possible, adhering to the terms of use of the GLUE and SuperGLUE benchmarks. As the datasets are commonly used in other related works and have been published in scientific works that went through an established review process, we do not check for the presence of any offensive content, as it was already removed by the authors of these publicly available datasets. Additionally, we do not collect or utilize any personally identifiable information or offensive content, and we do not engage in crowdsourcing for data annotation in any form. To our knowledge, we are not aware of any potential ethical harms or negative societal impacts of our work, apart from those related to the field of Machine Learning (i.e., the use of computational resources that consume energy and produce heat, resulting in indirect CO<sub>2</sub> emissions). We follow the license terms for the LLaMa-3.2-1B-Instruct model we use; all models and datasets permit their use as part of the research. As we transform conditional generation into the classification problem (generating only labels), in most cases, we minimize the problem of generating offensive or biased content.

Importantly, in line with the open-science spirit, PEFT-FACTORY is an open-source downstream fork of LLaMA-Factory, licensed under the Apache-2.0 license (we respect the license and add append headers to the files that we have added or modified).

**Impact Statement: CO<sub>2</sub> Emissions Related to Experiments.** The experiments in this paper require GPU computing resources as we train and evaluate 1 model for different methods (3) and datasets (4). Overall, the experiments, including evaluations (which did not require training but still utilized GPU resources for inference) and preliminary experiments (which are outside the scope of our work), were conducted using a private infrastructure with a carbon efficiency of 0.432 kgCO<sub>2</sub>eq/kWh. Approximately 50 hours of computation were performed on hardware of type A100 PCIe 40GB (TDP of 250W). Total emissions are estimated to be 9.24 kg CO<sub>2</sub>eq, of which 0% were directly offset. Whenever possible, we tried to re-

duce the computational costs.

## B Further Details

In this section, we include detailed information about PEFT-FACTORY that can be used by advanced users to further understand or extend our framework. In Table 3, we provide a comparison of different easy-to-use fine-tuning frameworks in terms of available PEFT methods, highlighting the undeniable contribution of PEFT-FACTORY.

### B.1 Preprocessing datasets

Out of 27 included datasets, we namely preprocess and adapt MultiRC, WiC, COPA, ReCoRD, WSC, MMLU, PIQA, SIQA, HellaSwag, Wingrande, OBQA, MathQA, and SVAMP datasets. We upload all our adapted and preprocessed dataset versions to HuggingFace Hub<sup>6</sup>. Additionally, some datasets contain numerical values by default, formatted as `class` values in the HuggingFace dataset class. We convert such formats to textual representations to ensure compatibility with autoregressive generative models. Therefore, we transform multiple input and output columns of a single dataset to just a two-column format, including only *input* and *output* for the LLM.

### B.2 Efficiency Metrics

The PSCP metric comprises a set of constants that must be configured to function properly. Specifically *pspc\_cp*, *pscp\_cf*, *pscp\_cm*, *pspc\_bp*, *pscp\_bf*, and *pscp\_bm*. The  $C$  values are set by the first three attributes, and the  $\beta$  values are set by the last three attributes. We also provide default values for these attributes.

The  $C$  values in PSCP calculation represent reference constants used for scaling the parameters (*pscp\_cp*), inference time (*pscp\_cf*), and peak memory usage (*pscp\_cp*). The  $\beta$  values are defaultly set to 1, but can be set to any positive number. The higher the number, the higher the importance of the number of parameters *pspc\_bp*, inference time *pspc\_bf*, and peak memory usage *pspc\_bm*. For detailed information on how to set these constants and the full equation, see Belanec et al. (2025).

### B.3 Graphical User Interface

PEFT-FACTORY utilized LLaMA-Board graphical user interface based on Gradio (Abid et al., 2019). In this section, we describe the changes to

<sup>6</sup><https://hf.co/collections/kinit/peft-factory>

	Reparametrized							Soft Prompt-Based					Adapter-Based			Selective		PEFT Extensibility		
	LoRA	QLoRA	DoRA	LoRA+	PiSSA	Galore	OFT	SVFT	Prompt Tuning	Prefix Tuning	P-Tuning	P-Tuning V2	MTP	LA <sup>3</sup>	Bottleneck Adapter	Sequential Bottleneck Adapter	Parallel Adapter		BitFit	LNtuning
PEFT-FACTORY	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
LLaMA-Factory	✓	✓	✓	✓	✓	✓	✓	✓												
FastChat	✓	✓	✓	✓	✓	✓	✓	✓												
LiGPT	✓	✓	✓	✓	✓	✓	✓	✓												
LMFlow	✓	✓	✓	✓	✓	✓	✓	✓												
Axolotl	✓	✓	✓	✓	✓	✓	✓	✓												
Open-Instruct	✓	✓	✓	✓	✓	✓	✓	✓												
H2O LLM Studio	✓	✓	✓	✓	✓	✓	✓	✓												

Table 3: Comparison of different PEFT methods available in PEFT-FACTORY with popular LLM fine-tuning frameworks. Current frameworks do not include support for other than reparametrized PEFT methods, while most of them are LoRA variations. These are all PEFT methods that were tested for functionality. PEFT Extensibility means that the framework also supports the modular addition of newly created PEFT methods, either by PEFT provider frameworks or directly by users.

the graphical user interface that enable fine-tuning LLMs with various PEFT methods.

During construction of the Gradio interface, PEFT-FACTORY takes the available PEFT methods and their configurations and constructs an interface for each configuration. Figure 2 shows the available PEFT methods to choose from the list. Each PEFT method contains default values that will be set automatically. However, the configuration can be further specified by the detailed configuration shown in Figure 3, which displays the configuration options for the Prompt Tuning method (Lester et al., 2021).

## C Custom PEFT Method Templates

We provide minimal templates for the `model.py` and `config.py` files to design a PEFT method compatible with the PEFT-FACTORY custom PEFT interface, as documented in our framework<sup>7</sup>.

Additionally, we provide an example directory structure (shown in Figure 5) that can be used to ensure compatibility with dynamic loading of PEFT-FACTORY.

### Custom Method Directory Structure

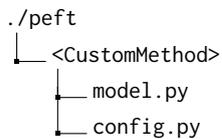
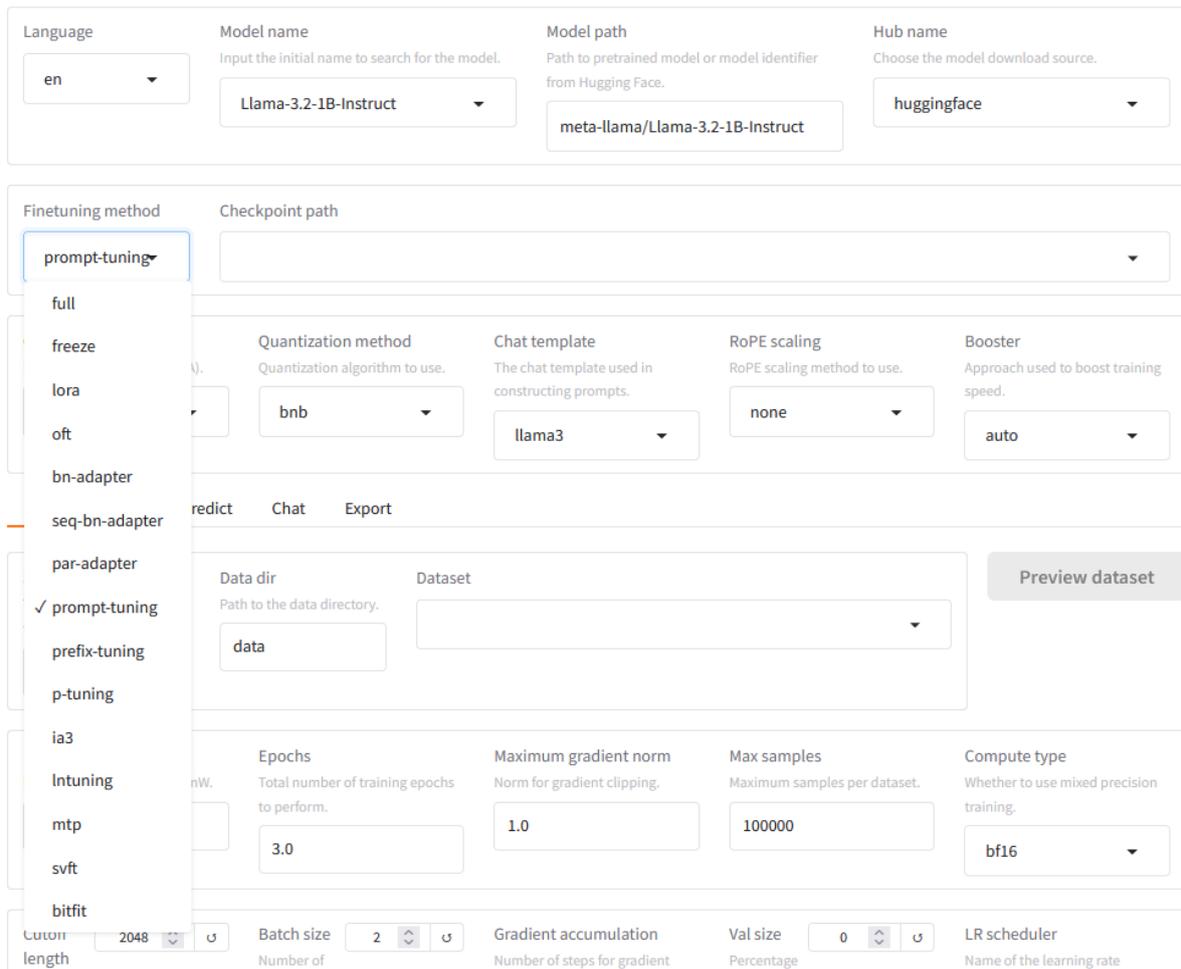


Figure 5: Example directory structure of custom PEFT interface used for dynamic loading of PEFT custom methods.

<sup>7</sup>For detailed information, please visit the [Templates section](#) of PEFT-FACTORY documentation.

## PEFT-Factory: Unified Parameter-Efficient Fine-Tuning of 100+ LLMs

Visit [GitHub Page](#)



The screenshot displays the PEFT-Factory web interface. At the top, there are four dropdown menus: Language (en), Model name (Llama-3.2-1B-Instruct), Model path (meta-llama/Llama-3.2-1B-Instruct), and Hub name (huggingface). Below these is a section for Finetuning method and Checkpoint path. The Finetuning method dropdown is open, showing a list of 19 options: full, freeze, lora, oft, bn-adapter, seq-bn-adapter, par-adapter, prompt-tuning (selected), prefix-tuning, p-tuning, ia3, lora, mtp, svft, and bitfit. The interface also includes fields for Quantization method (bnb), Chat template (llama3), RoPE scaling (none), and Booster (auto). There are buttons for Predict, Chat, and Export. A Data dir field contains 'data' and a Dataset dropdown is open. A Preview dataset button is present. At the bottom, there are input fields for Epochs (3.0), Maximum gradient norm (1.0), Max samples (100000), Compute type (bf16), Custom length (2048), Batch size (2), Gradient accumulation, Val size (0), and LR scheduler.

Figure 2: Selection of PEFT methods from Finetuning method dropdown menu. All 19 PEFT methods included in PEFT-FACTORY are available to choose.

prompt-tuning configurations

num\_virtual\_tokens

100

prompt\_tuning\_init

SAMPLE\_VOCAB

prompt\_tuning\_init\_text

tokenizer\_name\_or\_path

tokenizer\_kwargs

Figure 3: Configuration options for the Prompt Tuning method.

Maximum new tokens: 512

Top-p: 0.7

Temperature: 0.95

Output dir: eval\_2025-11-30-16-29-13

Preview command Start Abort

```

{
  "predict_accuracy": 0.3076923076923077,
  "predict_f1": 0.2506389193136181,
  "predict_flops": 1.65,
  "predict_memory": 5.6,
  "predict_model_preparation_time": 0.0011,
  "predict_params": 204800,
  "predict_pscp": 0.2,
  "predict_runtime": 31.768,
  "predict_samples_per_second": 3.274,
  "predict_steps_per_second": 1.637
}

```

Figure 4: Classification and PSCP results for prediction after training with Prompt Tuning.

# Similar, but why? A Toolkit for Explaining Text Similarity

Juri Opitz<sup>1\*</sup> Andrianos Michail<sup>1\*</sup> Lucas Möller<sup>2</sup> Sebastian Padó<sup>2</sup> Simon Clematide<sup>1</sup>

<sup>1</sup>University of Zurich, Switzerland

<sup>2</sup>IMS at University of Stuttgart, Germany

<sup>1</sup>{jurialexander.opitz, andrianos.michail, simon.clematide}@uzh.ch

<sup>2</sup>{lucas.moeller, pado}@ims.uni-stuttgart.de

## Abstract

Explaining text similarity and developing interpretable models are emerging research challenges (Opitz et al., 2025). We release XPLAINSIM, a Python package that unifies three complementary approaches for explaining textual similarity in an easily accessible way: 1. a token attribution method that explains how individual word interactions contribute to the predicted similarity of any embedding model; 2. a method for inferring structured neural embedding spaces that capture explainable aspects of text, and 3. a symbolic approach that explains textual similarity transparently through parsed meaning representations. We demonstrate the value of our package through intuitive examples and three focused empirical research studies. The first study evaluates interpretability methods for constructing cross-lingual token alignments. The second investigates how modern information retrieval methods handle stop words. The third sheds more light on a long-standing question in computational linguistics: the distinction between relatedness and similarity. XPLAINSIM is available at <https://github.com/flipz357/XPLAINSIM>.

## 1 Introduction and Background

Understanding semantic similarity is an important research question, both from an academic and practical perspective (Opitz et al., 2025). Partly, this is likely because it is a challenge in itself to express what makes two texts (dis)similar, even for humans (Fodor et al., 2024). Importantly, among the wide range of explainability methods (Sundararajan et al., 2020; Janizek et al., 2021), explanation of similarity represents a special case, since the assessment critically depends on *interactions between two inputs*, which significantly increases the complexity of interpretation and explanation (Figure 1). For example, when texts are represented as embeddings, the multiplicative interaction of

\*Equal contribution.

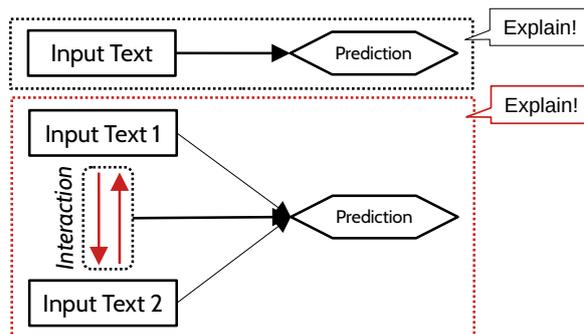


Figure 1: Complexity of explanation. Top: “Classic” single-input-prediction explanation problem. Bottom: A prediction for two inputs is not only influenced by either input, but also, first and foremost, by *their interaction*.

the two inputs through cosine distance yields the similarity score.

And such questions are not just theoretical. For instance, systems that search texts rely on text similarity—in sensitive domains, like law or medicine, understanding why a system delivers certain texts, but not others, may become especially crucial. Against the background of emerging AI laws (e.g., “right to explanation”; EU, 2024) the demand for transparency is expected to intensify.

In this work, we release an easy-to-use software package that combines three broadly applicable approaches for generating human-interpretable explanations of semantic textual similarity. The target audience covers both developers and end users interested in a better understanding of models of semantic similarity and the interpretation of their output. Finally, our toolbox can be a starting point for research in this area.

In its current version, our package includes three distinct explanation approaches. We select these three approaches since each of them is addressing the problem of similarity interpretability from a complementary perspective aligned with three explanation paradigms (Opitz et al., 2025): **Interaction attribution** (Moeller et al., 2023, 2024) is

Module	Tooling
attribution	Retrieve token-pair interactions for off-the-shelf models
spaceshaping	Learn structured neural representation of semantic aspects
symbolic	Meaning Representation similarity/ Parsing and matching graphs

Table 1: Simplified overview of included modules.

a *post-hoc approach* that allows to investigate off-the-shelf embedding models by explaining their decisions from an input perspective, focusing on token interactions. **Space shaping** (Opitz and Frank, 2022) lets us create *interpretable embeddings*. It enables the integration of custom semantic aspects into an embedding model, supporting fine-grained analysis of semantic decisions as well as efficient clustering and search. Finally, our **symbolic approach** (Banarescu et al., 2013; Opitz, 2023) uses Abstract Meaning Representation (AMR) graphs and grounds the similarity in the comparison of these *structured objects*.

An overview of all three available approaches is provided in Table 1. The modular design of our package supports future extensions and the integration of additional interpretability methods.

**Dependencies & License.** XPLAINSIM is released under the GPLv3 license and is publicly available at <https://github.com/flipz357/XPLAINSIM>. It can be installed via pip for Python and uses the import namespace `xplain` for brevity. Its main dependencies are `pytorch`<sup>1</sup> and `sentence-transformers`<sup>2</sup>. For symbolic similarity computation, we additionally rely on `amrlib`<sup>3</sup> and `smatchpp`<sup>4</sup>.

**Research study contribution.** In addition, our paper highlights our package’s value by contributing three focused research studies. The first experiment investigates the intrinsic cross-lingual alignment capabilities of multilingual embedding models. The second examines how IR-focused models handle stop words when matching queries to candidate documents. Both of these underlying phenomena are not directly observable in text embedding models and only become accessible through explainability methods. The third study uses mea-

<sup>1</sup><https://github.com/pytorch/pytorch>

<sup>2</sup><https://github.com/UKPLab/sentence-transformers>

<sup>3</sup><https://github.com/bjascob/amrlib>

<sup>4</sup><https://github.com/flipz357/smatchpp>

```
from xplain.attribution import ModelFactory
model= ModelFactory.build("gte-multilingual-base")
a= 'The dog runs after the kitten in the yard.'
b= 'Im Garten rennt der Hund der Katze hinterher.'
raw_A, tokens_a, tokens_b = \
    model.explain_similarity(a, b)
A_pp, tokens_a_pp, tokens_b_pp = \
    model.postprocess_attributions(...)
```

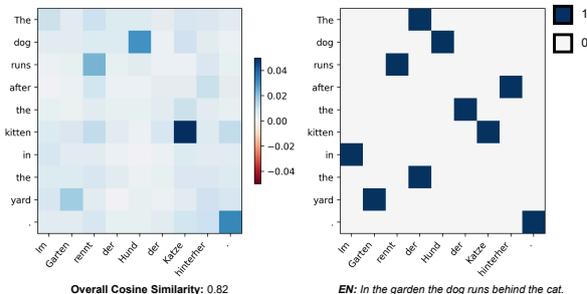


Figure 2: **Top:** Attribution generation for the off-the-shelf text embedding model `gte-multilingual-base`. **Left:** Attribution matrix, whose sum approximates the overall similarity score. **Right:** Sparsified attribution matrix generated by FlowAlign post-processing.

surements grounded in symbolic representations to shed light on structural differences between similarity and relatedness, two notions often conflated. These pilot studies demonstrate the practical benefits and distinct capabilities of our toolkit.

**Related work and context.** Several interpretability libraries exist for neural network-based NLP models (e.g., Kokhlikyan et al., 2020; Alammari, 2021; Attanasio et al., 2023; Fiotto-Kaufman et al., 2025). However, all these tools are focused on a single-input setting: they explain how input features contribute to a single model prediction. Text similarity is fundamentally different, because any prediction arises not only from each input in isolation but from the *interaction/comparison* between two inputs (cf. Figure 1). It is this gap that our XPLAINSIM package targets.

## 2 Implemented Methods

In this section, we introduce each of the three explainability methods included in XPLAINSIM and provide an example outlining its idea and usage.

### 2.1 Token Interaction Attributions

**Explanation mechanism.** For a given prediction, attribution methods assign importance values to input features, which, in our case, are pairs of words. That is, every pair of words (across two inputs) is

assigned a weight and the sum of all weights approximates the overall similarity score assigned by the embedding model. Moeller et al. (2023; 2024) have developed methods that compute such weights for embedding models. The interaction attributions generated by their method can be visualized as a matrix that decomposes the overall similarity score (cf. the left matrix in Figure 2).

**Implementation.** The method builds on the concept of integrated gradients, extending it to Siamese encoder architectures such as sentence transformers (Reimers and Gurevych, 2019). The resulting attributions are based on the *integrated Jacobians* of the embeddings of the two inputs with respect to their token representations. The Jacobians are computed with automatic differentiation; technical details can be found in the original publications.

We integrate the existing code base into our package and extend it in two ways: (1) We add support for multilingual and retrieval-focused embedding models. (2) We integrate optional post-processing for discretizing attribution matrices. This is motivated by the observation that raw attribution matrices tend to be relatively flat, making interpretation difficult. Discretization highlights the most relevant token-to-token alignments. These extensions provide the foundation for two demonstration experiments presented in Section 3.

**Post-attribution alignment sparsification.** The resulting attribution matrices indicate correspondences between the two inputs, similar to alignment matrices commonly used in machine translation. To enable comparison with alignment annotations, these matrices must be discretized and converted into a sparse binary format (Dou and Neubig, 2021). We implement two strategies:

**MaxAlign:** Each token is aligned to its strongest counterpart in the other sentence, but only if the link is mutual. Let  $AtoB(a)$  denote the token  $b$  in the second sentence that receives the highest attribution from token  $a$ , and  $BtoA(b)$  analogously for token  $a$ . Then:

$$\text{Align}_{ab} := \mathcal{I}[AtoB(a) = b \wedge BtoA(b) = a].$$

where  $\mathcal{I}[s]$  returns 1 if  $s$  is true, and 0 otherwise. This strategy yields a sparse, precision-oriented alignment by linking only mutually preferred token pairs and ignoring ambiguous or weak links.

**FlowAlign:** We also provide a more advanced alignment based on optimization. The attribution matrix is interpreted as a cost matrix by transforming attribution scores into costs (higher attribution means lower transport cost). Each token is assigned a weight of 1, and the Wasserstein distance is computed (aka Earth/Word Mover’s Distance, minimal transport Kusner et al., 2015). Intuitively, this distance expresses the minimal amount of work required to transform one set of embeddings to the other. A by-product of this computation is a sparse *Flow* matrix between embeddings (transportation plan), which we use as alignment  $\text{Align}_{ab} := \mathcal{I}[\text{Flow}_{ab} > \tau]$ , where  $\tau$  is a threshold.<sup>5</sup> The resulting alignment is slightly less sparse (n:m) and yields higher recall compared to MaxAlign.

**Example.** Figure 2 illustrates the process. We compute the similarity between an English and a German sentence using a multilingual GTE model (Li et al., 2023). Our code then generates the corresponding attribution matrix. Finally, we apply FlowAlign to discretize the attributions, revealing that the model aligns cross-lingual word pairs such as *dog* – *Hund*, or *kitten* – *Katze*. A systematic evaluation of this behavior is presented in Section 3.

## 2.2 Space Shaping

**Explanation mechanism.** Text embeddings are points in a high-dimensional vector space. This explanation method aims to decompose that space into lower-dimensional subspaces, each capturing a distinct semantic aspect (Opitz and Frank, 2022). For example, one subspace might represent named entities mentioned in a text, while another captures topical content. Similar to the attribution method, this approach decomposes the overall similarity score into multiple contributions—but at the level of abstract semantic aspects rather than individual token interactions. Such a decomposition enables explanations like: “two texts are similar in aspect X but differ in aspect Y”. In contrast to local attribution methods, which compute explanations during or after the similarity calculation, this approach embeds the explanation directly into the model. As a result, all explanatory structure is learned during training, and explanations are available at inference time without additional computational cost.

<sup>5</sup>An intuitive choice is  $\tau = 0$ , though we found that  $\tau = 0.029$  performs best across all language pairs in the dev set of the alignment task, and we set this as the default.

```

from sentence_transformers import InputExample
from xplain.spaceshaping import \
    PartitionedSentenceTransformer

# we need two lists with documents pairs
docs1, docs2 = ["abc",...], ["xyz",...]

# compute the training/partitioning signal
examples = []
for x, y in zip(docs1, docs2):
    similarities = []
    for metric in my_custom_metrics:
        similarities.append(metric.score(x, y))
    examples.append(InputExample(texts=[x, y], \
        label=similarities))

# instantiate model, we use 16 dimensions
# to express each metric
pt = PartitionedSentenceTransformer(
    feature_names=[metric.name for \
        metric in my_custom_metrics],
    feature_dims=[16]*len(my_custom_metrics))
train_examples, dev_examples = split(examples)
# train partitioning
pt.train_model(train_examples, dev_examples)

```

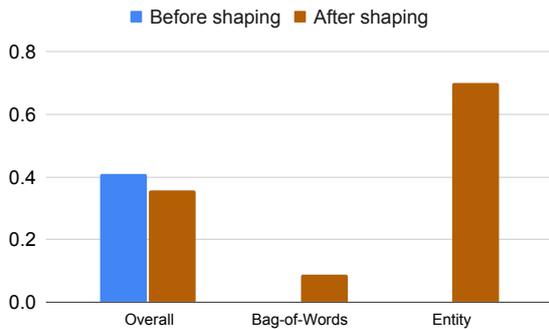


Figure 3: **Top:** Pseudo code to induce semantic subspaces. **Bottom:** Aspect similarities for the sentences “The kitten drinks milk” vs. “A cat slurps something.”

**Implementation.** Two challenges arise: Learning to partition the space, and preventing that the overall similarity of two texts deviates too much from a strong reference model<sup>6</sup>.

We extend the “S3BERT” method<sup>7</sup> from Opitz and Frank (2022) which used AMR-based metrics to measure aspectual similarities and generate a partitioning signal. We generalize this approach to enable the shaping of fully custom embedding spaces. Users can define their own interpretable similarity metrics to create distinct embedding subspaces. In the following paragraph, we provide an example of the training and inference process.

<sup>6</sup>Such a reference model can be any embedding model, from which the explainable partitioned model learns to maintain high prediction accuracy.

<sup>7</sup><https://github.com/flipz357/S3BERT>

```

from xplain.symbolic.model import AMRSimilarity
explainer = AMRSimilarity()
sent1 = ["Barack Obama holds a talk"]
sent2 = ["Hillary Clinton holds a talk"]
exp = explainer.explain_similarity(sent1, sent2)
# yields aspectual similarities:
# 'AGENT': 71.43, # (Ob. vs Cl. talking)
# 'NER': 60.0, # (entity: !=, type: ==)
# ..., # (various other stats)
# 'POLARITY': 100.0 # (both: no negation)
# 'global': 80.0 # (overall similarity)

# ::snt Barack Obama holds a talk.
(h / hold-04
 :ARG0 (p / person
 :name (n / name
 :op1 "Barack" :op2 "Obama" )))
:ARG1 (t / talk-01))

# ::snt Hillary Clinton holds a talk.
(h / hold-04
 :ARG0 (p / person
 :name (n / name
 :op1 "Hillary" :op2 "Clinton" )))
:ARG1 (t / talk-01))

```

Figure 4: **Top:** Using XPLAINSIM’s symbolic part to explain similarity via meaning graph difference statistics. **Bottom:** AMR graph differences make meaning disagreement overt.

**Example.** Figure 3 gives a high-level view of the space shaping process (the base model here is all-MiniLM-L12-v2, Reimers and Gurevych, 2019). First, the user defines interpretable aspects of interest using custom metrics. These metrics can be simple and approximate—for example, measuring word overlap between two texts to capture superficial structural similarity, or measure the similarity in named entity structure of the two input texts via SpaCy<sup>8</sup> NER tagging. A set of paired training texts is then created by assigning scores based on these custom metrics.

### 2.3 Symbolic Explanation with AMR

**Explanation mechanism.** A meaning representation (MR) is a symbolic encoding of the semantic structure of a text, typically grounded in linguistic theories of compositionality and discourse. An MR can take the form of a graph, where nodes represent entities and events mentioned in a text, and edges show their semantic relations (*agent, patient, instrument, cause*, etc.). With meaning expressed in such an explicit format, a metric between MRs can highlight (dis-)agreements with respect to specific parts and properties of semantic structure. Several

<sup>8</sup><https://spacy.io/>

papers have already explored MR measurements for, e.g., semantic similarity, language inference, generation evaluation, or cross-lingual analysis.<sup>9</sup> We provide the first ready-to-use tool that integrates parsing and aspect-based AMR similarity measurement in a single package. Ready-to-use specifically means that a strong default parser is integrated in our package, which is, to our knowledge, not the case for any other MR metric package.

**Implementation.** We need to select a type of meaning representation, a parser that generates such representations, and a metric that can compare such representations as well as (ideally) any of their subgraphs. For our package, we also offer a default way of measuring, based on two pillars:

1. Representation and Parsing: The AMR representation (Abstract MR, [Banarescu et al., 2013](#)) has broad applications and large resources ([Wein and Opitz, 2024](#); [Sadeddine et al., 2024](#)) as well as fairly accurate parsers ([Bai et al., 2022](#)). Concretely, we leverage the `amrLib` library that has pre-trained parsers for AMR representation. For an overview of the currently available parsing models we refer the reader to the Appendix, Table 5.

2. Measuring: Finding the largest common subgraph of two AMR graphs is an NP-complete problem ([Allen et al., 2008](#); [Nagarajan and Sviridenko, 2009](#); [Cai and Knight, 2013](#)). We adopt the `smatchpp` library ([Opitz, 2023](#)), a graph matching library that uses Integer Linear Programming. The similarity score of two graphs is the amount of shared nodes plus the amount of shared edges, normalized by the size of either graph, obtaining two directional similarities. For a symmetric similarity score, the directional similarities are averaged with harmonic mean. We compute this measure for several aspects elicited by AMR subgraphs, e.g., coreference, negations, named entities. For an overview of the currently available measurements we refer the reader to the Appendix, Table 4.

**Example.** Figure 4 shows an application of the AMR-based approach to two sentences. After parsing the inputs and measuring aspectual similarities, we find that both sentences have a similar event

<sup>9</sup>For a sample, we refer the reader to this list: [Manning and Schneider \(2021\)](#); [Opitz and Frank \(2021\)](#); [Wein and Schneider \(2024\)](#); [Müller and Kuwertz \(2022\)](#); [Opitz et al. \(2023\)](#); [Ghosh et al. \(2024\)](#); [Jayaweera et al. \(2024\)](#); [Kachwala et al. \(2024\)](#); [Sun and Xue \(2024\)](#); [Landes and Di Eugenio \(2024\)](#); [Park et al. \(2024\)](#); [de Vergnette et al. \(2025\)](#); [Thatikonda et al. \(2025\)](#)

Model	Discretiz.	X-Ling Word Alignment		
		Pr	Re	F1
<b>Baseline</b>	Diagonal	0.262	0.245	0.253
<b>XLM-R</b>	MaxAlign	0.527	0.064	0.114
	FlowAlign	0.520	0.470	0.494
<b>M-MPNet</b>	MaxAlign	0.797	0.334	0.471
	FlowAlign	0.636	0.578	0.606
<b>M-MiniLM</b>	MaxAlign	<b>0.799</b>	0.371	0.507
	FlowAlign	0.662	0.605	0.632
<b>M-E5-Base</b>	MaxAlign	0.775	0.477	0.591
	FlowAlign	0.668	<b>0.610</b>	<b>0.638</b>
<b>M-GTE</b>	MaxAlign	0.779	0.444	0.566
	FlowAlign	0.667	<b>0.610</b>	0.637

Table 2: Cross-lingual word-level alignment results aggregated across all languages. Full model ID’s: Table 3

with a different agent (`arg0`). Matching the full graphs with `smatchpp` yields a similarity of 0.8. Matching only the subgraphs capturing patients in events (`arg1`), yields a perfect match (essentially: Someone holds a talk), but the subgraphs that capture Named Entities match with only 0.6 (indeed, their only agreement is the type (`person`)).

### 3 Pilot Studies

We show the potential utility of our XPLAINSIM package in three empirical pilot studies.

#### 3.1 Cross-lingual Alignment

Multilingual embedding models are typically trained contrastively, on positive and negative pairs across languages, but are not explicitly supervised at the word level. In this experiment, we investigate to what extent multilingual models internally develop unsupervised cross-lingual word alignments.

**Experiment.** We use the ten English–{bg, da, es, et, hu, it, nl, pt, ru, sl} gold alignment datasets released by [Martelli et al. \(2023\)](#), see Appendix C for details. For each sentence pair, we compute interaction attributions from several multilingual embedding models. To compare these attributions with the gold alignment annotations, we discretize them into binary matrices. For this purpose, we apply both of our sparsification methods: the row- and column-wise max-pooling heuristic (MaxAlign), and the optimal transport-based approach. We evaluate the resulting binary alignment matrices using standard precision, recall, and F1 score. As a baseline, we include a diagonal alignment matrix.

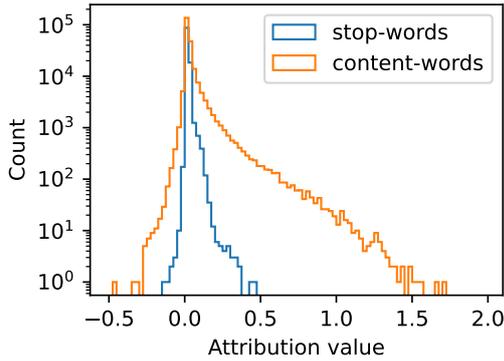


Figure 5: Histograms of attributions to stop and content words; MS MARCO validation split, IR model.

**Result.** Table 2 shows that M-GTE and M-E5-Base perform best in approximating a human-style discrete alignment. Notably, these models also perform well on multilingual retrieval data (e.g., MIR-ACL Zhang et al., 2023), suggesting that stronger alignment capabilities contribute to high performance in multilingual similarity and retrieval tasks.

FlowAlign post-processing consistently improves alignment performance across models. For users aiming to extract word-level alignments from attributions, we recommend combining a strong embedding model with FlowAlign for sparsification. In some cases, this combination yields substantial gains, e.g. XLM-R improves from an F1 score of 0.11 to 0.49, a 38-point increase. The impact of sparsification becomes smaller for higher-performing models, suggesting that stronger models rely on a latent discrete alignment.

### 3.2 Stop-words in Dense Retrieval

Semantic search is an application of text embedding models with high practical relevance. Dual encoder architectures can process queries and documents independently. The storage of document embeddings in vector databases enables efficient search in sublinear time complexity via approximate nearest neighbor algorithms. However, compressing entire documents into unified vector representations makes such dense retrieval results difficult to interpret. Using the interaction attribution method by Moeller et al. (2023), we study which parts of a given query and document a model matches.

**Experiment.** We build on the common distinction between stop and content words to evaluate whether an IR-optimized text embedding model effectively learns to suppress stop words. To this end,

we use the MS MARCO Passage Ranking dataset (Bajaj et al., 2016) (details in Appendix C) and evaluate a model fine-tuned on its training split<sup>10</sup>. For each positive query–passage pair in the validation split, we compute interaction attributions and sum the total attribution to all passage tokens. We then aggregate these contributions separately for stop words and content words, using NLTK’s stop word list for categorization.

**Result.** Figure 5 shows histograms of all stop- and content-word attributions, respectively. 97.6% of all attributions to stop-words fall within an interval of  $\pm 0.05$  around zero. This shows that the model learned to effectively ignore stop words. The distribution for content words, on the other hand, is much wider, ranging from  $-0.46$  to  $1.70$ . The model assigns both strongly positive and negative importance scores to words, suggesting that it can reward well-aligned document segments, penalize mismatches, and ignore large portions of document content, as indicated by the many content words whose attribution is near zero. Understanding which parts of queries and documents such models align—and where they are prone to errors or biases—is a promising direction for future research.

### 3.3 Relation Characterization with a Symbolic Approach

“Similarity” and “relatedness” are often treated as equivalent notions; however, this equivalence is not entirely accurate, as latent differences exist between them (Budanitsky and Hirst, 2006). To better understand the structural distinctions between similarity and relatedness, we employ symbolic, AMR-based aspectual measurements.

**Experiment.** For the similarity dataset, we use the validation partition of the STS benchmark (Cer et al., 2017); for relatedness, we use the SICK dataset (Marelli et al., 2014). Both datasets are well-established evaluation benchmarks, see Appendix C for details. We compute aspectual graph metrics for all sentence pairs in both datasets. Correlating these measures directly with the human similarity score is confounded by cases where particular semantic aspects are absent in both texts. Therefore, we compute Pearson correlation separately for each aspect, considering only sentence pairs where the graph metric detected an aspectual

<sup>10</sup>[sentence-transformers/msmarco-MiniLM-L12-cos-v5](https://huggingface.co/sentence-transformers/msmarco-MiniLM-L12-cos-v5)



Figure 6: Pearson correlation of aspectual differences, captured via symbolic meaning (sub-)graphs and graph metrics, with human-assessed similarity (STS) and relatedness (SICK).

difference.

**Result.** Figure 6 presents the obtained Pearson correlations. A notable difference is that variations in polarity and named entities are predictive of differences in human-assessed similarity (STS), but less so for relatedness (SICK). A commonality across both similarity and relatedness is their sensitivity to differences in shared concepts. Interestingly, similarity differences are more strongly associated with changes in *patients* (entities undergoing an action), whereas relatedness appears more influenced by the similarity of *agentive* structures.

**Discussion.** This analysis provides a rich contrastive characterization of similarity and relatedness. It also extends prior work, which has largely focused on individual words, by examining full semantic structures instead.

## 4 Summary

We release the XPLAINSIM package to support a better understanding of semantic similarity, and how it is computed in neural and non-neural models. In particular, we provide three types of methods, each with distinct trade-offs. Attribution-based methods reveal how interactions between words contribute to the final similarity score, from the perspective of neural embedding models. The most efficient approach is space shaping, which integrates explanations directly into the embedding space, effectively reducing the cost of generating explanations to zero. While it requires model training, it also enables customization of the explanation. Finally, the symbolic approach offers fine-grained comparisons based on meaning represen-

tations, adding an additional layer of transparency. Its limitations include reduced sensitivity to lexical nuance and dependence on parsing components.

Within this package description paper, we also contribute three research studies that highlight the value of the tools: 1. evaluating cross-lingual token alignment, 2. assessing structural token-weighting in information retrieval models, and 3. examining the distinction between similarity and relatedness.

We hope that XPLAINSIM lowers the barrier for both researchers and practitioners to explore explanations for why two texts are similar, beyond a single similarity score.

## Limitations

Our package is designed to be broadly applicable for all kinds of text similarity and explanation tasks. Potential constraints may arise from the specific characteristics of the included models: Space shaping requires custom design of measures and training, attribution analysis is applicable to the broader class of Siamese transformer models but comes at high computational costs that make it infeasible to apply to long-context tasks without further adaptations. Finally, meaning representations and their metrics provide highly interpretable and controllable similarity statistics, but they currently may not correlate as highly with human similarity ratings in STS annotation studies when compared to large neural models. This may not only be due to some structural insufficiencies, but also potentially due to noise from the parsing process, further aggravated by eventual domain shifts.

However, we emphasize that our package is designed to be easily usable, and readily extensible, such that implementation of other potential approaches, variation or improvement of current ones, and adding functionality, is straightforward. For future work, we also plan on expanding the library to include other interpretability methods adapted to semantic similarity, e.g., based on causal probing, counterfactual explanations, or Shapley values. We also plan on running human-centered studies for evaluating similarity with and through similarity interpretability methods.

For these and other future works with XPLAINSIM, we warmly invite the community to participate in contributing.

## Ethics statement

We do not identify direct ethical risks arising from the release of this toolkit. On the contrary, the primary purpose of our package is to explain models and model decisions. Since models from this domain are widely applied, also in document search contexts, we believe that better understanding their mechanism can only help to also better assess their risks, such as potential retrieval biases favoring certain documents.

## Acknowledgments

The Zurich authors received funding from the Swiss National Science Foundation (SNSF 213585) and the Luxembourg National Research Fund (17498891) through the *Impresso II* project.

## References

- J Alammari. 2021. [Ecco: An open source library for the explainability of transformer language models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 249–257, Online. Association for Computational Linguistics.
- James F. Allen, Mary Swift, and Will de Beaumont. 2008. [Deep semantic analysis of text](#). In *Semantics in Text Processing. STEP 2008 Conference Proceedings*, pages 343–354. College Publications.
- Giuseppe Attanasio, Eliana Pastor, Chiara Di Bonaventura, and Debora Nozza. 2023. [ferret: a framework for benchmarking explainers on transformers](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 256–266, Dubrovnik, Croatia. Association for Computational Linguistics.
- Xuefeng Bai, Yulong Chen, and Yue Zhang. 2022. [Graph pre-training for AMR parsing and generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6001–6015, Dublin, Ireland. Association for Computational Linguistics.
- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*.
- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. [Abstract Meaning Representation for sembanking](#). In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186, Sofia, Bulgaria. Association for Computational Linguistics.
- Alexander Budanitsky and Graeme Hirst. 2006. Evaluating WordNet-based measures of lexical semantic relatedness. *Computational Linguistics*, 32(1):13–47.
- Shu Cai and Kevin Knight. 2013. [Smatch: an evaluation metric for semantic feature structures](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 748–752, Sofia, Bulgaria. Association for Computational Linguistics.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. [SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.
- Rémi de Vergnette, Maxime Amblard, and Bruno Guillaume. 2025. [Evaluation framework for layered meaning representation](#). In *Proceedings of the Sixth International Workshop on Designing Meaning Representations*, pages 38–48, Prague, Czechia. Association for Computational Linguistics.
- Zi-Yi Dou and Graham Neubig. 2021. [Word alignment by fine-tuning embeddings on parallel corpora](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2112–2128, Online. Association for Computational Linguistics.
- EU. 2024. [Article 86: Right to explanation of individual decision-making](#).
- Jaden Fiotto-Kaufman, Alexander R. Loftus, Eric Todd, Jannik Brinkmann, Koyena Pal, Dmitrii Troitskii, Michael Ripa, Adam Belfki, Can Rager, Caden Juang, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Nikhil Prakash, Carla Brodley, Arjun Guha, Jonathan Bell, Byron C. Wallace, and David Bau. 2025. [Nnsight and ndif: Democratizing access to open-weight foundation model internals](#). In *The Thirteenth International Conference on Learning Representations*.
- James Fodor, Simon De Deyne, and Shinsuke Suzuki. 2024. [Compositionality and sentence meaning: Comparing semantic parsing and transformers on a challenging sentence similarity dataset](#). *Computational Linguistics*, pages 1–52.
- Sreyan Ghosh, Utkarsh Tyagi, Sonal Kumar, Chandira Kiran Evuru, Ramaneswaran S, S Sakshi, and Dinesh Manocha. 2024. [ABEX: Data augmentation for low-resource NLU via expanding abstract descriptions](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 726–748,

- Bangkok, Thailand. Association for Computational Linguistics.
- Joseph D Janizek, Pascal Sturmfels, and Su-In Lee. 2021. Explaining explanations: Axiomatic feature interactions for deep networks. *Journal of Machine Learning Research*, 22(104):1–54.
- Chathuri Jayaweera, Sangpil Youm, and Bonnie J Dorr. 2024. **AMREx: AMR for explainable fact verification**. In *Proceedings of the Seventh Fact Extraction and VERification Workshop (FEVER)*, pages 234–244, Miami, Florida, USA. Association for Computational Linguistics.
- Zoher Kachwala, Jisun An, Haewoon Kwak, and Filippo Menczer. 2024. **REMATCH: Robust and efficient matching of local knowledge graphs to improve structural and semantic similarity**. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 1018–1028, Mexico City, Mexico. Association for Computational Linguistics.
- Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. 2020. **Captum: A unified and generic model interpretability library for pytorch**. *Preprint*, arXiv:2009.07896.
- Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. 2015. From word embeddings to document distances. In *International conference on machine learning*, pages 957–966. PMLR.
- Paul Landes and Barbara Di Eugenio. 2024. **CALAMR: Component ALignment for Abstract Meaning Representation**. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2622–2637, Torino, Italia. ELRA and ICCL.
- Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*.
- Emma Manning and Nathan Schneider. 2021. **Referenceless parsing-based evaluation of AMR-to-English generation**. In *Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems*, pages 114–122, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014. **A SICK cure for the evaluation of compositional distributional semantic models**. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 216–223, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Federico Martelli, Andrei Stefan Bejgu, Cesare Campagnano, Jaka Cibej, Rute Costa, Apolonija Gantar, Jelena Kallas, Svetla Peneva Koeva, Kristina Koppel, Simon Krek, Margit Langemets, Veronika Lipp, Sanni Nimb, Sussi Olsen, Bolette Sandford Pedersen, Valeria Quochi, Ana Salgado, László Simon, Carole Tiberius, Rafael-J. Ureña-Ruiz, and Roberto Navigli. 2023. **XL-WA: a gold evaluation benchmark for word alignment in 14 language pairs**. In *CLiC-it*.
- Lucas Moeller, Dmitry Nikolaev, and Sebastian Padó. 2023. **An attribution method for Siamese encoders**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15818–15827, Singapore. Association for Computational Linguistics.
- Lucas Moeller, Dmitry Nikolaev, and Sebastian Padó. 2024. **Approximate attributions for off-the-shelf Siamese transformers**. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2059–2071, St. Julian’s, Malta. Association for Computational Linguistics.
- Almuth Müller and Achim Kuwertz. 2022. Evaluation of a semantic search approach based on amr for information retrieval in image exploitation. In *2022 Sensor Data Fusion: Trends, Solutions, Applications (SDF)*, pages 1–6. IEEE.
- Viswanath Nagarajan and Maxim Sviridenko. 2009. On the maximum quadratic assignment problem. *Mathematics of Operations Research*, 34(4):859–868.
- Juri Opitz. 2023. **SMATCH++: Standardized and extended evaluation of semantic graphs**. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1595–1607, Dubrovnik, Croatia. Association for Computational Linguistics.
- Juri Opitz and Anette Frank. 2021. **Towards a decomposable metric for explainable evaluation of text generation from AMR**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1504–1518, Online. Association for Computational Linguistics.
- Juri Opitz and Anette Frank. 2022. **SBERT studies meaning representations: Decomposing sentence embeddings into explainable semantic features**. 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 625–638, Online only. Association for Computational Linguistics.
- Juri Opitz, Lucas Moeller, Andrianos Michail, Sebastian Padó, and Simon Clematide. 2025. **Interpretable text embeddings and text similarity explanation: A survey**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 22314–22330, Suzhou, China. Association for Computational Linguistics.

- Juri Opitz, Shira Wein, Julius Steen, Anette Frank, and Nathan Schneider. 2023. [AMR4NLI: Interpretable and robust NLI measures from semantic graphs](#). In *Proceedings of the 15th International Conference on Computational Semantics*, pages 275–283, Nancy, France. Association for Computational Linguistics.
- Jinwoo Park, Hosoo Shin, Dahee Jeong, and Junyeong Kim. 2024. [Improving the representation of sentences with reinforcement learning and amr graph](#). In *2024 IEEE International Conference on Consumer Electronics (ICCE)*, pages 1–4.
- 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.
- Zacchary Sadeddine, Juri Opitz, and Fabian Suchanek. 2024. [A survey of meaning representations – from theory to practical utility](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2877–2892, Mexico City, Mexico. Association for Computational Linguistics.
- Haibo Sun and Nianwen Xue. 2024. [Anchor and broadcast: An efficient concept alignment approach for evaluation of semantic graphs](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 1052–1062, Torino, Italia. ELRA and ICCL.
- Mukund Sundararajan, Kedar Dhamdhere, and Ashish Agarwal. 2020. The shapley taylor interaction index. In *International conference on machine learning*, pages 9259–9268. PMLR.
- Ramya Keerthy Thatikonda, Wray Buntine, and Ehsan Shareghi. 2025. [Assessing the sensitivity and alignment of FOL closeness metrics](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 16775–16785, Suzhou, China. Association for Computational Linguistics.
- Shira Wein and Juri Opitz. 2024. [A survey of AMR applications](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6856–6875, Miami, Florida, USA. Association for Computational Linguistics.
- Shira Wein and Nathan Schneider. 2024. [Assessing the cross-linguistic utility of abstract meaning representation](#). *Computational Linguistics*, 50(2):419–473.
- Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamaloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2023. [MIRACL: A multilingual retrieval dataset covering 18 diverse languages](#). *Transactions*

Short Name	Text Embedding Model HF ID	Base Model HF ID
<b>XLM-R</b>	FacebookAI/xlm-roberta-base	FacebookAI/xlm-roberta-base
<b>M-MPNet</b>	sentence-transformers/paraphrase-multilingual-mpnet-base-v2	FacebookAI/xlm-roberta-base
<b>M-MiniLM</b>	sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2	microsoft/Multilingual-MiniLM-L12-H384
<b>M-E5-Base</b>	intfloat/multilingual-e5-base	FacebookAI/xlm-roberta-base
<b>M-GTE</b>	Alibaba-NLP/gte-multilingual-base	Alibaba-NLP/gte-multilingual-mlm-base

Table 3: Model Nomenclature: Mapping of short names to Hugging Face model IDs. The package supports most current embedding models by default. Users can load them by passing the corresponding Hugging Face model ID. Specialized architectures may require the addition of a custom subclass.

Aspect	Brief Description	Trigger Relation(s)	(Typed) Concept Nodes
AGENT	Actor / doer	:arg0	person, animal, nationality ...
CAUSE	Cause of event	:cause	cause-01
CONCEPT	Generic concept	:instance	—
FOCUS	Main predicate	:root	—
INSTRUMENT	Tool used	:instrument	instrument, tool
LOCATION	Place / path	:location, :path, :destination, :direction	city, state, river ...
MATHS	Mathematical entity	—	sum-of, product-of
NER	Named entity	:name	—
PATIENT	Entity affected	:arg1–:arg9	person, object
POLARITY	Negation / polarity	:polarity	—
POSSESSION	Ownership / possession	:poss	owner, possession
PURPOSE	Intended goal	:purpose	purpose-01
QUANTIFIER	Quantity / amount	:quant	monetary-quantity, distance-quantity, volume-quantity ...
QUESTION	Question structures	—	amr-unknown
SRL-core	Core semantic roles	:arg0–:arg9	person, object
TIME (temporal)	Temporal info	:time, :duration, :frequency	date-entity, date-interval
TOPIC	Subject / topic	:topic	topic-01
WIKI	Wikipedia link	:wiki	—

Table 4: Overview of graph aspects, trigger relations, and example concept nodes

## A Model Nomenclature

See Table 3.

## B Symbolic metrics and parsers

See Table 4 (metrics) and 5 (parsers).

## C Dataset Details

Further details on datasets from Section 3: Table 6.

Name	Size	Score	Speed
parse_xfm_bart_large	1.4GB	83.7 SMATCH	17/sec
parse_xfm_bart_base	492MB	82.3 SMATCH	31/sec
parse_spring	1.5GB	83.5 SMATCH	14/sec
parse_t5	785MB	81.9 SMATCH	11/sec
parse_gsii	787MB	76.8 SMATCH	28/sec

Table 5: Default `amrlib` parsing models for which our package automates installation. Table is taken from `amrlib`: <https://github.com/bjascob/amrlib-models>. “Speed is the inference speed on the AMR-3 test set (1898 graphs) using an RTX3090 with `num_beams=1` and `batch_size=32`. The units are sentences/second”.

Study	Dataset	Description	Size
Cross-lingual Alignment	XL-WA	Manual word alignment benchmark for 14 language pairs English–X	≈100 sent. in dev set, 200 sent. in test set (1500 word alignments in dev, 4000 word alignments in test) for each language pair. Taken from <a href="#">Martelli et al. (2023)</a>
Stopwords in Retrieval	MS-MARCO (v1.1)	Dataset for IR experiments including real-world questions, gold answers, and a set of candidate passages	9.65K queries in the test split, each with 10 candidate passages. Taken from <a href="#">Bajaj et al. (2016)</a> .
Relation Characterization	STS	Sentence pairs with similarity judgments on a Likert scale	1,371 sentence pairs coming from STS Benchmark. Taken from <a href="#">Cer et al. (2017)</a> .
Relation Characterization	SICK	Sentence pairs with relatedness judgments on a Likert scale	9,927 sentence pairs coming from SICK-R. Taken from <a href="#">Marelli et al. (2014)</a> .

Table 6: Dataset details for the use-case studies.

# ALIGNFIX: A Tool for Parallel Corpora Augmentation and Refinement

Samuel Frontull and Simon Haller-Seeber

Department of Computer Science, University of Innsbruck, Austria

{samuel.frontull, simon.haller-seeber}@uibk.ac.at

## Abstract

High-quality datasets are crucial for training effective state of the art machine translation systems. However, due to the data-intensive nature of these systems, they have to be trained on large amounts of text that can easily go beyond the scope of full human inspection. This makes the presence of noise that can degrade overall system performance a frequent and significant issue. While various approaches have been developed to identify and select only the highest-quality training examples, this is undesirable in scenarios where resources are limited. For this reason, we introduce AlignFix, an open-source tool for augmenting data, identifying and correcting errors in parallel corpora. Leveraging word alignments, AlignFix extracts consistent phrase pairs, enabling targeted replacements that can improve the dataset quality. Besides targeted replacements, the tool enables contextual augmentation by duplicating sentences and allowing users to substitute words with alternatives of their choice. The tool maintains and updates the underlying word alignments, thereby avoiding the costly recomputation. AlignFix runs locally in the browser, requires no installation, and ensures that all data remains entirely on the client side. It is released under Apache 2.0 license, encouraging broad adoption, reuse, and further development. A live demo is available at <https://ifi-alignfix.uibk.ac.at>.

## 1 Introduction

High-quality, carefully curated datasets are critical for the development of reliable machine translation (MT) systems. In an ideal scenario, only fully manually verified data would be available. However, neural MT systems are highly data-intensive (Koehn and Knowles, 2017; Gordon et al., 2021), necessitating the collection of as many texts as possible for training. Although modern architectures have enabled transfer learning for scenarios

with limited resources (Zoph et al., 2016), a sufficient amount of training data must still be accumulated (Gu et al., 2018).

For machine translation, the so-called *contextual augmentation* (Kobayashi, 2018; Wu et al., 2019; Gao et al., 2019) is an established technique for data augmentation and extends existing corpora by reusing existing sentences and replacing words in them. This technique is particularly effective for enhancing lexical coverage and ensuring the representation of rare words. However, this method requires a solid data foundation and relies on language models that provide sensible replacements, which are usually not available in low-resource scenarios.

A significantly more accessible and effective alternative is *back-translation* (Sennrich et al., 2016) which involves translating available monolingual target-language text into the source language using an auxiliary model. This approach can provide the (synthetic) parallel training data necessary to leverage state-of-the-art neural architectures for MT. In practice, however, datasets often become unwieldy when scaled to meet these requirements, resulting in a diminished insight into the data.

Synthesised data often contains numerous errors that can negatively impact the overall quality of a machine translation system (Hoang et al., 2018). Noisy data can for example lead to erroneous translations or hallucinations (Khayrallah and Koehn, 2018; Guerreiro et al., 2023). Consequently, it is crucial to filter and clean such datasets. Several tools have been developed for this purpose (Bogoychev et al., 2023; Zaragoza-Bernabeu et al., 2022; Aulamo et al., 2020). However, these tools primarily focus on data filtering, retaining only the highest-quality translation pairs and discarding the remainder. In contexts where data is scarce, aggressive filtering is not always desirable (Marashian et al., 2025). In such cases, it is preferable to identify and correct errors.

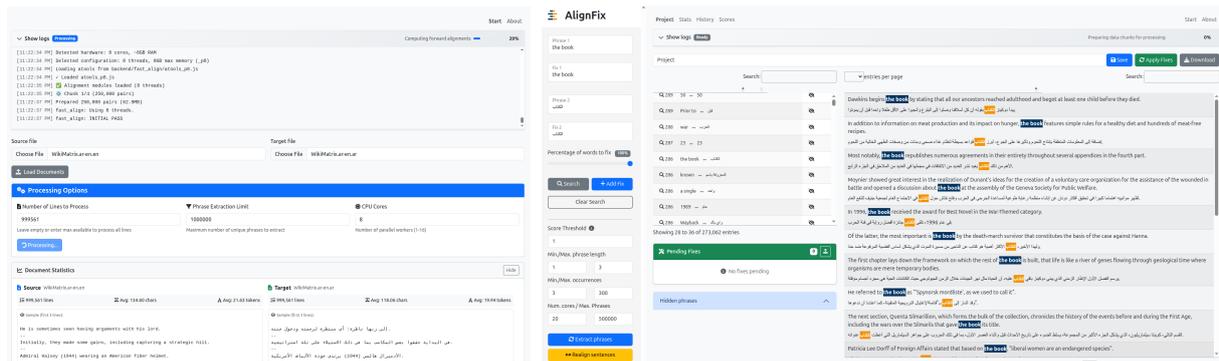


Figure 1: Left: Screenshot showing text alignment and phrase extraction. Right: Project view displaying extracted phrases. Both images are taken while working with the WikiMatrix Arabic–English corpus.

ALIGNFIX addresses this need by providing a tool for refining and augmenting parallel datasets through the extraction of aligned words and phrases, allowing for targeted, contextual interventions. The process works as follows:

- (i) the texts are tokenized to separate words from punctuation,
- (ii) (symmetric) word alignments are computed,
- (iii) phrase pairs that are translations of one another are extracted.

Fixes can be specified for these phrase pairs and applied selectively, allowing targeted adjustments only where intended. The tool maintains and incrementally updates the underlying word alignments when fixes are applied, thereby avoiding the costly recomputation of the alignments. The tool is designed to handle datasets up to one million samples efficiently and offers a user-friendly web interface. Figure 1 presents two screenshots of the interface. The left image illustrates text alignment and phrase extraction, while the right image shows the project view with the extracted phrases. Both screenshots were captured during work with the WikiMatrix Arabic–English corpus (Schwenk et al., 2021).

The primary use case for ALIGNFIX is thus to refine parallel corpora. However, it can also be used for other scenarios that require word-level alignment, such as augmenting corpora with glossary-enforced training data. In this work:

- we adapt and compile existing word-alignment and phrase-extraction methods so that they run directly in the browser, making them executable for everyone without installation, compilation, or technical expertise;

- we present ALIGNFIX, an open-source tool that allows for targeted data augmentation and refinement of parallel corpora and demonstrate its performance on different datasets;
- we demonstrate its practical utility in a low-resource domain scenario using two novel datasets of meteorological forecasts, which we make publicly available.

ALIGNFIX is available at <https://ifi-alignfix.uibk.ac.at> and is demonstrated in a supplementary video<sup>1</sup>. The source code<sup>2</sup> is provided under the Apache 2.0 open-source licence.

## 2 Related Work

Several toolkits have been developed to automate the cleaning and preparation of bitexts. OpusCleaner and OpusTrainer (Bogoychev et al., 2023) are widely adopted open-source toolkits that streamline downloading, preprocessing, and mixing data for large-scale neural MT. Similarly, OpusFilter (Aulamo et al., 2020) offers a modular toolbox for filtering, language identification, and alignment, allowing users to chain custom heuristic filters. For noise detection, Bicleaner and its successor Bicleaner AI (Zaragoza-Bernabeu et al., 2022) identify and discard noisy sentence pairs. While Bicleaner relies on heuristics, Bicleaner AI utilizes transformer-based models for a more accurate text classification. However, these approaches primarily function as filters. In low-resource scenarios, as highlighted by Marashian et al. (2025), data scarcity makes the rejection of "imperfect" sentence pairs undesirable. Discarding data that

<sup>1</sup>[https://youtu.be/F\\_7fyWc4vZo](https://youtu.be/F_7fyWc4vZo)

<sup>2</sup><https://github.com/alignfix/alignfix>

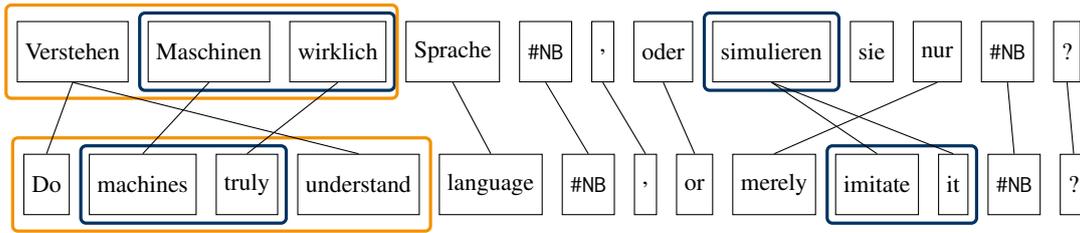


Figure 2: Alignment between a German and an English sentence with the consistent phrase pairs up to length 3.

contains recoverable errors can further starve an already data-poor system.

Beyond simple filtering, some tools provide functionalities to clean and fix problematic elements in corpora. Bifixer, part of the Bitextor project, focuses on technical repairs such as fixing encoding errors and removing near-duplicate sentence pairs (Ramírez-Sánchez et al., 2020). For the creation and management of alignments, SentAlign (Steingrímsson et al., 2023) utilizes LaBSE embeddings to identify semantically similar sentence pairs, employing dynamic programming for optimal alignment recovery. Once corpora are created, tools like InterText (Vondřička, 2014) provide a flexible editor for managing and manually aligning parallel texts.

While the aforementioned tools excel at either bulk filtering or sentence-level management, there is a lack of lightweight tools designed for corpus refinement and *contextual* word augmentation.<sup>3</sup> AlignFix addresses this by leveraging word alignments to allow targeted, manual replacements without discarding samples, thereby preserving valuable training data.

### 3 AlignFix

In this section, we describe our method and provide implementation details. In the first part, we describe the steps involved in extracting the phrase pairs from parallel corpora. In the second part, we discuss how fixes can be applied and how samples are augmented.

#### 3.1 Method Overview

In this section, we explain the individual steps of tokenization, alignment, and phrase extraction. All three components are implemented in (parallelized) C and compiled to WebAssembly (WASM), enabling efficient execution directly in the browser.

<sup>3</sup>OpusTrainer also offers data augmentation, but focuses on surface-level text manipulations (e.g., casing, all-caps) to improve model robustness rather than contextual refinement.

The resulting tokenized texts, word alignments, and extracted phrase pairs are persisted in an in-browser SQLite database.

**Tokenization** The computation of the word-alignments, requires tokenized texts (tokens are separated by blanks) as input. Therefore, as first step, we tokenize each sentence by explicitly separating punctuation from the surrounding text. Whenever punctuation is attached directly to a word without an intervening space, we insert a dedicated non-blank marker token (#NB) to ensure that the resulting tokenized text is reversible. For example, the sentence *The corpus is small, but valuable.* is tokenized to *The corpus is small #NB , but valuable #NB ..* This allows us to later reconstruct the original text (with possible fixes).

**Word-Alignments** To compute word alignments, we rely on the tokenized texts as produced in the preprocessing step. For the computation of the word-alignments, we use `fast_align` (Dyer et al., 2013). Figure 2 shows an example of word-alignments between two sentences.

We adapted the original C++ implementation and compiled it to WebAssembly using `emcc` (Zakai, 2011). This required several modifications: (i) we exposed the main function and key entry points to `emcc` so they could be invoked directly from JavaScript; (ii) we replaced the original OpenMP-based parallelisation with a WebAssembly-compatible setup using `pthread`s, enabling multi-threaded execution inside the browser; and (iii) we adjusted the build configuration to allow the model parameters to be loaded from in-memory buffers rather than from the local file system. These changes produced a browser-executable alignment tool with efficient parallel processing, allowing us to run alignment entirely in the browser. Symmetrization is carried out using `atools`, which we likewise compiled to WebAssembly following the same procedure.

**Phrases Extraction** Using the word alignments, we extract all sequences of aligned words up to a predefined length and record every occurrence of these phrases. To do this, we have implemented the method introduced in Och et al. (1999) and refined in Koehn et al. (2003) to extract *consistent* phrase pairs in C++ and compiled it to WebAssembly using emcc, enabling parallel execution via pthreads. This allows fast retrieval and inspection of their occurrences in the corpus. To make this process memory-efficient, the corpus is processed in batches, with batch sizes configurable based on available system memory.

We only collect *consistent* phrase pairs as they can be fully replaced without affecting other words that may occur in between otherwise. For the example illustrated in Figure 2, beside the single word pairs, the extracted phrases up to a maximum length of three would be:

1. ⟨Maschinen; machines⟩, ⟨wirklich; truly⟩, ⟨Sprache; language⟩, ⟨oder; or⟩, ⟨nur; merely⟩
2. ⟨Maschinen wirklich; machines truly⟩, ⟨simulieren; imitate it⟩.
3. ⟨Verstehen Maschinen wirklich; Do machines truly understand⟩

We do not include ⟨Maschinen wirklich Sprache; machines truly understand language⟩, because the word *understand* is not aligned to any word in *Maschinen wirklich Sprache*. We trim punctuation and non-blank symbols (e.g. we treat ⟨Sprache #NB; language #NB⟩ as ⟨Sprache; language⟩). We also remove pairs only consisting of punctuation symbols, e.g. ⟨?; ?⟩, ⟨,; ,⟩ or ⟨#NB ;; #NB ‘;’⟩.

**Managing Extraction Scale** There is one drawback of extracting all phrases in the corpus: the number of extracted phrases grows rapidly. For example, in a corpus of 100k sentences, the number of phrase pairs can easily exceed one million with a maximum phrase length of three. Storing all of them would introduce substantial overhead. Therefore, we allow the user to specify an upper limit on the number of phrase pairs to collect (default: 500k). Based on the maximum phrase length, we determine an appropriate batch size for processing the corpus, so that we can guarantee to stay below a peak memory usage of 4GB<sup>4</sup>. After each batch, we check whether the number of collected phrase pairs exceeds the user-defined limit. If so, we prune

phrase pairs with a single occurrence within the processed batch. If the limit is still exceeded, we iteratively remove pairs with two occurrences, three occurrences, and so forth, until the total number of phrase pairs falls below the threshold. We then proceed with the next batch.<sup>5</sup> Users who are interested in phrase pairs that rarely occur can still search for them directly in the full corpus.

## 3.2 Corpus Augmentation and Refinement

This section describes the implemented features for corpus augmentation and refinement.

**Data Augmentation** ALIGNFIX supports contextual data augmentation. The user can duplicate existing sentence pairs and selectively replace aligned words to generate new training examples. For instance, whenever the word *car* occurs, it can be substituted with *automobile*, along with the corresponding replacement in the target language. Similarly, even substitutions that change the meaning but still yield coherent sentences can be applied when appropriate. For example, in a sentence such as “*She touched her ear.*” the word *ear* may be replaced with *nose* to form “*She touched her nose.*” While not all contexts support such substitutions, ALIGNFIX enables users to perform them in a controlled manner, thereby enriching the corpus with additional valid sentence pairs. This enables controlled lexical diversification and allows users to introduce examples for terms that are underrepresented or entirely absent from the original corpus.

**Refinement** Word alignments are leveraged to enable users to correct phrases in both the source and target sentences. In ALIGNFIX, these corrections can be applied either to every occurrence of a given phrase, or to a selected subset of occurrences across the corpus. For instance, based on the phrase pairs extracted in Figure 2, a user could choose to replace *Do machines truly understand* with *Can machines understand* throughout the corpus, wherever it aligns with *Verstehen Maschinen wirklich*. Computing word alignments is computationally expensive. Therefore, ALIGNFIX preserves alignment consistency when replacements are applied. In cases where a single token is replaced by multiple tokens, the original aligned to-

<sup>4</sup>See discussion on memory considerations in Section 4

<sup>5</sup>In the worst case, this procedure may discard (rare but important) phrase pairs that would have appeared  $\text{corpus\_size} / \text{batch\_size}$  times across the entire dataset but this pruning is necessary to avoid hitting memory limits.

Corpus	Size	Cores	Tokenization (T)		Alignment (A)		Phrase Extraction (P)		DB Insert Time (s)	Efficiency (s / 1k lines)		
			Time (s)	Mem (MB)	Time (s)	Mem (MB)	Time (s)	Mem (MB)		T	A	P
$S_{6k}$	6k	8	0.67	9.70	13.16	1.10	7.57	125.00	1.22	0.11	2.19	1.26
$S_{6k}$	6k	16	0.66	11.44	7.16	7.51	5.19	19.94	0.57	0.11	1.15	0.84
$D_{100k}$	100k	8	5.39	121.96	223.92	108.61	66.90	49.57	13.34	0.05	2.24	0.67
$D_{100k}$	100k	16	3.14	113.07	103.94	109.56	86.24	499.80	24.48	0.03	1.04	0.86
$W_{418k}$	418k	8	26.07	298.93	579.44	310.97	412.40	984.75	44.02	0.06	1.39	0.99
$W_{418k}$	418k	16	9.92	358.96	326.05	9.54	306.72	872.06	68.24	0.02	0.78	0.73
$W_{1M}$	1M	8	78.33	843.10	2073.67	499.00	1329.66	2082.20	162.34	0.08	2.07	1.33
$W_{1M}$	1M	16	40.99	1005.00	1182.35	384.35	985.88	2332.95	133.19	0.04	1.18	0.99

Table 1: Benchmark results for tokenization, alignment, and phrase extraction across corpora, including efficiency normalized per 1k sentence pairs.

kens are distributed across the new tokens to maintain the alignment structure.<sup>6</sup>

Even after the pruning of phrases described above, the remaining set of phrase pairs may still be large. Therefore, to facilitate the identification of potential error candidates, ALIGNFIX provides two additional filtering mechanisms that can substantially reduce the number of phrase pairs: (i) ignore known phrase pairs. (ii) filter based on translation quality scores. In (i), the user may upload a list of phrase pairs to be excluded from extraction which could for example be a list of verified translations. In (ii), the user may upload quality scores for the sentence pairs in the corpus. Scores should range from 0 (low-quality translation) to 1 (high-quality translation) – for instance, 0 for back-translated data and 1 for expert translations. The user can then define a threshold: only phrase pairs occurring exclusively in translations with a score below the threshold are retained.<sup>7</sup>

## 4 Experimental Setup

The aim of our experiments is to systematically evaluate the effectiveness of the tool introduced in this work, both with respect to its operational performance and the impact it can have on downstream machine translation.

### 4.1 Tool Performance

To assess the performance, we conducted experiments across multiple corpora and computing environments. In Table 1 we report the results.

<sup>6</sup>This heuristic alignment repair strategy provides a practical approximation; more accurate alignments could be obtained by recomputing them after each replacement.

<sup>7</sup>ALIGNFIX provides experimental metrics to estimate translation quality. These metrics are currently under development and have not yet been fully validated; their evaluation and refinement are planned as part of future work.

**Benchmark** The benchmark suite covers three publicly available corpora of different sizes and linguistic characteristics.  $S_{6k}$  represents the 6k-sentence Seed corpus (Maillard et al., 2023) for Italian (Ferrante, 2024) and serves as a controlled small-scale reference for evaluating baseline throughput on a low-volume dataset.  $D_{100k}$  corresponds to a 100k-sentence heterogeneous general-domain corpus for Uzbek–Karakalpak (Mamasaidov and Shopulatov, 2024), enabling analysis on medium-sized data. To assess performance on substantially larger material,  $W_{418k}$  uses a 418k-sentence subset of the WikiMatrix German–Spanish corpus (Schwenk et al., 2021). Finally,  $W_{1M}$  contains one million sentence pairs from the WikiMatrix Arabic–English corpus (Schwenk et al., 2021), selected as a large-scale and cross-family dataset to stress-test the tool under substantial data volume. Together, these corpora enable a comprehensive evaluation of runtime, memory usage, and throughput across variation in size and languages.

**Hardware and Performance** The performance was evaluated on two systems running Chromium v142. The primary Mini-PC features a 12th Gen Intel® Core™ i9-12900H (20 cores) with 64 GB DDR4-3000 memory, used for 16-thread tests with 16 GB WASM memory and a 4 GB JS heap. A lower-resource laptop<sup>8</sup> with an 8th Gen Intel® Core™ i7-8565U (4 cores, 8 threads) and 40 GB DDR4-2667 memory was used for 8-thread tests with 8 GB WASM memory and a 4 GB JS heap. Across both machines, the runtime data demonstrate that the full pipeline scales efficiently with available parallelism, and that large corpora are processable within browser constraints.

<sup>8</sup>Due to OS-level power scaling, the effective CPU frequency (and thus execution time) may vary on battery or under background load.

**Scalability and Efficiency** The efficiency values in Table 1 show that processing time (in seconds) per 1k sentence pairs decreases as corpus size increases, demonstrating favorable scaling of the pipeline. From the 6k corpus to the 100k corpus, tokenization time drops from 0.11s to 0.05s per 1k lines, while alignment throughput remains effectively constant, changing only slightly from 2.19s to 2.24s despite the larger dataset. Phrase extraction also becomes more efficient at scale, decreasing from over 1s per 1k lines on small data to 0.67s for the 100k corpus and reaching 0.73s per 1k lines on the 418k corpus, before stabilizing on the 1M corpus. These results indicate that one-time initialization and model-loading overheads are quickly amortized, and that the pipeline benefits substantially from multithreaded execution.

**Memory Considerations** The memory measurements in Table 1 report only JS heap usage (usedJSHeapSize), which is managed by V8’s garbage collector. WebAssembly linear memory, allocated separately, is not included; although its size can be obtained via WebAssembly.Memory.buffer.byteLength, including it would mix separate memory regions and could be misleading. The WASM modules were compiled with a maximum linear memory of 1 GB per CPU core, matching the 8-thread (8 GB) and 16-thread (16 GB) configs used in our experiments.

Across all experiments, JS heap usage remained below 2.5 GB, safely within the browser’s 4 GB limit, with minor fluctuations due to garbage collection. Combined WASM+JS memory usage therefore stayed below the effective upper bounds of 12 GB (4+8) or 20 GB (4+16), even on the largest 1M-sentence corpus. Overall, memory consumption grows moderately with corpus size.

## 4.2 Impact of Targeted Corrections

To illustrate the applicability of ALIGNFIX in a realistic low-resource scenario, we conducted experiments on weather forecast texts provided by the *Amt für Meteorologie und Lawinenwarnung* of the Autonomous Province of Bolzano – South Tyrol<sup>9</sup>. The corpus consists of 689 parallel Ladin (Val Badia)–German (VB–DE) reference translations<sup>10</sup> and additional 15,969 VB-only weather forecast

<sup>9</sup>Datasets released by the authors with permission of the *Amt für Meteorologie und Lawinenwarnung*.

<sup>10</sup><https://huggingface.co/datasets/sfrontull/south-tyrol-weather-1ld-deu>

Model	BLEU	COMET
Ladin (Val Badia) → German		
gemi-ni-2.5-flash-lite	15.9±1.2	67.0
Helsinki-NLP/opus-mt-it-de fine-tuned with backtranslations + 138 fixes with ALIGNFIX	17.0±1.2 <b>18.6±1.2</b>	67.8 69.6
German → Ladin (Val Badia)		
Helsinki-NLP/opus-mt-de-it fine-tuned with backtranslations + 138 fixes with ALIGNFIX	30.5±1.6 <b>32.3±1.5</b>	55.7 56.6

Table 2: Comparison of translation quality ( $\mu \pm 95\%$  CI) for German–Ladin weather forecasts, highlighting the gains achieved by applying 138 targeted corrections.

texts<sup>11</sup>. This setup reflects a typical low-resource condition, where the lack of parallel corpora necessitates the synthesis of training data through backtranslation.

**Data Augmentation via Backtranslation** We generated a synthetic parallel dataset by translating the 15,969 Ladin monolingual texts into German using Gemini 2.5 Flash-Lite (Comanici et al., 2025) in a zero-shot setting.<sup>12</sup> Following the backtranslation paradigm, we then fine-tuned a DE → VB model on this synthetic corpus. As no pre-trained German–Ladin model is available and Ladin is closely related to Italian, we used the Helsinki-NLP/opus-mt-de-it model (Tiedemann et al., 2024; Tiedemann and Thottingal, 2020) as base model for this experiment. The model was trained for up to 20 epochs with a batch size of 8 and learning rate  $2 \cdot 10^{-5}$ , with early stopping set to 3 epochs.

**Targeted Corrections** After establishing this baseline, we used ALIGNFIX to identify and correct systematic errors in the German backtranslations produced by the large language model (LLM). In total, we applied 138 targeted phrase-level fixes and fine-tuned the model on this refined data (with the same configuration).

To quantify the scale of the applied corrections: the 138 fixes modified 6,677 of the 15,969 synthetic sentences (41.8%). In total, ALIGNFIX introduced 56,906 character-level edits, corresponding to an average of 85 edits per changed sentence and an overall edit intensity of 37.3% relative to the

<sup>11</sup><https://huggingface.co/datasets/sfrontull/south-tyrol-weather-1ld>

<sup>12</sup>Prompt: Translate the following sentence from Ladin to German: <Ladin\_Text>

original text. These figures highlight that a moderate number of targeted interventions (requiring roughly 1–2 hours of manual effort) can propagate broadly across a domain corpus, producing substantial systematic improvements. A complete list of the applied fixes is provided in Appendix A.

**Results** Table 2 reports the BLEU and COMET scores for the evaluated models. The BLEU scores were computed using sacreBLEU (Post, 2018). We employed paired bootstrap resampling (--paired-bs) to assess statistical significance; values in bold denote a significant improvement over the baseline. The COMET scores were computed using the Unbabel/wmt22-comet-da (Rei et al., 2022) model. The results clearly demonstrate the positive effect of our interventions on translation quality. Fine-tuning on the backtranslated data results in 30.5 BLEU. Applying targeted corrections with ALIGNFIX increases performance to 32.3 BLEU (+1.8 BLEU). Despite Ladin being unsupported by COMET, the 0.9 increase suggests improvements in semantic adequacy as well.

We also examined the effect of these fixes in the opposite translation direction (Ladin → German). Fine-tuning on the synthetic corpus already improves over the zero-shot Gemini baseline (+1.1 BLEU and +0.8 COMET). The refined corpus (in this case, with target-side corrections) yields further improvements, reaching 18.6 BLEU and 69.6 COMET (additional +1.6 BLEU and +1.8 COMET).

## 5 Conclusion

We presented ALIGNFIX, a tool for improving parallel corpora by leveraging word alignments to propagate corrections consistently across sentence pairs. ALIGNFIX enables users to modify individual tokens or phrases while automatically maintaining alignment integrity, even when a single token is replaced by multiple tokens. Through its combination of browser-executable algorithms and phrase-based repair operations, the system offers a flexible, scalable, and user-friendly framework for enhancing translation corpora across a wide range of practical scenarios.

Our experiments highlight an important mechanism in low-resource machine translation where training data is synthesized. Errors in the synthetic texts can systematically remove or distort domain-specific terminology. If key terms are mistranslated or omitted during backtranslation, they never appear aligned with their correct counterparts in

the synthetic parallel data. As a result, the model fails to learn these correspondences and may later hallucinate or substitute more frequent but incorrect alternatives at inference time. By restoring correct terminology and phrase structure on the synthetic source side, ALIGNFIX allows to reintroduce these missing lexical links, strengthening the learned cross-lingual mapping and reducing errors in translation.

**Future Work** We consider the automated identification of potential errors to be a crucial feature. While the current functionality supports user-defined lists of phrase-pairs to exclude (e.g., to filter out correct pairs that do not require review), this is not a scalable solution. Potential errors could also be detected intrinsically. In future work, we would like to explore such methods and provide users with suggestions for possible fixes to substantially reduce the amount of manual work required.

Our current experiments and implementation support corpora of up to approximately one million sentence pairs. Larger datasets may exceed the memory limitations of the underlying database and browser execution environment. In future work, we further aim to improve memory prediction, adapt batch sizing to corpus characteristics, and optimize I/O and storage efficiency (e.g., OBFS) to handle even larger corpora.

## Limitations

While ALIGNFIX is largely language-agnostic, the current implementation relies on whitespace-based tokenization and existing word alignment tools, which limit direct applicability to languages without explicit word boundaries (e.g., Chinese, Japanese, Thai). Lightweight, WASM-compatible tokenization strategies could be integrated to support scripts without whitespace segmentation, applying language-specific tokenizers only when necessary while preserving a unified, aligner-compatible output format.

ALIGNFIX assumes pre-aligned parallel data. This design choice reflects the primary target use case, where alignment is implicitly provided by construction. In scenarios where sentence alignment is unavailable or noisy, additional preprocessing is required. Moreover, the effectiveness of this tool depends critically on the quality of word alignments. If a corpus is too small to support robust statistical alignment, the approach may fail to produce satisfactory results.

## Acknowledgments

This work was carried out as part of the research project *Intelligent Writing Assistant for Ladin* at the University of Innsbruck, in collaboration with the Ladin Cultural Institute “Micurá de Rù”. We thank the reviewers for their valuable comments and the Amt für Meteorologie und Lawinenwarnung of the Autonomous Province of Bolzano/Bozen – South Tyrol for providing the data.

## References

- Mikko Aulamo, Sami Virpioja, and Jörg Tiedemann. 2020. [OpusFilter: A Configurable Parallel Corpus Filtering Toolbox](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 20–26. Association for Computational Linguistics.
- Nikolay Bogoychev, Jelmer van der Linde, Graeme Nail, Barry Haddow, Jaume Zaragoza-Bernabeu, Gema Ramírez-Sánchez, Lukas Weymann, Tudor Nicolae Mateiu, Jindřich Helcl, and Mikko Aulamo. 2023. [OpusCleaner and OpusTrainer: Open Source Tools for Training Machine Translation and Large Language Models](#). *arXiv preprint arXiv:2311.14838*.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobson, Idan Szpektor, Nan-Jiang Jiang, and 3416 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. [A Simple, Fast, and Effective Reparameterization of IBM Model 2](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.
- Edoardo Ferrante. 2024. [A High-quality Seed Dataset for Italian Machine Translation](#). In *Proceedings of the Ninth Conference on Machine Translation*, pages 567–569, Miami, Florida, USA. Association for Computational Linguistics.
- Fei Gao, Jinhua Zhu, Lijun Wu, Yingce Xia, Tao Qin, Xueqi Cheng, Wengang Zhou, and Tie-Yan Liu. 2019. [Soft Contextual Data Augmentation for Neural Machine Translation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5539–5544, Florence, Italy. Association for Computational Linguistics.
- Mitchell A Gordon, Kevin Duh, and Jared Kaplan. 2021. [Data and Parameter Scaling Laws for Neural Machine Translation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5915–5922, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jiatao Gu, Hany Hassan, Jacob Devlin, and Victor O.K. Li. 2018. [Universal Neural Machine Translation for Extremely Low Resource Languages](#). 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 344–354, New Orleans, Louisiana. Association for Computational Linguistics.
- Nuno M. Guerreiro, Elena Voita, and André Martins. 2023. [Looking for a Needle in a Haystack: A Comprehensive Study of Hallucinations in Neural Machine Translation](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1059–1075, Dubrovnik, Croatia. Association for Computational Linguistics.
- Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. 2018. [Iterative Back-Translation for Neural Machine Translation](#). In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 18–24, Melbourne, Australia. Association for Computational Linguistics.
- Huda Khayrallah and Philipp Koehn. 2018. [On the Impact of Various Types of Noise on Neural Machine Translation](#). In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83, Melbourne, Australia. Association for Computational Linguistics.
- Sosuke Kobayashi. 2018. [Contextual Augmentation: Data Augmentation by Words with Paradigmatic Relations](#). 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 452–457, New Orleans, Louisiana. Association for Computational Linguistics.
- Philipp Koehn and Rebecca Knowles. 2017. [Six Challenges for Neural Machine Translation](#). In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.
- Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003. [Statistical Phrase-Based Translation](#). In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 127–133.
- Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzmán.

2023. **Small Data, Big Impact: Leveraging Minimal Data for Effective Machine Translation**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.
- Mukhammadsaid Mamasaidov and Abror Shopulatov. 2024. **Open Language Data Initiative: Advancing Low-Resource Machine Translation for Karakalpak**. In *Proceedings of the Ninth Conference on Machine Translation*, pages 606–613, Miami, Florida, USA. Association for Computational Linguistics.
- Ali Marashian, Enora Rice, Luke Gessler, Alexis Palmer, and Katharina von der Wense. 2025. **From Priest to Doctor: Domain Adaptation for Low-Resource Neural Machine Translation**. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7087–7098, Abu Dhabi, UAE. Association for Computational Linguistics.
- Franz Josef Och, Christoph Tillmann, and Hermann Ney. 1999. **Improved Alignment Models for Statistical Machine Translation**. In *1999 Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*.
- Matt Post. 2018. **A Call for Clarity in Reporting BLEU Scores**. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.
- Gema Ramírez-Sánchez, Jaume Zaragoza-Bernabeu, Marta Bañón, and Sergio Ortiz-Rojas. 2020. **Bifixer and Bicleaner: Two Open-Source Tools to Clean Your Parallel Data**. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT 2020)*, pages 1875–1879. European Association for Machine Translation.
- Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. **COMET-22: Unbabel-IST 2022 submission for the metrics shared task**. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021. **WikiMatrix: Mining 135M Parallel Sentences in 1620 Language Pairs from Wikipedia**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361, Online. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. **Improving Neural Machine Translation Models with Monolingual Data**. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.
- Steinþór Steingrímsson, Hrafn Loftsson, and Andy Way. 2023. **SentAlign: Accurate and Scalable Sentence Alignment**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 256–263. Association for Computational Linguistics.
- Jörg Tiedemann, Mikko Aulamo, Daria Bakshandaeva, Michele Boggia, Stig-Arne Grönroos, Tommi Niemen, Alessandro Raganato, Yves Scherrer, Raúl Vázquez, and Sami Virpioja. 2024. **Democratizing neural machine translation with OPUS-MT**. *Language Resources and Evaluation*, 58(2):713–755.
- Jörg Tiedemann and Santhosh Thottingal. 2020. **OPUS-MT — Building open translation services for the World**. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal.
- Pavel Vondříčka. 2014. **Aligning Parallel Texts with InterText**. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC 2014)*, pages 1875–1879. European Language Resources Association (ELRA).
- Xing Wu, Shangwen Lv, Liangjun Zang, Jizhong Han, and Songlin Hu. 2019. **Conditional BERT Contextual Augmentation**. In *Computational Science – ICCS 2019: 19th International Conference, Faro, Portugal, June 12–14, 2019, Proceedings, Part IV*, page 84–95, Berlin, Heidelberg. Springer-Verlag.
- Alon Zakai. 2011. **Emscripten: an LLVM-to-JavaScript compiler**. In *Proceedings of the ACM International Conference Companion on Object Oriented Programming Systems Languages and Applications Companion, OOPSLA ’11*, page 301–312, New York, NY, USA. Association for Computing Machinery.
- Jaume Zaragoza-Bernabeu, Marta Bañón, Gema Ramírez-Sánchez, and Sergio Ortiz-Rojas. 2022. **Bicleaner AI: Bicleaner Goes Neural**. In *Proceedings of the 13th Language Resources and Evaluation Conference (LREC 2022)*, pages 824–831.
- Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. **Transfer Learning for Low-Resource Neural Machine Translation**. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1568–1575, Austin, Texas. Association for Computational Linguistics.

## A Detailed Phrase-Level Refinements

Table 3 presents all 138 targeted interventions, their frequency in the backtranslated data, and a comparison of baseline German output with our refined translations. Note, for example, the different translations of the Ladin word *niores* (clouds) in the German texts hallucinated by the LLM, ranging from *Blumen* to *Mädchen*.

	Ladin (VB)	Original (DE)	Fixed (DE)	#		Ladin (VB)	Original (DE)	Fixed (DE)	#	
1	indô	wieder	erneut	915	70	Tres plü variabl	Drei variabel	Zunehmend unbeständig	60	
2	sovënz	oft	meist	694	71	cuntra le Südtirol	gegen Südtirol	Richtung Südtirol	60	
3	danmisdé	am Nachmittag	am Vormittag	431	72	niores	Mädchen	Wolken	59	
4	presciun bassa	Tiefdruckgebiet	Tief	370	73	Tres plü da nio	Drei mehr als nichts	Zunehmend bewölkt	59	
5	niores	Nebel	Wolken	367	74	niores	Nächte	Wolken	58	
6	temperatôres mascimes	Höchsttemperaturen	Höchstwerte	346	75	Sorëdl y niores	Sonne und Blumen	Sonne und Wolken	58	
7	Da doman	Von morgen an	In der Früh	316	76	indlunch	wieder	überall	58	
8	Sorëdl y niores	Sonnenschein und Blumen	Sonne und Wolken	302	77	Da sorëdl	Von der Sonne	Sonnig	57	
9	en pert	teilweise	teils	279	78	manco tömies	weniger dicht	weniger feucht	55	
10	Domisdé	Heute	Nachmittag	277	79	Da sorëdl y cialt	Von Sonne und Wärme	Sonnig und warm	54	
11	bel	schön	freundlich	269	80	da doman	von morgen	in der Früh	54	
12	instabil	instabil	unbeständig	259	81	sorëdl	oben	Sonne	53	
13	en gran pert	größtenteils	überwiegend	241	82	süd dla provinzia	Süden der Provinz	Süden des Landes	53	
14	Na presciun alta	Ein hoher Druck	Ein Hoch	238	83	döt sarëgn	ganz klar	wolkenlos	52	
15	pert dla provinzia	Teil der Provinz	Teil des Landes	236	84	dadoman	morgen	in der Früh	51	
16	limit dla nëi	Schneegrenze	Schneefallgrenze	233	85	vënt da nord	Wind aus Norden	Nordwind	50	
17	La presciun alta	Der hohe Druck	Das Hoch	217	86	Dantadöt sorëdl	Gib mir die Sonne	Überwiegend sonnig	49	
18	da nio	von Schnee	bewölkt	209	87	gnanca na niora	nicht einmal eine Wolke	wolkenlos	43	
19	por intant	vorerst	vorübergehend	206	88	Valgamia	Wir gehen	Recht	42	
20	Sön munt	Auf dem Berg	Auf den Bergen	203	89	Plülere sorëdl	Die Schwestern	Recht sonnig	42	
21	Da sorëdl y da nio	Von Sonne und von Schnee	Sonne und Wolken	202	90	niores zënza faz- iun	Wolken ohne Nieder- schlag	harmlose Wolken	38	
22	mascimes	Höchsttemperaturen	Höchstwerte	199	91	Sön la	Auf der Alpenhaupt- tkamm	Am Alpenhauptkamm	37	
23	niores	Blumen	Wolken	197	92	Sön la Ciadëna	Auf der Alpenhaupt- tkamm	Am Alpenhauptkamm	37	
24	Ciadëna	Alpenkette	Alpenhauptkamm	184	93	niores a gröm	Wolken	Quellwolken	37	
25	i crëps	den Gipfeln	den Bergen	179	94	Dër da nio	Sehr gut	Sehr bewölkt	37	
26	Dadoman	Morgen	In der Früh	178	95	Da nio	Von nichts	Bewölkt	36	
27	Tres	Drei	Zunehmend	177	96	dantadöt	hauptsächlich	überwiegend	34	
28	Da doman	Von morgen	In der Früh	174	97	Ciadëna centrala dles	der zentralen Alpenhaupt- tkamm	dem Alpenhauptkamm	33	
29	la Ciadëna	Kette	Alpenhauptkamm	172	98	dantadöt	vor allem	überwiegend	32	
30	da sorëdl	von der Sonne	sonnig	171	99	Plü variabl	Mehr variabel	Wechselhafter	30	
31	Domisdé	Morgen	Am Nachmittag	166	100	meste	meistens	mild	30	
32	manco da nio	weniger von nichts	weniger bewölkt	161	101	sön la	auf der Alpenhaupt- tkamm	am Alpenhauptkamm	29	
33	bonamënter	meist	voraussichtlich	144	102	sön la Ciadëna	auf der Alpenhaupt- tkamm	am Alpenhauptkamm	29	
34	Da sorëdl y	Von der Sonne und	Sonnig und	137	103	gnanca na niora	nicht einmal eine Stunde	wolkenlos	29	
35	Domisdé	Vormittags	Nachmittags	132	104	zënza fazium	ohne Auflösung	harmlos	28	
36	Sön la	Auf der zentralen Alpen- hauptkamm	Am Alpenhauptkamm	132	105	vignitant	bald	zeitweise	28	
37	moscedoz	Mix aus	Mischung aus	131	106	Danmisdé	Morgen	Am Vormittag	28	
38	Domisdé	Heute Morgen	Am Nachmittag	130	107	condiziuns	Bedingungen	Verhältnisse	27	
39	te tröc posc	an vielen Orten	verbreitet	129	108	naota	mehr	zunächst	26	
40	tröp	viel	viel	129	109	niores a gröm	Haufen	Quellwolken	25	
41	meste	Nebel	mild	128	110	dër meste	sehr traurig	sehr mild	24	
42	arbassa	sinken	gehen zurück	125	111	Tröpes niores	Kleine Tropfen	Viele Wolken	24	
43	Domisdé	Übermorgen	Am Nachmittag	122	112	Dër da sorëdl	Sehr von Sonne	Sehr sonnig	22	
44	I valurs mascimai	Die maximalen Werte	Höchstwerte	117	113	Ciarü alt	Schau hoch	Hochnebel	20	
45	Sön i crëps	Auf den Gipfeln	Auf den Bergen	114	114	aboc sorëdl	viel Sonne	zeitweise Sonne	19	
46	da nio	von nichts	bewölkt	106	115	plü tömia	kältere	feuchtere	19	
47	raiun dles Alpes	Alpenregion	Alpen	105	116	Variabl y da nio	Variable von nichts	Wechselhaft und be- wölkt	19	
48	niores	Schneefälle	Wolken	102	117	bel plan	freundlich langsam	allmählich	18	
49	niores a gröm	größere Wolken	Quellwolken	97	118	bel plan	gut	allmählich	18	
50	minimes	Tiefsttemperaturen	Tiefstwerte	97	119	aboc	meistens	zeitweise	17	
51	cresta de confin	dem Kamm	dem Alpenhauptkamm	96	120	naota	noch	zunächst	17	
52	Ciadëna centrala	zentralen Alpenhaupt- tkamm	Alpenhauptkamm	95	121	niores	Schneefelder	Wolken	17	
53	de transiziun	Übergangsdruck	Zwischenhoch	90	122	stopa sovënz la	beeinträchtigen oft die	behindern oft die	14	
54	ciarü	klar	Hochnebel	85	123	Sorëdl y niores a	Sonnenschein und Schnee in	Sonne und Quell- wolken	14	
55	y danmisdé	und übermorgen	und am Vormittag	84	124	niores	Jüngeren	Wolken	14	
56	niores	Berge	Wolken	83	125	niores a slaiër	Wolken zum Anpflanzen	Schleierwolken	12	
57	tömia	kühle	feuchte	82	126	Valgamia da sorëdl	Recht sonnig aus	Recht sonnig	12	
58	Tres plü instabil	Drei instabiler	Zunehmend unbeständig	77	127	banc de ciarü	Schneebänke	Nebelfelder	12	
59	domisdé	heute Morgen	heute Nachmittag	76	128	indlunch	später	überall	12	
60	romagn variabl	bleibt variabel	bleibt wechselhaft	76	129	Da nio	Von Schnee	Bewölkt	12	
61	Al romagn variabl	Es bleibt variabel	Es bleibt wechselhaft	76	130	Sön la Ciadëna centrala dles Alpes	Auf der zentralen Alpen- hauptkamm der Alpen	Am Alpenhauptkamm	11	
62	plöiüdes	Schauern	Regenschauern	76	131	niores a gröm	Wolken in Haufen	Quellwolken	10	
63	Mioramënt dl	Erinnerung an die Zeit	Wetterbesserung	73	132	Dantadöt da nio	Dank von nichts	Wolken überwiegen	10	
64	no	Schnee	Wolke	72	133	ciarü alt	klare Höhe	Hochnebel	9	
65	y cialt	und Wärme	und Warm	71	134	gnanca na niora	nicht einmal eine Wolke	wolkenlos	9	
66	Le tēmp	Die Zeit	Das Wetter	69	135	a gröm	Quellwolken Haufen	Quellwolken	8	
67	Sö por munt	Oben auf dem Berg	Auf den Bergen	69	136	N pice mioramënt	Ein kleinerer Fortschritt	Leichte	Wet- terbesserung	6
68	Sorëdl y niores	und Blumen	und Wolken	66	137	Da sorëdl y da nio	Von Sonne und von Nichts	Sonne und Wolken	6	
69	bones condiziuns	guten Bedingungen	gute Verhältnisse	62	138	Na presciun alta temporanea	Ein hoher temporärer Gefängnisaufenthalt	Ein Zwischenhoch	5	

Table 3: All 138 fixes applied to the synthesised Ladin–German corpus.

# PromptLab: A Collaborative Platform for Prompt Engineering and Dataset Curation

Maged S. Al-shaibani<sup>1</sup> Zaid Alyafeai<sup>2</sup> Dania Refai<sup>3</sup> Nawaf Alomari<sup>3</sup> Ahmed Ashraf<sup>3</sup>  
Mais Alheraki<sup>3</sup> Mustafa Alturki<sup>4</sup> Hamzah Luqman<sup>1,3</sup> Irfan Ahmad<sup>1,3</sup>

<sup>1</sup>SDAIA-KFUPM JRC for AI <sup>2</sup>KAUST <sup>3</sup>KFUPM <sup>4</sup>HUMAIN AI

{maged.alshaibani, irfan.ahmad, hluqman}@kfupm.edu.sa

zaid.alyafeai@kaust.edu.sa, malturki@thefutureai.co

{g202391270, g201931050, g202411740, g202401480}@kfupm.edu.sa

## Abstract

PromptLab is a web-based platform for collaborative prompt engineering across diverse natural language processing tasks and datasets. The platform addresses primary challenges in prompt development, including template creation, collaborative review, and quality assurance through a comprehensive workflow that supports both individual researchers and team-based projects. PromptLab integrates with HuggingFace and provides AI-assisted prompt generation via OpenRouter<sup>1</sup>, and supporting real-time validation with multiple Large Language Models (LLMs). The platform features a flexible templating system using Jinja2, role-based project management, peer review processes, and supports programmatic access through RESTful APIs. To ensure data privacy and support sensitive research environments, PromptLab includes an easy CI/CD pipeline for self-hosted deployments and institutional control. We demonstrate the platform’s effectiveness through two evaluations: a controlled comparison study with six researchers across five benchmark datasets and 13 models with 90 prompts; and a comprehensive case study in instruction tuning research, where over 350 prompts across 80+ datasets have been developed and validated by multiple team members. The platform is available at <https://promptlab.up.railway.app> and the source code is available on GitHub at <https://github.com/KFUPM-JRC AI/PromptLab>.

## 1 Introduction

Prompt engineering is a fundamental technique for effectively utilizing large language models across diverse natural language processing tasks (Minaee et al., 2024). The practice of designing natural language instructions to guide model behavior has proven notable for achieving performance gains in zero-shot and few-shot settings. However, the process of creating, refining, and managing prompts at

<sup>1</sup><https://openrouter.ai/>

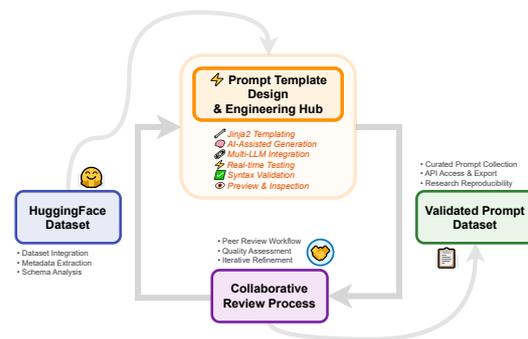


Figure 1: PromptLab General Pipeline

scale presents significant challenges, particularly for collaborative research environments where consistency, quality, and reproducibility are paramount (Schulhoff et al., 2024b; Sahoo et al., 2024a).

Recent comprehensive surveys have catalogued over 58 distinct prompting techniques (Schulhoff et al., 2024a) and established systematic taxonomies for prompt engineering methodologies (Sahoo et al., 2024b), underscoring the maturation of prompt engineering as a research discipline while revealing the complexity and diversity of approaches required for different tasks and domains. The growing sophistication of prompting techniques demands platforms that can support collaborative development, systematic evaluation, and reproducible research practices. As research teams increasingly work across institutional boundaries and language barriers, the need for a comprehensive collaborative environment becomes even more demanding.

The challenges facing prompt engineering research extend beyond individual technical difficulties to encompass broader issues of research coordination, quality control, and knowledge sharing. Current approaches to prompt development often rely on ad-hoc methodologies and informal sharing mechanisms that limit reproducibility and hinder

scientific progress. The lack of standardized workflows for collaborative prompt development has led to duplicated efforts, inconsistent quality standards, and missed opportunities for cross-institutional research collaboration.

Existing prompt engineering tools primarily focus on individual use cases and lack comprehensive support for collaborative workflows. While platforms like PromptSource (Bach et al., 2022) have demonstrated the value of structured prompt creation environments, they often fall short in providing the project management, peer review, and real-time evaluation proven useful for large-scale research initiatives. PromptLab addresses these limitations by providing a comprehensive platform designed specifically for collaborative prompt engineering. It combines flexible prompt template development with robust project management features. The platform’s design philosophy centers on three core principles. First, prompt creation should be accessible to researchers with varying technical backgrounds and different domains while maintaining the expressiveness necessary for complex tasks. This requires carefully designed interfaces that maintain a balance between simplicity and templating power. Second, collaborative prompt development requires advanced workflow management that goes beyond simple sharing mechanisms to include role-based access control, systematic review processes, and progress tracking. Third, quality assurance in prompt engineering demands automated validation and human expertise, live testing with different LLMs, and an iterative refinement workflow.

## 2 Background & Related Work

The development of PromptLab builds upon extensive research in prompt engineering and instruction tuning. Below, we briefly survey this literature.

### 2.1 Prompt Engineering and Optimization Techniques

Prompt engineering has rapidly evolved from basic few-shot learning (Brown, 2020) into advanced optimization frameworks (Schulhoff et al., 2024a; Sahoo et al., 2024b). Chain-of-Thought prompting (Wei et al., 2022b) demonstrated significant reasoning improvements, leading to advanced variants including Self-Consistency (Wang et al., 2022a), Tree-of-Thoughts (Yao et al., 2023), and Plan-and-Solve approaches (Wang et al., 2023a).

Automated prompt optimization evolved as an exciting research direction. The Automatic Prompt Engineer (APE) framework (Zhou et al., 2022) treats instruction generation as natural language synthesis, discovering superior zero-shot prompts compared to human-engineered alternatives. Optimization by PROMpting (OPRO) (Yang et al., 2023) leverages language models as optimizers, achieving up to 50% improvements on Big-Bench Hard tasks through iterative refinement. PromptAgent (Wang et al., 2023b) introduces strategic planning via Monte Carlo Tree Search, systematically exploring expert-level prompt spaces through domain knowledge integration and error feedback mechanisms. Another notable work in this direction are DSPy (Khattab et al., 2023). Greater-Prompt (Zheng et al., 2025) is another prompt optimization-based work that integrates diverse techniques under a customizable API, supporting both text feedback-based and gradient-based optimization across different model scales.

On another dimension, interactive prompt development has been explored through visual interfaces. (Strobelt et al., 2022) introduced PromptIDE for real-time experimentation and performance visualization, while (White et al., 2023) developed comprehensive prompt pattern catalogs analogous to software design patterns. Workflow (Wang et al., 2024) introduces a social prompt engineering paradigm, enabling users to collaboratively create, share, and discover prompts within a community-driven platform. Prompterator (Sučik et al., 2023) provides a human-in-the-loop environment for iteratively refining prompts based on human feedback. These works establish precedents for user-friendly collaborative prompt development.

Another direction focusing on parameter-efficient methods includes HyperPrompt (He et al., 2022), which uses HyperNetworks (Ha et al., 2016; Chauhan et al., 2024) to generate task-conditioned prompts for multi-task learning with only 0.14% additional parameters. Complementary techniques include AutoPrompt (Shin et al., 2020) for gradient-guided prompt search, Prefix Tuning (Li and Liang, 2021), and Prompt Tuning (Lester et al., 2021) for continuous prompt optimization, alongside Prompt-OIRL (Zhang et al., 2023) applying reinforcement learning principles to query-dependent prompt generation.

## 2.2 Instruction Tuning Datasets and Collaborative Infrastructure

Large-scale datasets have been fundamental to advancing instruction tuning research. P3 (Public Pool of Prompts) (Bach et al., 2022) established template-based approaches with 2,000+ prompts across 270+ English datasets, directly influencing PromptLab’s design. The multilingual extension xP3 (Muennighoff et al., 2022) spans 46 languages and 16 NLP tasks, enabling cross-lingual models like BLOOMZ and mT0.

Super-NaturalInstructions (Wang et al., 2022c) provides 1,616 diverse tasks with expert-written instructions across 76 task types, establishing evaluation frameworks for cross-task generalization. BIG-Bench (Srivastava et al., 2022) represents collaborative benchmark development across 450+ authors and 132 institutions, providing architectural insights for community-driven platforms while focusing on tasks beyond current model capabilities.

Instruction tuning foundations began with InstructGPT (Ouyang et al., 2022) and evolved through the Flan family (Wei et al., 2022a; Chung et al., 2024) and T0 (Sanh et al., 2021), demonstrating unified multitask approaches. Synthetic data generation through Self-Instruct (Wang et al., 2022a) and cost-effective approaches like Alpaca (Taori et al., 2023) complement human-curated datasets, while quality-focused methods (Zhou et al., 2024) emphasize careful curation over quantity.

## 2.3 The Landscape of Prompt Engineering Tools

The tooling ecosystem of prompt engineering evolves rapidly to meet diverse needs from research and development domains. This remains a very active area of research and development, with new tools and frameworks appearing regularly as the field matures (Wei et al., 2022b; Brown, 2020).

Current tools can be broadly categorized into two main paradigms: web-based platforms and command-line interface (CLI) tools. Web-based platforms are particularly popular among downstream LLM application developers, those building web and mobile applications that interact with LLMs through APIs and primarily work with direct API integrations. These include playground environments backed by LLMs providers like **OpenAI**

**Playground**<sup>2</sup> and **Anthropic Console**<sup>3</sup>, as well as collaborative platforms designed for prompt development and testing. Table 1 presents and compares the features of PromptLab compared to other web-based platforms and tools.

CLI-based tools, while also popular among developers, tend to attract researchers and those who translate research into usable toolkits for broader adoption. This category includes frameworks like **OpenPrompt** (Ding et al., 2021), **DSPy** (Khattab et al., 2023), **GreaterPrompt** (Zheng et al., 2025), **LangChain**<sup>4</sup>, **LlamaFactory** (Zheng et al., 2024), **LlamaIndex**<sup>5</sup>, **Mirascope**<sup>6</sup>, and **promptwright**<sup>7</sup>. These tools often emphasize programmatic approaches when interacting with LLMs, including prompt engineering, treating it as a software engineering discipline with modules, signatures, and systematic optimization approaches.

PromptLab distinguishes itself within this landscape by focusing on collaboration-oriented research workflows. Unlike existing tools that primarily target application developers and individual prompters, PromptLab is designed specifically for the research community, emphasizing collaborative prompt development, systematic evaluation, and reproducible research practices. This positioning addresses a gap in the current ecosystem where collaboration and research-oriented features are often secondary considerations.

## 3 Platform Architecture and Design

PromptLab builds upon the foundational work of PromptSource (Bach et al., 2022), a pioneering work in template-based prompting research. However, PromptSource requires initial setup and local deployment. Furthermore, the collaboration setup via GitHub pull requests creates barriers for domain experts lacking a computing background. PromptLab addresses these limitations through a comprehensive web-based architecture that eliminates technical overhead while introducing seamless collaborative features. At a high level, PromptLab organizes work around *projects*, each associated with one or more HuggingFace datasets. Within a project, users operate under one of three

<sup>2</sup><https://platform.openai.com/playground/prompts>

<sup>3</sup><https://console.anthropic.com/>

<sup>4</sup><https://www.langchain.com>

<sup>5</sup>[https://github.com/run-llama/llama\\_index](https://github.com/run-llama/llama_index)

<sup>6</sup><https://github.com/Mirascope/mirascope>

<sup>7</sup><https://github.com/StacklokLabs/promptwright>

Table 1: Comparison of prompt engineering tools across key dimensions relevant to collaborative research environments. Rows with "Limited Teams" collaboration are for the freemium tier.

Platform	Audience	Collaboration	HF Datasets	Pricing
PromptSource	Research	GitHub PRs	Native	Free
Agenta	Developers	Limited teams	No	Freemium
PromptHub	Developers	Teams	No	Freemium
LLMs Playground	Individuals	N/A	No	Pay-per-use
PromptLayer	Developers	Limited Teams	No	Freemium
ChainForge	Research	Link Sharing	No	Free
Langfuse	Developers	Limited Teams	No	Freemium
<b>PromptLab</b>	<b>R&amp;D</b>	<b>Teams</b>	<b>Native</b>	<b>Free</b>

roles: *prompters*, *reviewers*, and *administrators*, and prompts progress through a structured lifecycle: draft creation, peer review, revision, and final approval. Figure 1 illustrates this end-to-end pipeline.

### 3.1 Core System Architecture

The platform migrated from Streamlit (as in promptsource) to employ Django as a Python backend framework. Database was setup with PostgreSQL for persistence. Redis was utilized for caching HuggingFace communications, saving time during dataset-intensive calls. Celery enables background task processing for computationally intensive operations, including dataset synchronization, and batch HuggingFace dataset refresh operations. The system exposes RESTful APIs with project-specific authentication tokens, enabling programmatic access from any programming environment. Docker-based deployment and CI/CD pipelines support both public research collaborations and private institutional use cases. PromptLab, instead of prompt review processes via GitHub pull requests, implemented role-based project management with access controls for prompters, reviewers, and administrators. Additionally, administrators can optionally set a minimum prompting workload where the datasets will be randomly distributed among the prompters.

### 3.2 Enhanced Templating and User Experience

PromptLab follows PromptSource practices when designing prompts, maintaining metadata fields like the name, and answer in classification tasks. For the prompt template, PromptLab extends

Jinja2<sup>8</sup> templating with an optional alternative intelligent field insertion through double backslash triggers that display a dropdown list of available dataset fields. Selected fields render as visual tag components, improving the user experience while maintaining the Jinja2 `{{dataset_key}}` syntax under the hood. PromptLab adds a tag field that enables prompt organization and filtering by specific properties like their style (Chain-of-Thought, role-playing, formal), and task type (reasoning, classification).

The integrated template testing view provides real-time visual feedback through color-coded validation indicators. Figure 4 demonstrates this validation across different tasks, showing how the platform identifies well-formed templates versus those with syntax errors or logical inconsistencies, following the original prompting guidelines from PromptSource, and reducing review iteration cycles required for prompt refinement.

### 3.3 AI-Assisted Generation and LLMs Real-Time Evaluation

As primarily inspired by (Wang et al., 2022b; Taori et al., 2023), PromptLab allows prompters to generate AI-assisted prompts using ChatGPT as a strong pretrained language model. Prompters can generate diverse prompt variations from seed templates and carefully check the generated prompts before submitting for review. Real-time evaluation with OpenRouter models enables immediate testing across multiple language models without leaving the platform interface.

<sup>8</sup><https://jinja.palletsprojects.com/en/stable/>

### 3.4 Collaborative Workflow: Review Lifecycle and Communication

A central contribution of PromptLab is its structured multi-stage collaborative workflow, replacing ad hoc coordination systems. Once a prompter submits a draft, it enters a review queue where reviewers can *approve* it, *return it for modification* with inline comments, or *reject* it. Each revision is tracked as a new version, preserving the full development history of a prompt. Administrators retain override capabilities at any stage and can monitor workload distribution across prompters. Figure 5 illustrates the complete state transitions from prompt creation to final approval, ensuring reviewer approvals for considered prompts.

RESTful APIs enable programmatic access with project-specific authentication, supporting diverse programming environments beyond Python-specific bindings. Docker-based deployment options and CI/CD pipelines accommodate both public research collaborations and private institutional projects.

## 4 Platform Evaluation

We conducted two complementary evaluations to evaluate PromptLab for practical considerations: prompt quality assessment comparing PromptLab-generated and manually-created prompts, and human usability analysis.

### 4.1 Prompt Quality Evaluation

To evaluate whether PromptLab facilitates higher-quality prompt development, we conducted a controlled comparison study where we compared prompts created by the platform with those created without using the tool.

#### 4.1.1 Experiment Setup

We selected five benchmark datasets spanning classification (IMDB (Maas et al., 2011)), Toxic-Chat (Lin et al., 2023), MMLU (Hendrycks et al., 2020)) and generation tasks (IWSLT2017 (Cettolo et al., 2017), XSUM (Narayan et al., 2018)). Six researchers<sup>9</sup> with a strong background in prompt engineering were divided into two groups: the **Manual Group** developed prompts without PromptLab assistance, while the **PromptLab Group** used PromptLab. Each researcher created three prompts per dataset, yielding a total of 90 prompts (45

<sup>9</sup>3 Masters, 1 PhD student, 1 Post-Doc, and 1 Professor, all in computer science-related domains

per group) employing diverse strategies including direct instruction, chain-of-thought (Wei et al., 2022b), and role-playing.

We evaluated all prompts across 13 LLMs spanning small ( $\leq 1$ B), medium (1–10B), and large ( $> 10$ B) parameter scales. Each prompt was tested on a stratified subset of 1,000 samples from the test split of each dataset, where samples were randomly shuffled with a random seed for reproducibility. To account for evaluation variance, we conducted 10 independent runs per model using deterministic sampling (temperature=0.0) and averaged the results. Models were evaluated locally using vLLM (Kwon et al., 2023). For classification tasks, we explicitly instructed models to generate only the class label without additional text; violations were penalized. Similarly, for generation tasks, models were instructed to output only translations or summaries. Metrics followed task conventions: standard classification metrics (accuracy, precision, recall, and macro F1-score, F1-score is reported in the figures) for classification tasks, BLEU (Papineni et al., 2002) for translation, and ROUGE-L (Lin, 2004) for summarization. The evaluation code is available in GitHub<sup>10</sup>.

#### 4.1.2 Results and Analysis

Results are reported with normalized prompt performance. We normalized scores by setting the best-performing prompt to 100, with all other prompts scaled proportionally:  $(x/y) \times 100$ , where  $x$  is the prompt’s score and  $y$  is the best score.

Figure 2 presents normalized prompt performance across small models ( $\leq 1$ B parameters). Medium and large parameter figures are provided in the appendices (Figures 7 and 8). PromptLab prompts (green), in the small LLMs, consistently achieve higher performance compared to manual prompts (red) across all five datasets and models. Manual prompts also exhibit higher variance. At medium and large scales, this difference diminishes as models become more capable and less sensitive to prompt formulation.

Figure 3 aggregates results across all model size categories, showing per-model mean performance and variance for each dataset. PromptLab prompts demonstrate a performance advantage across small and medium models. The consistency of this advantage across model scales suggests that PromptLab’s features, particularly real-time validation and

<sup>10</sup><https://github.com/KFUPM-JRCAI/promptlab-evaluation>

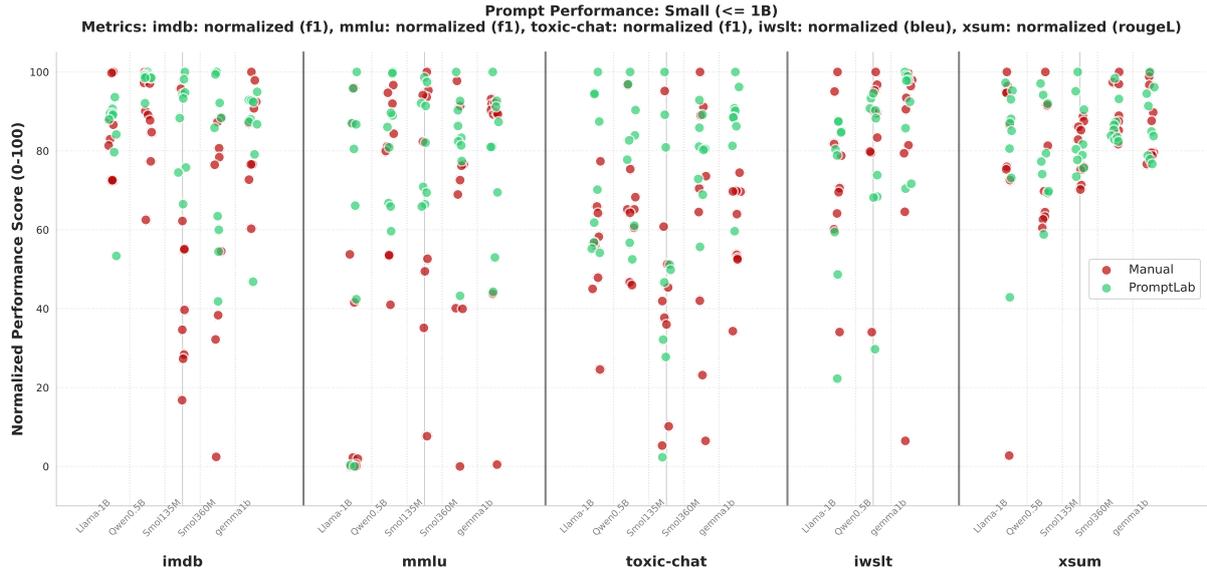


Figure 2: Prompt performance for small models (<1B parameters) across five datasets. Each point represents a single prompt’s normalized performance. PromptLab prompts (green) consistently outperform manually-created prompts (red).

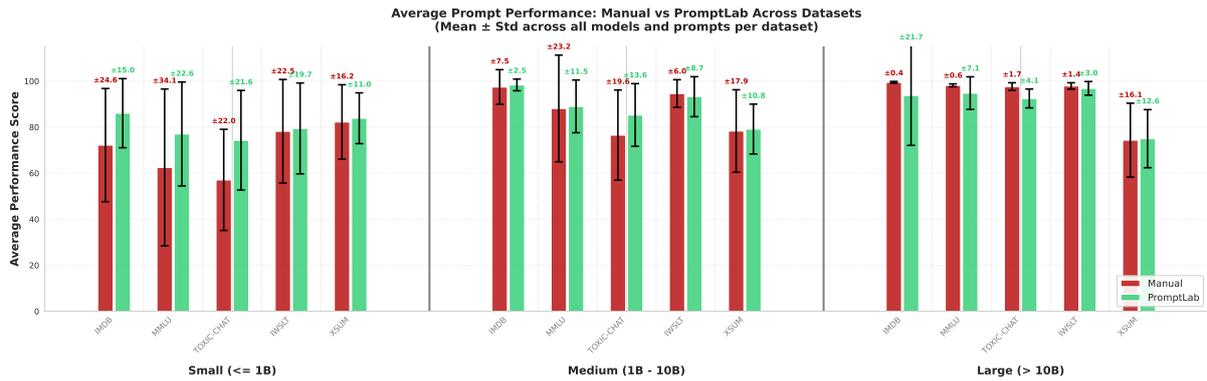


Figure 3: Per-model mean performance comparing manual versus PromptLab prompts across model size categories. Each bar represents one model’s (averaged) performance across all the prompts in that group.

template feedback, enable researchers to identify and iterate on more effective prompting.

## 4.2 Usability Evaluation

To validate PromptLab’s usability effectiveness, we conducted a user evaluation study<sup>11</sup> with 14 participants representing diverse backgrounds: 4 undergraduates, 8 graduate students/engineers, and 2 senior researchers. Participants ranged from beginners to advanced practitioners in NLP research, with varying prompt engineering experience levels. The evaluation assessed platform usability across key dimensions using 5-point Likert scales. Results demonstrate strong user satisfaction with mean scores of 4.3/5 for interface intuitiveness,

<sup>11</sup>We used the following evaluation form <https://forms.gle/nG9jXHRStUkgvLqE9>

4.4/5 for AI-assisted validation, 4.5/5 for template feedback clarity, 4.2/5 for navigation ease, and 4.6/5 for collaborative workflow support. Participants rated PromptLab favorably compared to existing methods (4.4/5), with high likelihood to use (4/5) and recommend (4.2/5) the platform. Appendix G provides more in-depth details and insights on the evaluation results.

## 5 Case Study: Arabic NLP Instruction Tuning Research

To demonstrate PromptLab’s effectiveness in real-world research scenarios, we present a case study from an Arabic instruction tuning project that represents one of the most substantial collaborative prompt engineering efforts for Arabic NLP to date.

The project involved seven researchers and three reviewers who developed over 350 manually crafted prompt templates spanning 20+ Arabic NLP tasks across 80+ datasets (refer to Appendix C for the full datasets listing), totaling more than 6.3 million samples. Figure 6 in Appendix D presents the prompts distributions over the datasets. Tasks covered a broad spectrum including dialect identification, sentiment analysis, sarcasm detection, natural language inference, machine translation, and summarization, among others. Researchers employed both manual and AI-assisted prompt creation workflows. The platform’s assignment and review interfaces allowed members to track individual progress and maintain a full revision history for each prompt submission. As a side result of this research, a framework for scoring prompts across similarity, performance, efficiency, and consistency dimensions was developed and presented in (Refai et al., 2025).

## 6 Conclusion and Future Work

PromptLab addresses the gaps in research-based collaborative prompt engineering through a comprehensive platform that supports the complete prompt life cycle from creation to validation and publication. The platform’s design balances accessibility with technical sophistication, combining intuitive visual interfaces with powerful programmatic capabilities to enable broader participation while preserving expressiveness. The Arabic instruction tuning case study demonstrates the platform’s effectiveness, with over 350 prompts across 80+ datasets successfully developed and validated through systematic collaborative workflows involving multiple researchers and reviewers.

The integration of AI-assisted prompts development with human oversight establishes a productive collaborative paradigm that enhances researcher productivity without compromising quality control. Future work will focus on integrating automated prompt optimization techniques, model evaluation, model fine-tuning, and developing comprehensive analytics for collaborative pattern analysis. The platform’s open-source availability and deployment flexibility position it as a foundation for continued progress in prompt engineering, contributing to the democratization of advanced NLP research across diverse individual and institutional environments.

## Acknowledgments

We gratefully acknowledge the support of the SDAIA-KFUPM Joint Research Center for Artificial Intelligence (JRCAI) for providing the resources and infrastructure that made this work possible through grant number JRC-AI-UCG-07. We extend our thanks to the 14 participants who contributed their time and expertise to the platform usability evaluation study, as well as the six researchers who participated in the controlled prompt quality comparison experiment. We also thank the team members involved in the Arabic instruction tuning case study for their dedicated efforts in developing and validating over 350 prompt templates. Finally, we appreciate the open-source community and the developers of the tools and libraries upon which PromptLab is built.

## References

- Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, and 1 others. 2022. Promptsource: An integrated development environment and repository for natural language prompts. *arXiv preprint arXiv:2202.01279*.
- Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. [Overview of the IWSLT 2017 evaluation campaign](#). In *Proceedings of the 14th International Conference on Spoken Language Translation*, pages 2–14, Tokyo, Japan. International Workshop on Spoken Language Translation.
- Vinod Kumar Chauhan, Jiandong Zhou, Ping Lu, Soheila Molaei, and David A Clifton. 2024. A brief review of hypernetworks in deep learning. *Artificial Intelligence Review*, 57(9):250.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, and 1 others. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.
- Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Hai-Tao Zheng, and Maosong Sun. 2021. Openprompt: An open-source framework for prompt-learning. *arXiv preprint arXiv:2111.01998*.
- David Ha, Andrew Dai, and Quoc V Le. 2016. Hypernetworks. *arXiv preprint arXiv:1609.09106*.

- Yun He, Huaixiu Steven Zheng, Yi Tay, Jai Gupta, Yu Du, Vamsi Aribandi, Zhe Zhao, Yaguang Li, Zhaoji Chen, Donald Metzler, and 1 others. 2022. Hyperprompt: Prompt-based task-conditioning of transformers. *arXiv preprint arXiv:2203.00759*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, and 1 others. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*, pages 611–626.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang. 2023. Toxicchat: Unveiling hidden challenges of toxicity detection in real-world user-ai conversation. *Preprint*, arXiv:2310.17389.
- 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.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196*.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, and 1 others. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *ArXiv*, abs/1808.08745.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 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.
- Dania Refai, Maged S Al-Shaibani, and Irfan Ahmad. 2025. Is this the best prompt? scoring prompts for arabic nlp across llms. *IEEE Access*.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024a. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927*.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024b. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927*.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, and 1 others. 2021. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.
- Sander Schulhoff, Michael Ilie, Neel Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Fu, Barnabás Póczos, and 1 others. 2024a. The prompt report: A systematic survey of prompt engineering techniques. *arXiv preprint arXiv:2406.06608*.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, and 1 others. 2024b. The prompt report: A systematic survey of prompting techniques. *arXiv preprint arXiv:2406.06608*.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, and 1 others. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.

- Hendrik Strobelt, Albert Webson, Victor Sanh, Benjamin Hoover, Johanna Beyer, Hanspeter Pfister, and Alexander M Rush. 2022. Interactive and visual prompt engineering for ad-hoc task adaptation with large language models. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):1146–1156.
- Samuel Sučik, Daniel Skala, Andrej Švec, Peter Hraška, and Marek Šuppa. 2023. **Prompterator: Iterate efficiently towards more effective prompts**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 471–478, Singapore. Association for Computational Linguistics.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023a. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. *arXiv preprint arXiv:2305.04091*.
- Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P Xing, and Zhiting Hu. 2023b. Promptagent: Strategic planning with language models enables expert-level prompt optimization. *arXiv preprint arXiv:2310.16427*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022a. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022b. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, and 1 others. 2022c. 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.
- Zijie Wang, Aishwarya Chakravarthy, David Munechika, and Duen Horng Chau. 2024. **Workflow: Social prompt engineering for large language models**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 42–50, Bangkok, Thailand. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022b. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2023. Large language models as optimizers. *arXiv preprint arXiv:2309.03409*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Hao Sun Zhang, Alihan Hüyük, and Mihaela van der Schaar. 2023. Query-dependent prompt evaluation and optimization with offline inverse rl. *arXiv preprint arXiv:2309.06553*.
- Wenliang Zheng, Sarkar Snigdha Sarathi Das, Yusen Zhang, and Rui Zhang. 2025. **GreaterPrompt: A unified, customizable, and high-performing open-source toolkit for prompt optimization**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 405–415, Vienna, Austria. Association for Computational Linguistics.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. **Llamafactory: Unified efficient fine-tuning of 100+ language models**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, and 1 others. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*.

## A PromptLab templating pipeline

Figure 4 and Figure 5 show the general templating pipeline utilized by PromptLab.

<p>Between the following two sentences, which is more realistic?            1. {{first_sentence}}            2. {{second_sentence}}            Select the correct option: {{answer_choices   join(' or ')}}.                           {{answer_choices[label]}}</p>	Valid
<p>Among the following pairs of sentences, one is more likely or makes more sense:            1. {{first_sentence}}            2. {{second_sentence}}            Evaluate and determine which one is more reasonable: {{label   join(' or ')}}                           {{label[label]}} ✘</p>	Invalid Does not render!
<p>You have the following Arabic sentence: {{arabic}}.            Based on this sentence, select the dialect from these options:            {{answer_choices   join(', ')}}, and provide the appropriate dialect.                           {{answer_choices[label]}}</p>	Valid
<p>Review the following text: '{{Text}}'. Determine the dialect from the following options:{{answer_choices   join(', ')}}     ✘</p>	Invalid No output!

Figure 4: Examples for **valid** and **invalid** templates for common sense validation and dialect identification tasks.

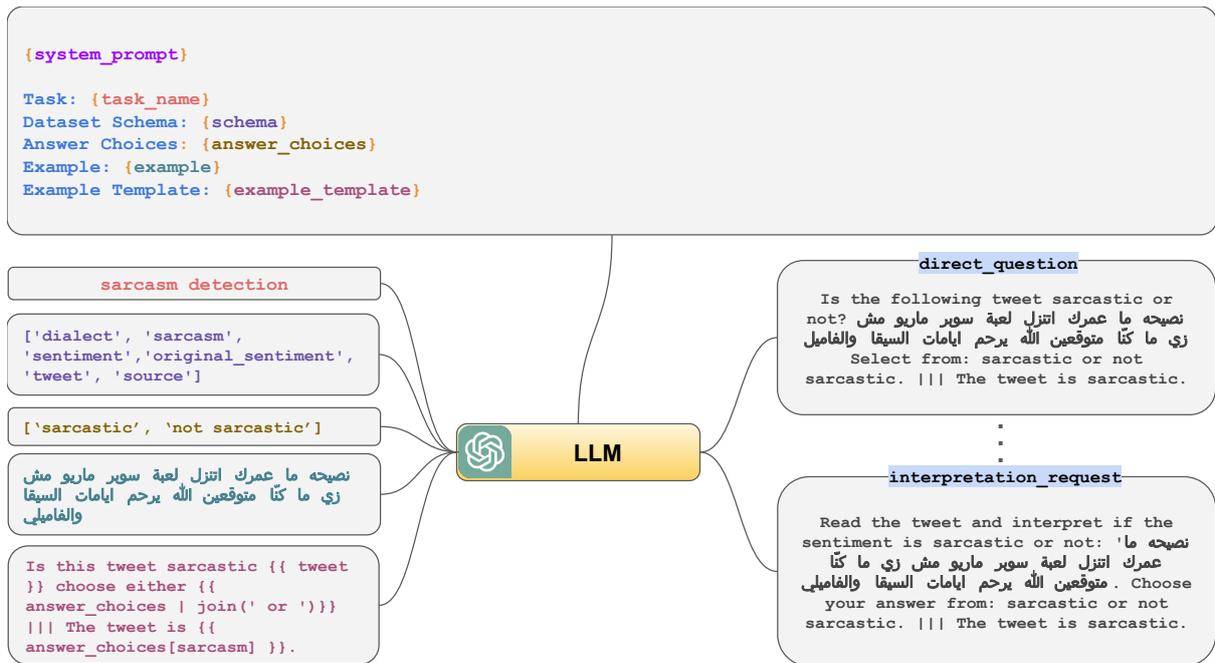


Figure 5: LLM template creation pipeline.

## B AI generated prompts

We prompt the LLM to generate templates for a given dataset. The most important part is the system prompt. We use the following system prompt to create the templates for a given dataset. Texts shown in red are variables. For example, we can replace English with Arabic to create templates in Arabic instead of English.

**System Prompt:** You are a prompt template creator. Given a task, dataset schema, and answer\_choices. You should create a template using Jinja that can be applied to an example in the dataset. The prompt and completion are separated by |||. You should create 5 different templates in English language as a json with a key that represents each template content. Please choose creative templates with enough variation. The order of the jinja variables can be changed. Do Not use a general name of the template like "template", USE more representative name. Do NOT print any other text except the json. Do NOT use any integer features. If there are answer\_choices, use as is do NOT change. If there are answer\_choices use the variable as is, do NOT introduce any new answer choices. All the jinja variables must be from the schema. Do NOT introduce new variable names. If the answer\_choices in the example template exists use it in the completion without any changes. This is an important test. Please respect all the mentioned points.

## C Prompts from Arabic NLP Datasets

Table 2 presents a listing of 80+ Arabic NLP datasets currently available on the Hugging Face platform that we utilized for our case study. This collection, totaling over 6.3 million samples, spans various tasks and linguistic phenomena. The datasets exhibit significant variation in size, from small specialized collections like PAAD (206 samples) to extensive corpora such as APCD (1,831,770 samples) and APCD2 (1,657,003 samples). As shown in Table 2, the datasets cover fundamental text processing tasks like diacritization (Arabic Text Diacritization, Shakkelha) as well as complex semantic tasks such as

commonsense validation and natural language inference (ArEntail). Many datasets focus on specific dialect variations or regional Arabic variants, such as the Shami dataset (66,251 samples) and the Arabic Dialects Dataset (9,992 samples). Others target particular applications like sentiment analysis (HARD, LABR) and text classification (SANAD, Ultimate Arabic News Dataset). While this collection represents a substantial resource for Arabic NLP research, the distribution of samples across different tasks and the varying quality of annotations present both opportunities and challenges for instruction-tuning approaches. Each dataset in Table 2 is accompanied by its Hugging Face URL, facilitating easy access and integration into research workflows. Each of these datasets has at least one prompt associated with it.

Table 2: A comprehensive list of Arabic NLP datasets curated for prompts.

<b>Dataset</b>	<b>Size</b>	<b>HuggingFace ID</b>
Commonsense validation	11,000	arbml/Commonsense_Validation
Arabic Text Diacritization	55,000	arbml/arabic_text_diacritization
Shakkelha	533,384	arbml/shakkelha
AraBench	42,113	arbml/AraBench_dev
Arabic Dialects Dataset	9,992	arbml/Arabic_Dialects_Dataset
araData	12,231	arbml/araData
Shami	66,251	arbml/Shami
Twt15DA_Lists	3,037	arbml/Twt15DA_Lists
Emotional-Tone	10,065	emotone-ar-cicling2017/emotone_ar
SemEval-2018 Task 1	4,381	SemEvalWorkshop/sem_eval_2018_task_1
ArMATH	6,000	arbml/ArMATH
APCD	1,831,770	arbml/APCD_remap
APCD2	1,657,003	arbml/APCDv2
ashaar	212,499	arbml/Ashaar_dataset
ArabicMMLU	14,575	arbml/ArabicMMLU
cidar-mcq	100	arbml/CIDAR-MCQ-100
belebele	900	facebook/belebele
QA4MRE	160	arbml/qa4mre
EXAMS	562	OALL/Arabic_EXAMS
AQMAR	4,188	arbml/AQMAR_batched
CANERCorpus	16,139	arbml/caner_batched
Cross-lingual NER	1,208	arbml/Zero_Shot_Cross_Lingual_NER_ar_batched
Disease NER	3,906	arbml/Disease_NER_batched
Named Entities Lexicon	48,753	arbml/Named_Entities_Lexicon
ArEntail	6,000	arbml/ArEntail
Textual Entailment	422	arbml/ArabicTE
Arabic OSACT4	7,837	arbml/OSACT4_hatespeech
MLMA hate speech	3,353	arbml/MLMA_hate_speech_ar
MPOLD	4,000	arbml/MPOLD
OffensEval 2020	9,666	arbml/offenseval_2020
Dangerous Speech Dataset	5,009	arbml/Dangerous_Dataset
Arabic Hate Speech 2022	9,823	arbml/Arabic_Hate_Speech

Continued on next page

Table 2 – continued from previous page

<b>Dataset</b>	<b>Size</b>	<b>HuggingFace ID</b>
Arabic POS Dialect	1,400	arbml/QCRI_arabic_pos_dialect
Arabic senti-lexicon	3,941	arbml/Senti_Lexicon
AQAD	17,911	arbml/AQAD
Arabic RC datasets	1,008	arbml/Arabic_RC_AQA
ARCD	1,395	hsseinmz/arc
tydiqa-goldp	15,726	khalidalt/tydiqa-goldp
MKQA	10,000	apple/mkqa
xquad	1,190	google/xquad
TYDIQA	15,726	asas-ai/tydiqa-ar
ACVA	9,000	arbml/ACVA
ArSarcasm	10,547	iabufarha/ar_sarcasm
Sa`7r	19,804	arbml/SaudiIrony
ArSarcasm-v2	15,548	arbml/ArSarcasm_v2
iSarcasmEval	1,400	arbml/iSarcasmEval_task_A
nsurl	3,715	arbml/nsurl_2019_task8_test
Quran Hadith Datasets	8,144	arbml/Quran_Hadith
HARD	105,698	Elnagara/hard
LABR	14,695	mohamedadaly/labr
OCLAR	3,916	arbml/oclar
BRAD 1.0	510,598	arbml/BRAD
AJGT	1,800	komari6/ajgt_twitter_ar
ArSAS	19,897	arbml/ArSAS
ArSentiment	8,364	hadyelsahar/ar_res_reviews
ASTAD	68,070	arbml/Sentiment_Analysis_Tweets
ASTD	9,694	arbml/ASTD
BBN Blog Posts	1,200	arbml/BBN_Blog_Posts
ElecMorocco2016	10,254	arbml/ElecMorocco
ATT	2,154	arbml/ATT
MSAC	1,829	arbml/MSAC
NileULex	5,953	arbml/NileULex
Sudanese Dialect tweets	2,119	arbml/Sudanese_Dialect_Tweet
Sudanese Telecom tweets	5,346	arbml/Sudanese_Dialect_Tweet_Tele
Syria Tweets	2,000	arbml/Syria_Tweet_Sentiment
TSAC	11,871	arbml/TSAC
AT-ODTSA	3,000	arbml/AT_ODSTA
AraStance	4,063	arbml/arastance
Mawqif	3,502	arbml/Mawqif
ANS CORPUS	3,786	arbml/ANS_stance
ArCovidVac	9,988	arbml/ArCovidVac
AraSum	49,603	arbml/AraSum

Continued on next page

Table 2 – continued from previous page

Dataset	Size	HuggingFace ID
WikiLingua	9,995	esdurmus/wiki_lingua
XLSum	46,897	GEM/xlsum
AGS-Corpus	141,467	FahdSeddik/AGS-Corpus
Goud-sum	158,282	Goud/Goud-sum
ARGEN	1,200	arbml/ARGEN_title_generation
ANTCORPUS	10,161	arbml/antcorpus
Khaleej-2004	5,690	arbml/khaleej_2004
OSAC	5,070	arbml/OSAC_CNN
PAAD	206	arbml/PAAD
SANAD	141,807	arbml/SANAD
Ultimate Arabic News	196,279	arbml/ultimate_arabic_news
Watan-2004	20,291	arbml/watan_2004
<b>Total</b>	<b>6,324,527</b>	

#### D Arabic instructions tuning dataset distribution

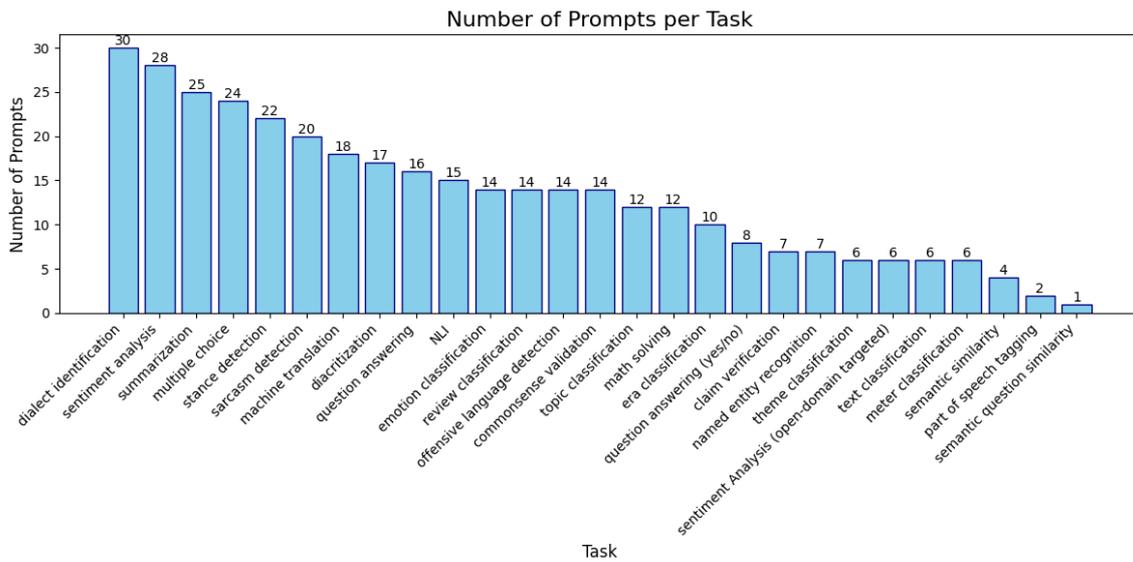


Figure 6: Prompts distribution across different Arabic NLP tasks.

```

1 import requests
2 import json
3
4 # The URL for the API endpoint
5 url = "https://promptlab.up.railway.app/api/prompt/create"
6
7 # The headers for the request
8 headers = {
9     "Content-Type": "application/json"
10 }
11
12 # The data payload
13 data = {
14     "name": "<prompt-name>",
15     "template": "Translate {text} to {language}",
16     "dataset_huggingface_name": "<dataset-name>",
17     "dataset_subset": "", # optional
18     "project_secret_key": "<prompting-project-secret-key>",
19     "created_by": "<created-by>",
20     "tags": ["AI generated", "AI translated"], # optional
21     "text_direction": "ltr", # optional, defaults to ltr
22     "answer_choices": json.dumps([{"value": "choice1"}])
23 }
24
25 # Send the POST request
26 response = requests.post(url, headers=headers, json=data)
27
28 # Check the response
29 if response.status_code == 201:
30     print("Prompt created successfully!")
31     print("Response:", response.json())
32 else:
33     print("Failed to create prompt")
34     print("Status code:", response.status_code)
35     print("Response:", response.text)

```

Listing 1: Python code for creating prompts via PromptLab API

## E REST API communication

To facilitate programmatic access, we implemented a RESTful API to easily integrate with any existing pipelines. We mainly implemented two primary endpoints:

- **Prompt Creation Endpoint:** Enables programmatic creation of prompts. The endpoint implements proper validation and authentication mechanisms to ensure data integrity.
- **Prompt Retrieval Endpoint:** Provides filtered access to the prompt repository related to a prompting project.

These endpoints are useful as they:

- Automate prompt generation.
- Automate prompt sharing and merging on different datasets when appropriate.
- Integrate the platform into existing research pipelines.
- Help in implementing prompt evaluation workflows.

Listings 2,1 provide Python code snippets to interact with these endpoints.

To retrieve existing prompts, the prompt listing endpoint can be utilized to list prompts related to a project as in Figure 2:

```
1 import requests
2
3 url = 'https://promptlab.up.railway.app/api/prompt/list'
4 api_response = requests.get(
5     url=f'{url}?project_secret_key=<project-secret-key>',
6 )
7
8 if api_response.ok:
9     prompts = api_response.json()
10 else:
11     print(api_response.text)
```

Listing 2: Python code for retrieving prompts via PromptLab API

## F Platform Prompt Quality Evaluation

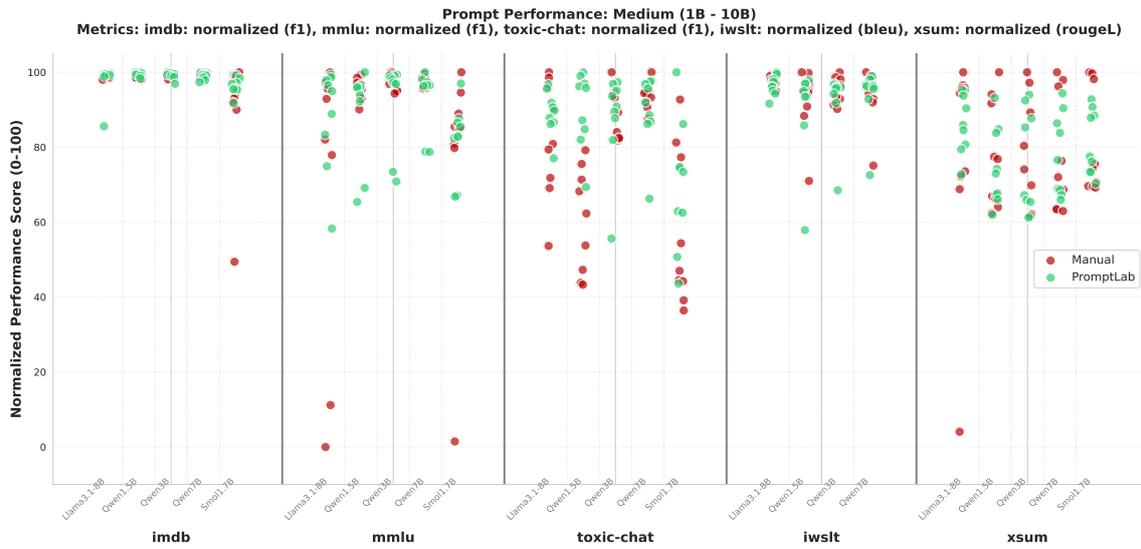


Figure 7: Prompt performance for medium models (10B<1B parameters) across five datasets. Each point represents a single prompt’s normalized performance. PromptLab prompts (green) consistently outperform manually-created prompts (red).

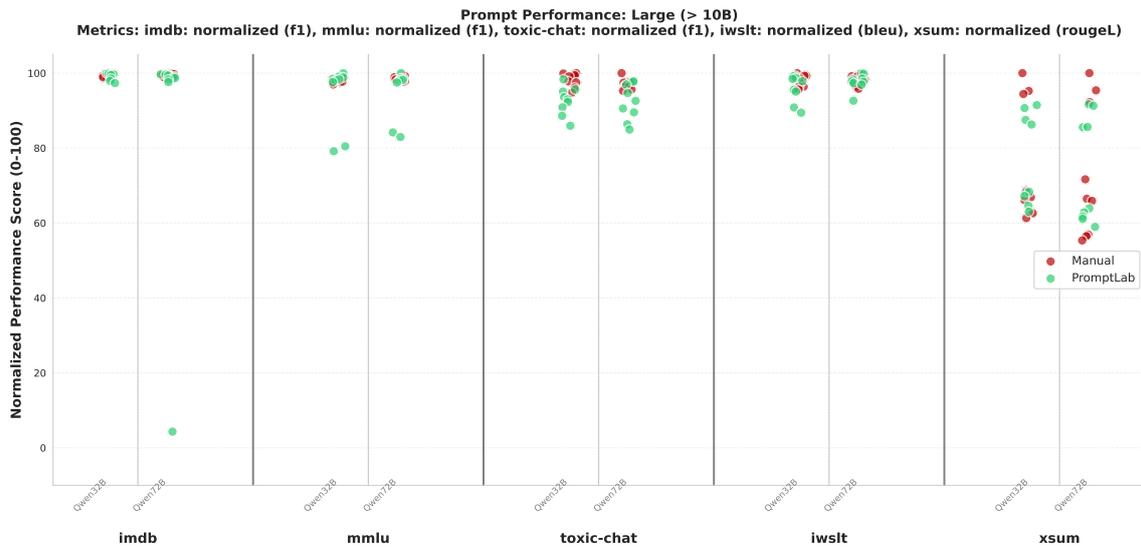


Figure 8: Prompt performance for large models (>10B parameters) across five datasets. Each point represents a single prompt’s normalized performance. PromptLab prompts (green) consistently outperform manually-created prompts (red).

## G Platform Usability Evaluation Results

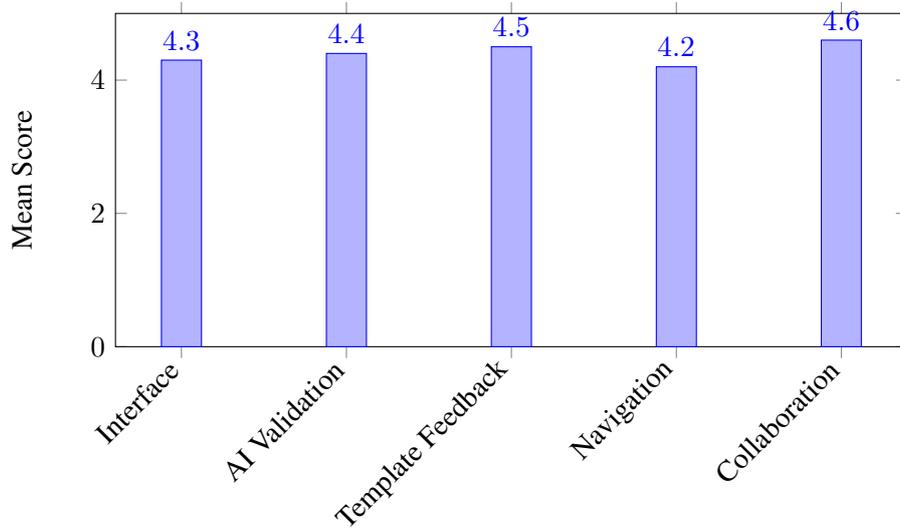


Figure 9: Platform usability scores across key dimensions (5-point Likert scale).

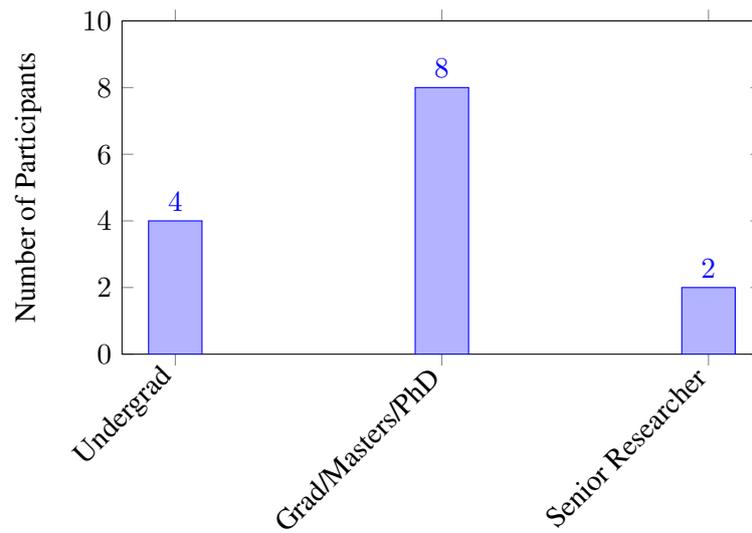


Figure 10: Participant distribution by education level (N=14).

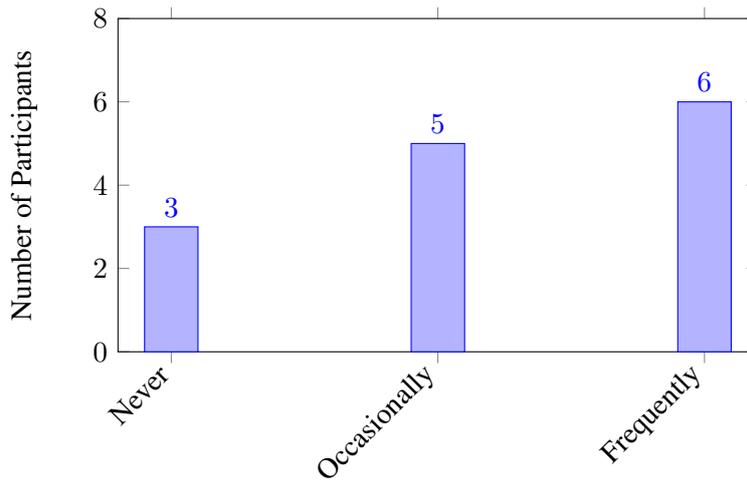


Figure 11: Prior prompt engineering experience distribution (N=14).

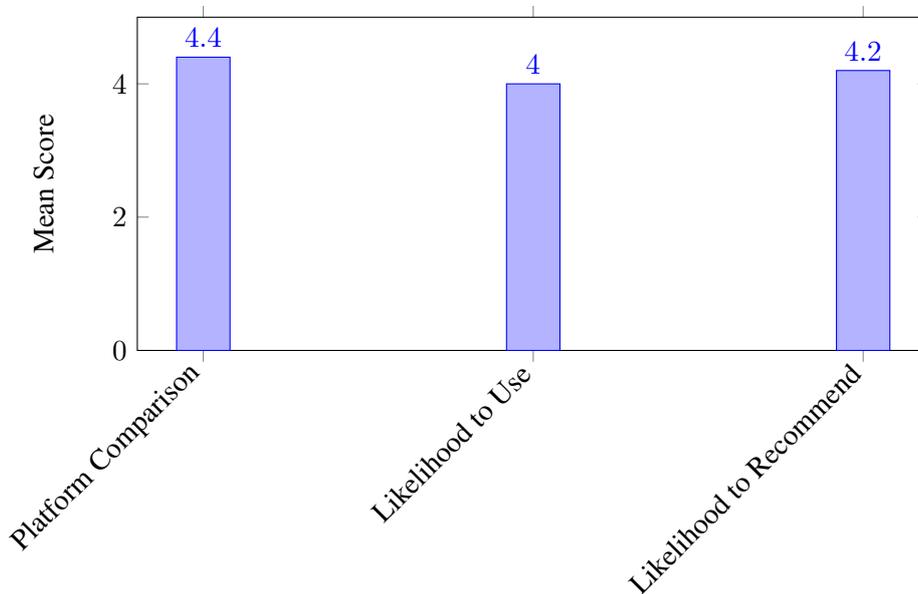


Figure 12: Overall platform satisfaction metrics (5-point Likert scale).

## H PromptLab: Tool Documentation

PromptLab is designed to streamline and enhance the process of creating, refining, and managing prompts for Large Language Models (LLMs). Unlike existing prompt engineering platforms, PromptLab integrates a comprehensive set of functionalities—ranging from dataset exploration and prompt editing to review workflows and AI-assisted prompt generation. By offering a centralized environment where researchers, linguists, and developers can collaboratively develop and maintain prompts, PromptLab aims to foster an environment for preparing high-quality datasets.

Furthermore, PromptLab is structured to support an iterative feedback loop, guiding users through the entire lifecycle of prompt creation—from initial brainstorming to final approval. By providing dedicated interfaces for dataset inspection, prompt variation, structured reviews, and automated transformations (such as translation or AI-generated expansions), the tool facilitates a more dynamic and data-driven approach to prompt engineering. The platform’s detailed record-keeping of revision histories, reviewer decisions, and dataset metadata further ensures reproducibility and accountability in collaborative research settings. In doing so, PromptLab empowers the NLP community with a scalable, user-friendly resource that promotes best practices, accelerates model training preparation, and ultimately contributes to the

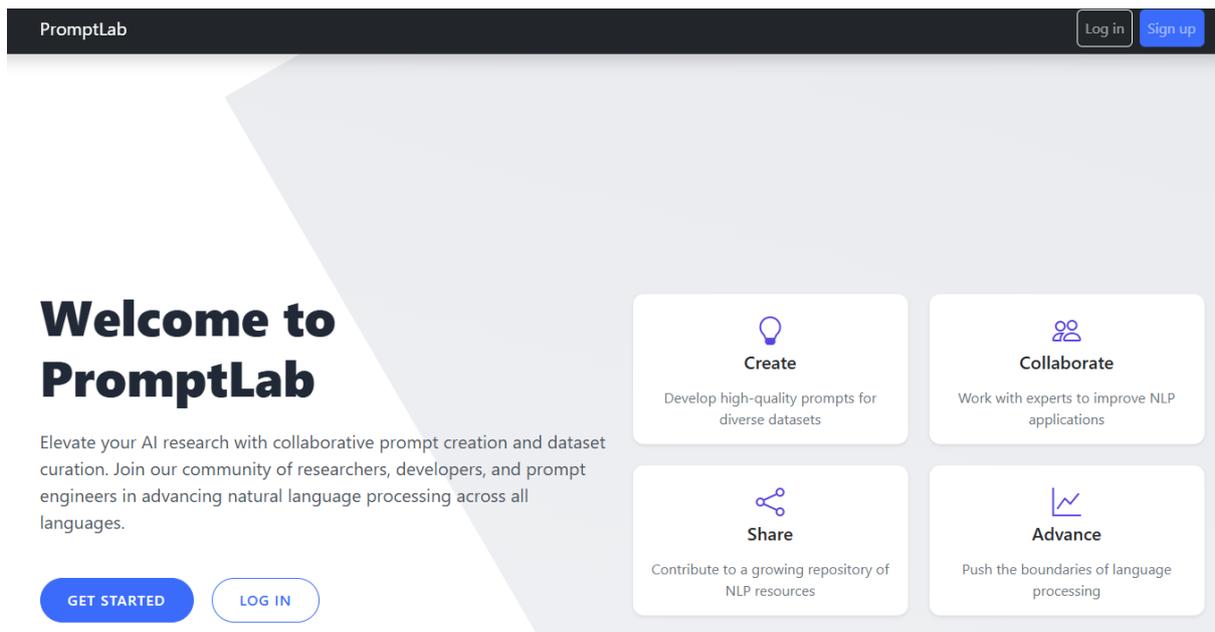


Figure 13: PromptLab: The main dashboard.

advancement of language processing technologies.

In the next section, we will present the main functionalities of the tool along with sample screens.

## H.1 Main Features of the Tool

Upon logging into the platform, users are presented with a primary dashboard designed to facilitate streamlined navigation and highlight the tool’s core functionalities. The dashboard is shown in Figure 13.

The top portion of the interface displays the platform’s logo and a concise header menu offering essential actions, such as “Log In” or “Sign Up,” depending on the user’s authentication state. The central section of the screen prominently features a welcoming message (“Welcome to PromptLab”) accompanied by a brief explanatory tagline emphasizing the platform’s focus on elevating NLP research through collaborative prompt creation and dataset curation.

Below the heading, four distinct feature panels are arranged horizontally, each represented with an icon and descriptive label:

- **Create:** Encourages the development of high-quality prompts.
- **Collaborate:** Invites users to engage with experts and peers, fostering an environment that promotes shared learning and collective improvements in NLP.
- **Share:** Provides pathways for contributing to a communal repository of NLP resources, thereby expanding the overall dataset and knowledge base accessible to the community.
- **Advance:** Focuses on the progressive enhancement of natural language processing, guiding users toward advanced practices and cutting-edge methodologies.

PromptLab organizes work into prompting projects that group related datasets, prompts, and collaborators under a single umbrella. From the My Projects page, users can browse, search, and sort all accessible projects, each shown as a card summarizing its title, description, number of datasets, and team size, with quick actions to view the workspace, open tasks and datasets, or edit project details as illustrated in Figure 14. Opening a project leads to a project overview screen displaying high-level metadata (name, description, owner, creation details, and minimum-prompt requirements), project-level API information including a secret key and example usage code, and a Team Members panel listing collaborators and their roles with search and pagination controls as shown in Figure 15. Together, these views provide the

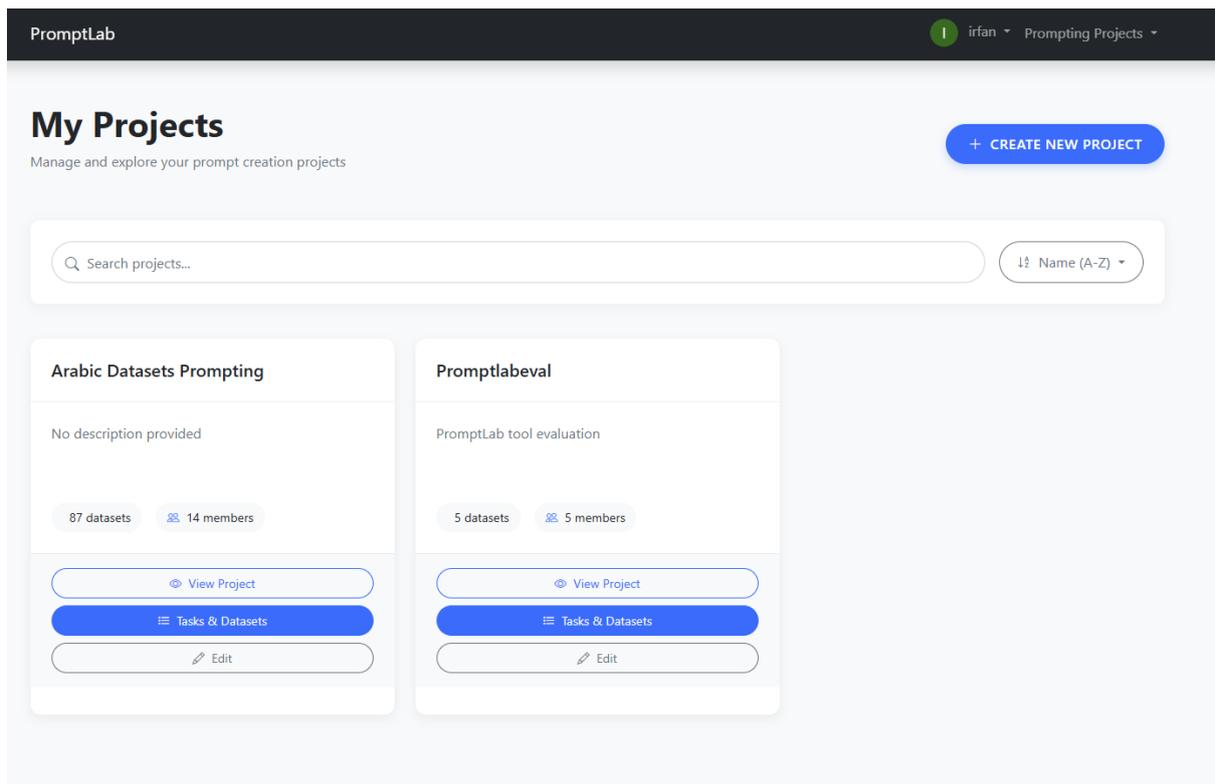


Figure 14: PromptLab: List of projects.

main project-management layer in PromptLab, from which users configure access, connect automated workflows, and then proceed to detailed task and dataset management.

A prominent “Get Started” call-to-action button is positioned beneath the introductory text, directing users to begin creating or exploring content immediately. The overall visual presentation employs a clean, modern design with intuitive iconography and a balanced color palette, enhancing ease of use and navigability. This layout ensures that both newcomers and experienced researchers can efficiently locate key functionalities and engage with the platform’s core features.

Prominent call-to-action buttons direct new users to begin exploring available resources, while returning researchers can effortlessly resume their work. Transitioning from this home environment, the “NLP Tasks” interface (Figure 16) provides a structured overview of diverse language-related tasks, each aligned with relevant datasets and associated prompt counts. The clean tabular layout simplifies the discovery process by allowing users to quickly scan the catalog of tasks, gauge their complexity through prompts and dataset availability, and select actions such as “View Datasets” to delve deeper into a particular dataset.

Building upon the task-level insights, the “PromptLab Prompts List” interface, as shown in Figure 17, offers a detailed registry of prompts associated with chosen datasets for a specific user. Here, users can track prompt submission status (e.g., draft or submitted), review semantic tags that characterize the prompt’s nature (e.g., “CoT” for chain-of-thought reasoning), and make informed decisions about which prompts require refinement and validation. This structured presentation supports a workflow where researchers iterate on prompt quality and scope.

The “My Assigned Datasets” interface, as in Figure 18, personalizes the experience by displaying only those datasets a user is currently assigned to work on. The interface quantifies progress through a simple count (e.g., “2/5” prompts created), helping users maintain momentum and accountability in their prompt development process. This tailored view closes the loop from broad task exploration to individual dataset responsibilities, encouraging sustained engagement and focused contributions to the platform’s collaborative research ecosystem.

As researchers advance deeper into the platform’s resources, they encounter specialized interfaces

PromptLab irfan Prompting Projects

# PromptLabEval

PromptLab tool evaluation Edit Project

Owner: majed.alshaibani Created: Minimum prompts: 5

## API Information

Project Secret Key (for API access)

**DyDck** Copy

Use this key for API access to this project. Keep it secure.

### Example API Usage

```
import requests
import json

# API endpoint for creating prompts
url = "https://promptlab.up.railway.app/api/prompt/create"

# Request data with your project secret key
data = {
    "name": "Example Prompt",
    "template": "Translate {text} to {language}",
    "dataset_huggingface_name": "arbm1/watan_2004",
    "project_secret_key": "DyDck",
    "created_by": "username",
    "tags": ["translation", "api-test"]
}

# Send the request
response = requests.post(url,
    headers={"Content-Type": "application/json"},
    json=data)
```

## Team Members (5)

Search members...

- maged.alshaibani** (Owner)
- irfan** (Reviewer, Prompter)
- zaid1** (Reviewer, Prompter)
- irfan9** (Reviewer, Prompter)
- zaid1** (Reviewer, Prompter)

Navigation: < 1 2 >

## Project Statistics

Figure 15: PromptLab: Project management dashboard.

PromptLab irfan Prompting Projects

Home > NLP Tasks

# NLP Tasks

Search tasks... View All Datasets

#	Task Name	Datasets	Prompts	Actions
1	<a href="#">claim verification</a>	1	10	<span>View Datasets</span>
2	<a href="#">commonsense validation</a>	1	14	<span>View Datasets</span>
3	<a href="#">diacritization</a>	2	17	<span>View Datasets</span>
4	<a href="#">dialect identification</a>	5	31	<span>View Datasets</span>
5	<a href="#">emotion classification</a>	2	14	<span>View Datasets</span>
6	<a href="#">era classification</a>	2	10	<span>View Datasets</span>
7	<a href="#">machine translation</a>	4	27	<span>View Datasets</span>
8	<a href="#">math solving</a>	1	12	<span>View Datasets</span>
9	<a href="#">meter classification</a>	3	6	<span>View Datasets</span>
10	<a href="#">multiple choice</a>	6	33	<span>View Datasets</span>

Navigation: 1 2 3 >

Figure 16: PromptLab: List NLP Tasks.

PromptLab irfan Prompting Projects

## PromptLab Prompts List Sync with Hugging Face

#	Name	Dataset	Subset	Status	Tags
1	<a href="#">A Simple Test Prompt</a>	AraBench_dev		DRAFT	No tags
2	<a href="#">Elicit all options</a>	sem_eval_2018_task_1		DRAFT	Example Prompt
3	<a href="#">To the point</a>	imdb		SUBMITTED	No tags
4	<a href="#">Focusing on the main message</a>	imdb		SUBMITTED	Focusing on the main message overall sentiment
5	<a href="#">Summarize and judge</a>	imdb		SUBMITTED	Summarize judge
6	<a href="#">Straight</a>	iwslt2017		SUBMITTED	No tags
7	<a href="#">literal translation</a>	iwslt2017		SUBMITTED	No tags
8	<a href="#">Arabic translation in MSA using meaning, context, and tone</a>	opus_infopankki		SUBMITTED	No tags
9	<a href="#">basic translation as Arabic expert</a>	opus_infopankki		SUBMITTED	No tags
10	<a href="#">DialectSarcasmInquiry</a>	iSarcasmEval_task_A		SUBMITTED	AI generated

1 2 3 4 5 »

Figure 17: PromptLab: An interface to monitor prompts created by a user.

## My Assigned Datasets

Task	Dataset	Prompts Created	Actions
NLI	<a href="#">ArabicTE</a>	2 / 5	<a href="#">My Prompts on this dataset</a>
math solving	<a href="#">ArMATH</a>	1 / 5	<a href="#">My Prompts on this dataset</a>
summarization	<a href="#">wiki_lingua</a>	3 / 5	<a href="#">My Prompts on this dataset</a>
diacritization	<a href="#">arabic_text_diacritization</a>	2 / 5	<a href="#">My Prompts on this dataset</a>
informativeness	<a href="#">ArCovidVac</a>	0 / 5	<a href="#">My Prompts on this dataset</a>
multiple choice	<a href="#">Arabic_EXAMS</a>	1 / 5	<a href="#">My Prompts on this dataset</a>
stance detection	<a href="#">Mawqif</a>	1 / 5	<a href="#">My Prompts on this dataset</a>
sarcasm detection	<a href="#">SaudiIrony</a>	0 / 5	<a href="#">My Prompts on this dataset</a>
claim verification	<a href="#">ANS_stance</a>	5 / 5	<a href="#">My Prompts on this dataset</a>
era classification	<a href="#">APCD_remap</a>	1 / 5	<a href="#">My Prompts on this dataset</a>

1 2 3 »

Figure 18: PromptLab: Explore the progress of a user for the datasets assigned to him.

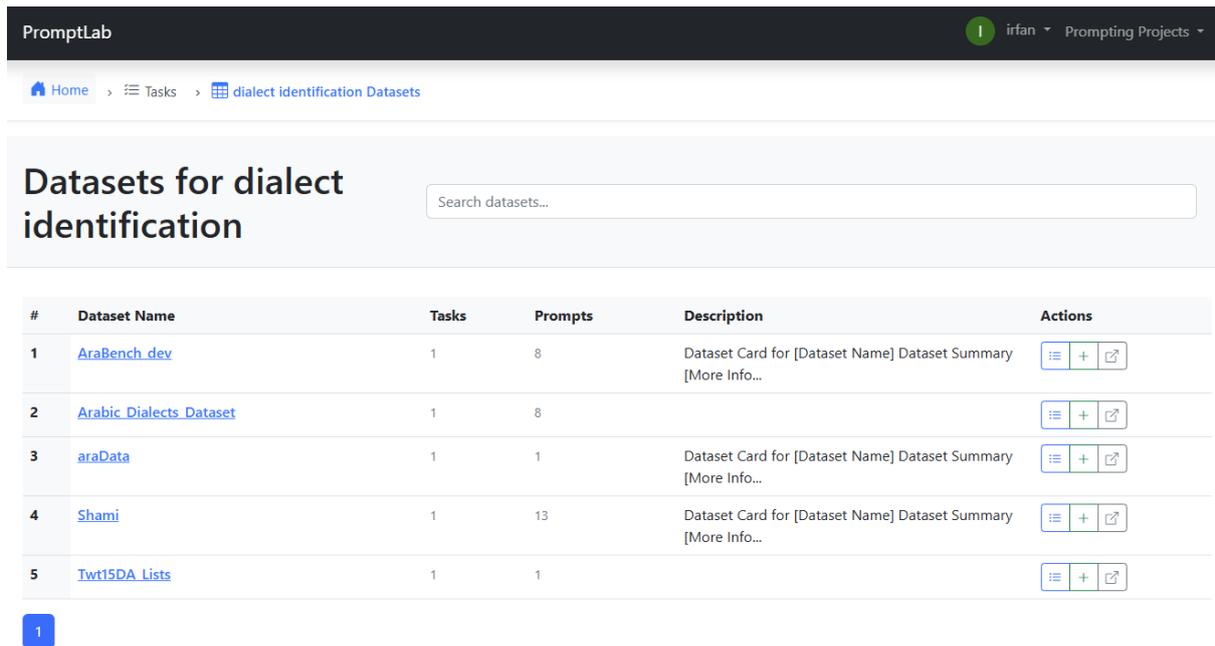


Figure 19: PromptLab: Explore the datasets for a specific task.

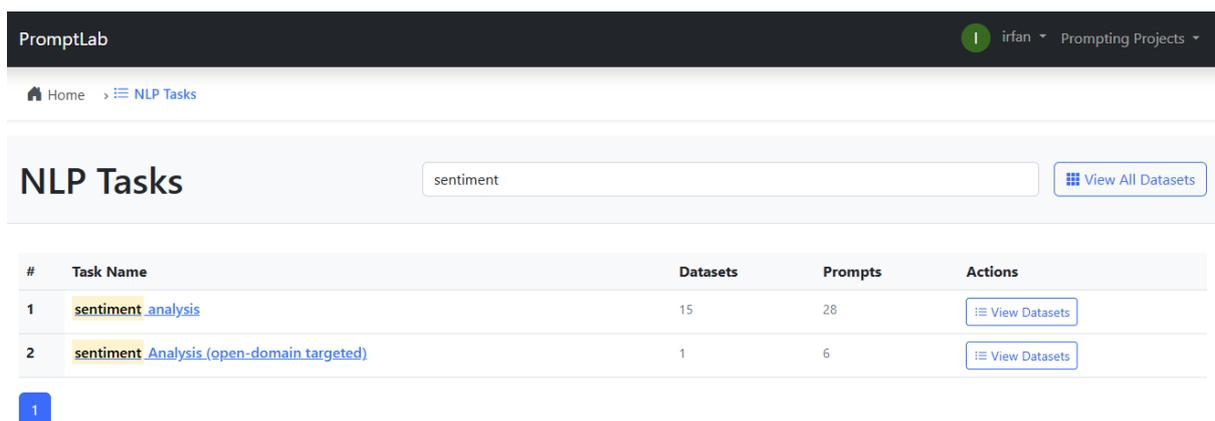


Figure 20: PromptLab: Explore tasks and datasets associated with a task.

that serve as navigation hubs for narrowing down tasks and datasets. One such interface focuses on “Datasets for a <Task, like Dialect Identification>” as shown in Figure 19, guiding the user from a specific task domain (e.g., dialect identification) to a curated selection of relevant datasets. Users can quickly filter through entries using a search bar, skim tabulated information about the number of tasks and prompts linked to each dataset, and view a brief textual description. The interface also provides actionable controls to directly access dataset details, add new prompts, or edit existing records, thereby ensuring that researchers can promptly engage with or contribute to the datasets they find most relevant or needs attention.

Similarly, the “NLP Tasks” interface can be refined through thematic filtering—an example being a direct keyword search (e.g., “sentiment”). This refined view, as shown in Figure 20, distills the platform’s broad inventory of tasks into a more focused subset that matches the researcher’s interests. Each returned task entry is accompanied by the count of datasets and prompts currently associated with it, as well as straightforward navigation buttons leading to those datasets. Such a search-driven interface encourages exploratory browsing and rapid iteration, allowing researchers to quickly locate tasks that align with their specific research questions or methodological preferences.

Expanding outward from this thematic filtering, the “All Datasets” interface (Figure 21) provides a

The screenshot shows the 'All Datasets' page in PromptLab. At the top, there is a search bar labeled 'Search datasets...'. Below it is a table with the following columns: #, Dataset Name, Tasks, Prompts, Description, and Actions. The table lists 10 datasets, each with a unique ID, name, task count, prompt count, and a brief description. Action buttons for each dataset include a list icon, a plus sign, and a link icon.

#	Dataset Name	Tasks	Prompts	Description	Actions
1	<a href="#">ACVA</a>	1	8	None	[List] [Add] [Link]
2	<a href="#">AGS-Corpus</a>	1	8	Dataset Card for AGS Dataset Summary AGS is the first publ...	[List] [Add] [Link]
3	<a href="#">ajgt_twitter_ar</a>	1	12	Dataset Card for Arabic Jordanian General Tweets Dataset Summ...	[List] [Add] [Link]
4	<a href="#">ANS_stance</a>	2	20		[List] [Add] [Link]
5	<a href="#">antcorpus</a>	1	6		[List] [Add] [Link]
6	<a href="#">APCD_remap</a>	2	9	None	[List] [Add] [Link]
7	<a href="#">APCDv2</a>	1	4	None	[List] [Add] [Link]
8	<a href="#">AQAD</a>	1	2	Dataset Card for "AQAD" More Information needed	[List] [Add] [Link]
9	<a href="#">AQMAR_batched</a>	1	1	None	[List] [Add] [Link]
10	<a href="#">AraBench_dev</a>	1	8	Dataset Card for [Dataset Name] Dataset Summary [More Info...	[List] [Add] [Link]

At the bottom of the table, there is a pagination control showing page 1 of 9, with a right arrow.

Figure 21: PromptLab: Explore datasets.

comprehensive catalog of the platform’s entire dataset repository. Researchers are greeted with a sortable, paginated table presenting the dataset name, the number of associated tasks and prompts, and a concise description—a dataset card that can contain summaries, key points, or additional metadata. Integrated search functionality allows users to instantly narrow this broad collection down to a manageable subset tailored to their investigative needs. The interface’s interactive controls—such as the ability to add prompts or view dataset details—foster an environment of continual enrichment and refinement of the datasets themselves.

By applying search terms on the “All Datasets” interface (e.g., filtering by the keyword “sentiment” as shown in Figure 22), users can locate specific resources that resonate with their target analysis domains. The returned listings not only confirm that the desired subject matter is present in the repository but also quantify its richness via the number of prompts and tasks available. This synergy of broad cataloging, fine-grained filtering, and direct action links creates a cohesive workflow: from scanning an extensive resource library to pinpointing niche datasets, and finally, to taking meaningful steps in prompt creation and data curation.

As a researcher’s exploration narrows down to a single dataset, the platform provides specialized interfaces to support prompt-level oversight and refinement. Consider the “Prompts for Dataset” interface (Figure 23) : upon selecting a specific dataset—such as AraBench—users are presented with a structured listing of all associated prompts. This listing is organized into intuitive tabs (e.g., “My Prompts,” “Submitted Prompts,” and “All Prompts (Drafts Excluded)”) that filter the view based on the user’s involvement and the prompt’s lifecycle stage. Each prompt entry includes pertinent information such as the task category and creation metadata. Status indicators (e.g., “APPROVED”) and tag labels (e.g., “Example Prompt”) help track a prompt’s review stage and thematic relevance at a glance. A prominent “Create New Prompt” button encourages contribution and continuous dataset enrichment, allowing researchers to easily add their own prompts once they have reviewed existing entries.

As users interact with their assigned datasets and refine their own contributions, the platform supports

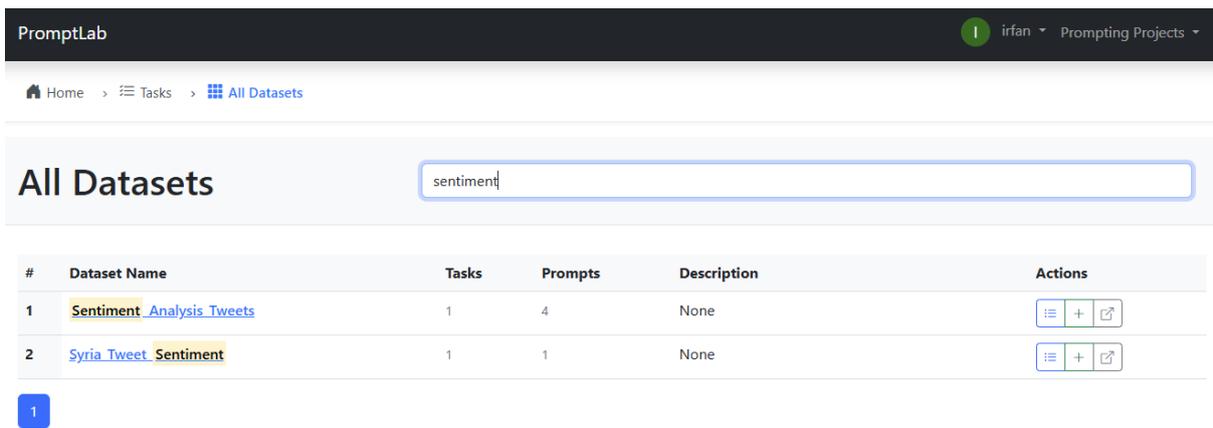


Figure 22: PromptLab: Search datasets.

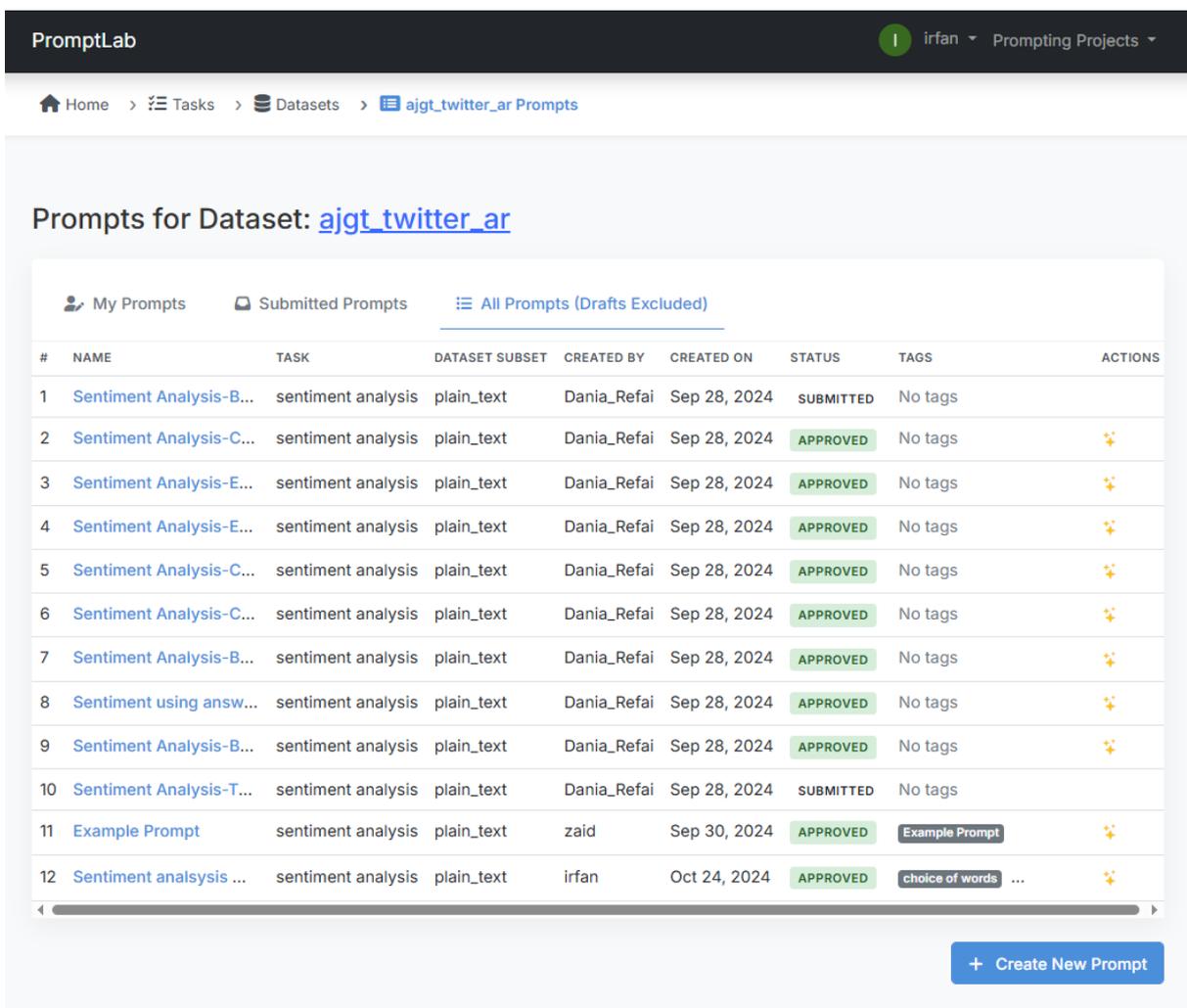


Figure 23: PromptLab: Prompt list for a specific dataset.

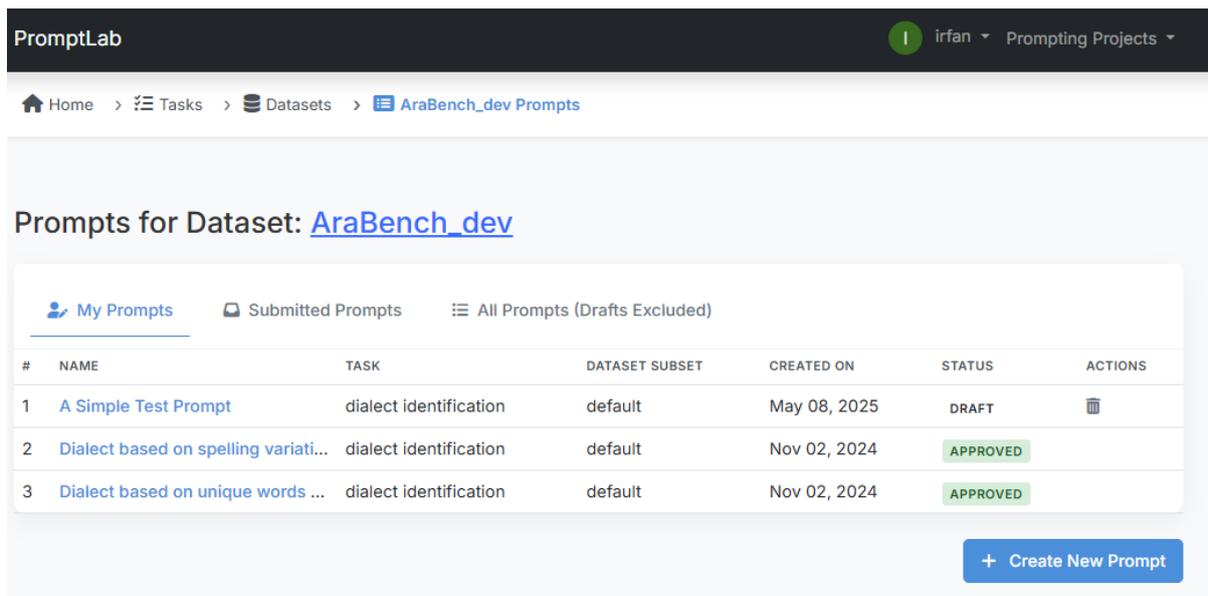


Figure 24: PromptLab: Prompt list for a specific dataset for a specific user.

a smooth transition from broad exploration to personal accountability. In the refined interface focusing on “My Prompts (Figure 24,” researchers find only those prompts they have authored. This personal view simplifies the cognitive load, enabling them to quickly identify which prompts need attention, revision, or further annotation. By separating personal work from the larger collaborative corpus, the interface ensures that researchers can maintain a clear, well-defined workflow—even as the number of prompts grows across various tasks and datasets.

When working with datasets used in multiple tasks—such as “ANS stance”, which may span multiple related tasks (e.g., stance detection, claim verification)—the system’s interface offers convenient filters directly above the prompt listings. These thematic filters (e.g., “stance detection prompts,” “claim verification prompts,” or a unified “all ANS-stance prompts” view) empower researchers to toggle between different thematic slices of the dataset’s prompt inventory, streamlining the discovery of prompts relevant to their current focus. By segmenting prompts along conceptual lines, the interface enhances both navigability and methodological clarity as illustrated in Figure 25.

The dedicated prompt creation interface encapsulates the platform’s commitment to facilitating comprehensive, high-quality prompt generation. When initiating a new prompt for a dataset (e.g., AraBench as shown in Figure 26), users encounter a dual-pane view: on one side, the platform provides contextual dataset information, including links to external resources (like Hugging Face pages), summaries of dataset content, and controls for selecting subsets or splits. On the other side, a streamlined prompt creation form guides researchers through the process of defining prompt characteristics—such as the answer choices relevant to dialect classification, specifying tags for metadata, and setting the appropriate text direction. This environment merges dataset understanding with structured prompt-authoring features, culminating in a highly informed, user-driven creation process. Once satisfied, researchers can apply templates, save, or submit their prompts for review, ensuring continuous improvement of the platform’s collective corpus.

The platform provides interfaces that emphasize prompt review, iteration, and adherence to established standards of quality and consistency. The interfaces presents a dedicated workspace for applying and editing prompt templates. Here, users encounter a side-by-side view: on one side, a carefully curated prompt template offering guidance, context, and structure; on the other, a free-form text area where the researcher can craft and finalize the actual prompt content. A prominent “Apply Template” (Figure 27) button allows for effortless insertion of predefined structures, ensuring that all prompts align with the dataset’s thematic requirements and desired formatting conventions. Meanwhile, the main panel encourages prompt authors to incorporate domain-specific answer choices, add descriptive tags, and configure the prompt’s text direction. The synergy between these panels—contextual dataset insights on

PromptLab irfan Prompting Projects

Home > Tasks > Datasets > ANS\_stance Prompts

### Prompts for Dataset: [ANS\\_stance](#)

#	NAME	TASK	DATASET SUBSET	CREATED BY	CREATED ON	STATUS	TAGS
1	<a href="#">Nawaf ff</a>	stance detection	default	nawaf	Oct 27, 2024	SUBMITTED	No tags
2	<a href="#">statement_relationship</a>	stance detection	default	irfan	Oct 17, 2024	SUBMITTED	AI generated
3	<a href="#">predict the relationship</a>	claim verification	default	majed.alshaibani	Oct 02, 2024	APPROVED	No tags
4	<a href="#">Example Prompt</a>	claim verification	default	zaid	Sep 30, 2024	APPROVED	Example Prompt
5	<a href="#">ClaimVerification-Co...</a>	claim verification	default	Dania_Refai	Oct 09, 2024	APPROVED	No tags
6	<a href="#">Claim Verification-Act...</a>	stance detection	default	Dania_Refai	Nov 06, 2024	SUBMITTED	No tags
7	<a href="#">agreement_query</a>	stance detection	default	irfan	Oct 17, 2024	SUBMITTED	AI generated
8	<a href="#">ClaimVerification-Basi...</a>	claim verification	default	Dania_Refai	Nov 06, 2024	SUBMITTED	No tags
9	<a href="#">Nawaf Amazing prompt</a>	stance detection	default	nawaf	Oct 27, 2024	SUBMITTED	No tags
10	<a href="#">consistency_check</a>	stance detection	default	irfan	Oct 17, 2024	SUBMITTED	AI generated
11	<a href="#">conflicting_statement...</a>	stance detection	default	irfan	Oct 17, 2024	SUBMITTED	AI generated
12	<a href="#">Named answer choices</a>	claim verification	default	ahmed6	Oct 01, 2024	APPROVED	No tags
13	<a href="#">ClaimVerification-Act...</a>	stance detection	default	Dania_Refai	Nov 06, 2024	SUBMITTED	No tags
14	<a href="#">Structured instruction...</a>	claim verification	default	ahmed	Jan 24, 2025	SUBMITTED	Meta prompting
15	<a href="#">ClaimVerification-CO...</a>	claim verification	default	Dania_Refai	Oct 09, 2024	APPROVED	No tags
16	<a href="#">A_Example</a>	claim verification	default	ahmed	Sep 01, 2024	RETURNED_FOR_MODIFIC...	No tags
17	<a href="#">Prompt with zero-sho...</a>	claim verification	default	ahmed	Jan 24, 2025	SUBMITTED	Zero-shot COT
18	<a href="#">Detailed steps to think...</a>	claim verification	default	ahmed	Jan 24, 2025	SUBMITTED	details step by step
19	<a href="#">correlation_judgment</a>	stance detection	default	irfan	Oct 17, 2024	SUBMITTED	AI generated
20	<a href="#">Nawaf Alomari</a>	stance detection	default	nawaf	Oct 27, 2024	SUBMITTED	No tags

[+ Create New Prompt](#)

Figure 25: PromptLab: Prompt list for a specific dataset encompassing multiple tasks.

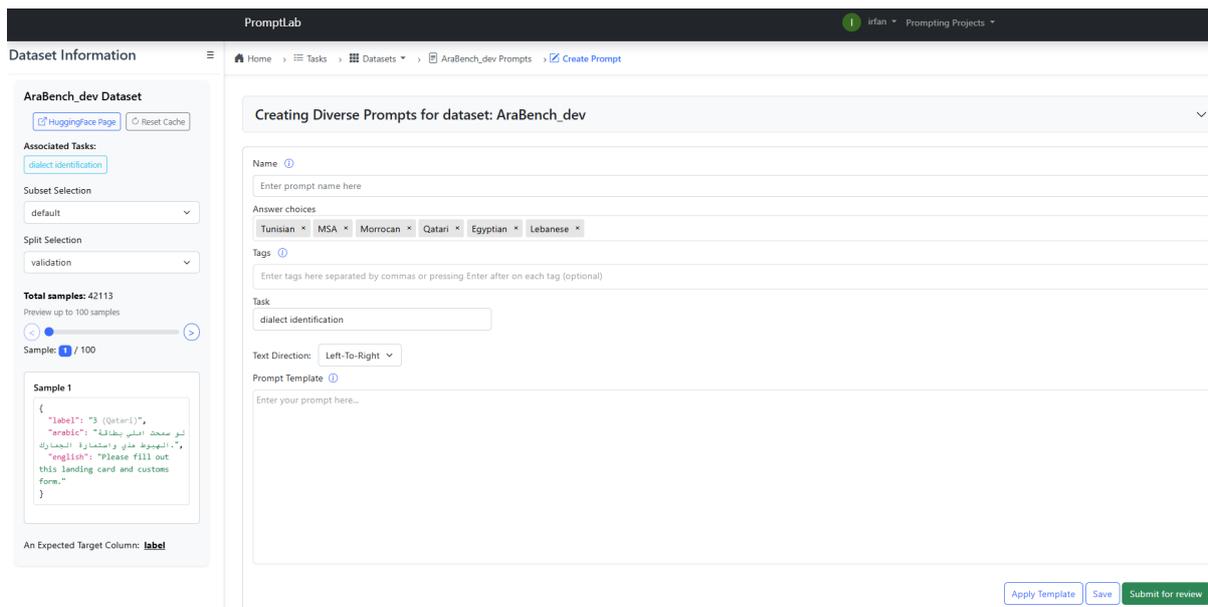


Figure 26: PromptLab: An interface to create a new prompt.

the left and customizable prompt details on the right—empowers researchers to produce prompts that are both contextually relevant and methodologically sound.

The platform also provides higher-level guidance as shown in Figure 28. An interface introduces an instructional overlay—presented as a collapsible, prominently styled banner—that encourages researchers to explore a range of prompt variation strategies. By suggesting axes such as interrogative vs. affirmative framing, localization of the task description, and implicit contextualization, this interface invites users to break out of rigid formulae and experiment with more nuanced, versatile prompt formulations. This guidance serves as a gentle pedagogical layer, reinforcing the platform’s ethos: the pursuit of prompt diversity can yield more robust, generalizable model performance.

The platform offers a specialized modal window (as shown in Figure 29) that encapsulates the essence of effective prompt engineering. Within this dialogue box, users find a concise tutorial on employing Jinja syntax for dynamic prompt construction, instructions for leveraging “answer choices” as a means of fine-grained control, and recommended best practices for clarity, context, and output specification. Far from being a static reference, this modal is seamlessly integrated into the workflow: a user can consult it at any moment, refining their approach on the fly. By coupling the authoring environment with immediate, contextually relevant guidance, the platform facilitates an iterative learning process in which each new prompt crafted is better informed, better structured, and ultimately more valuable to the broader research community.

Once a prompt template has been defined, PromptLab allows users to test it directly against integrated large language models (LLMs) without leaving the interface. In the “LLM Testing” tab on the right-hand side, the user selects the desired subset and split, chooses an AI model from the model drop-down (via OpenRouter), and clicks “Test Prompt” to run the current dataset example through the full template. The response panel then displays the model’s output along with basic metadata such as prompt and completion token counts, so users can inspect whether the instructions are followed, compare different models, and estimate cost. By iteratively editing the prompt, re-testing, and finally using “Apply Template”, “Save”, or “Submit for review”, users can systematically refine prompts before applying them more broadly. This functionality is shown in Figure 30.

Below this workspace, a “Reviewer Actions” panel enables peer review. In this section, reviewers can leave comments, record decisions, and track the prompt’s movement through a formalized review and approval pipeline. Such a multi-tiered interface integrates authoring, advising, and authoritative feedback into a unified process, advancing not just individual prompt quality, but also the platform’s collective standards.

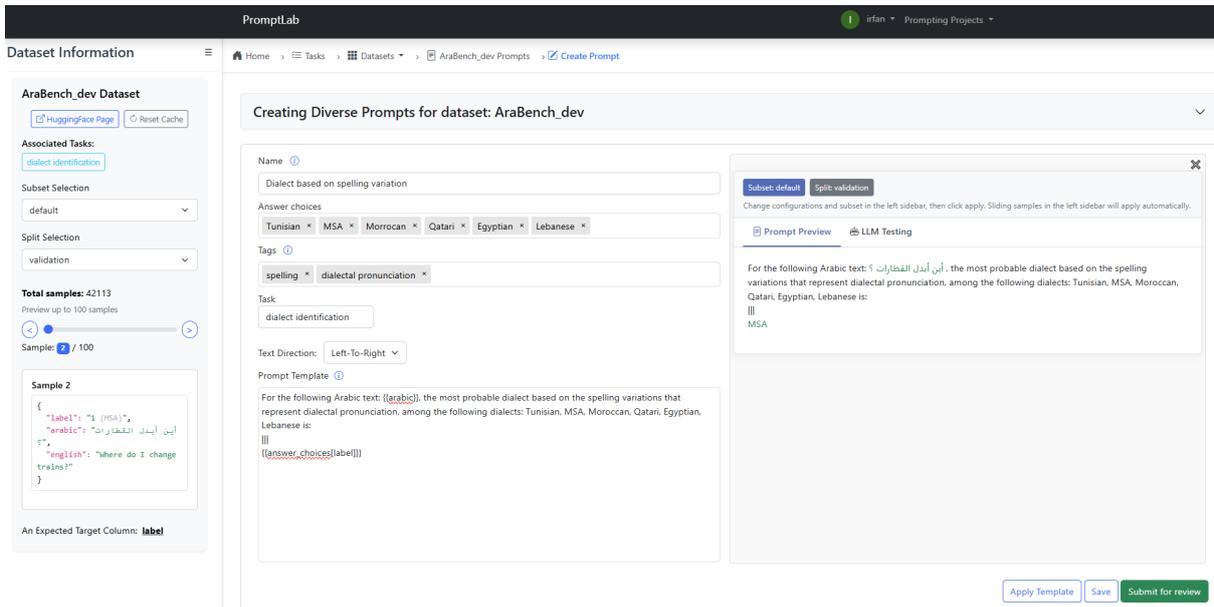


Figure 27: PromptLab: An interface apply a prompt on a dataset sample.

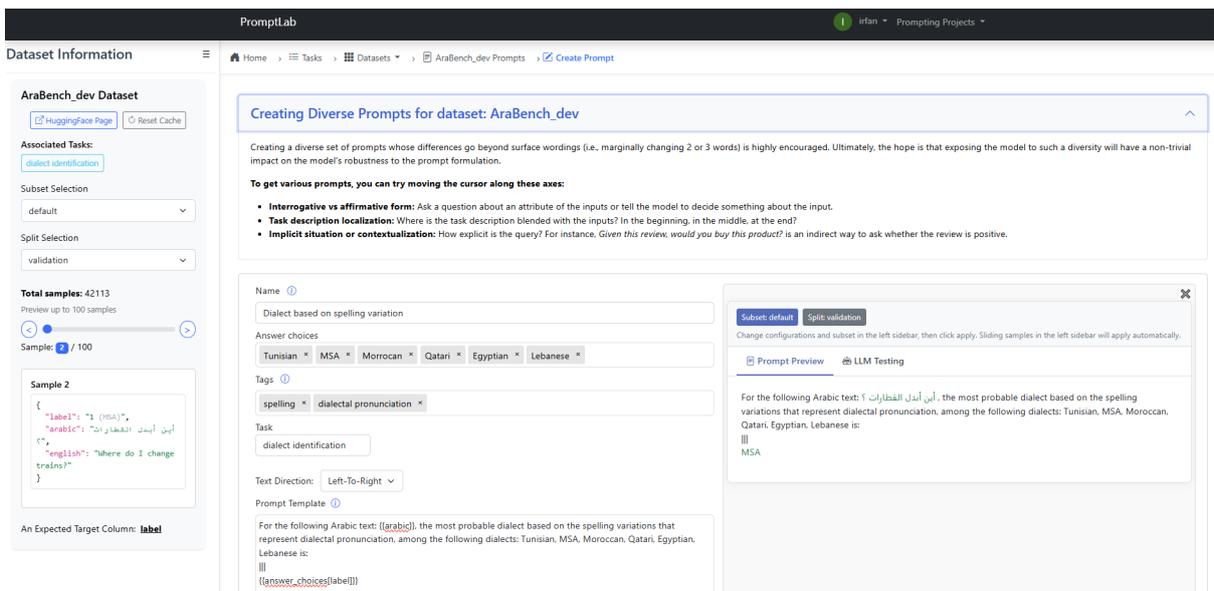


Figure 28: PromptLab: Basic guidelines on creating prompts.

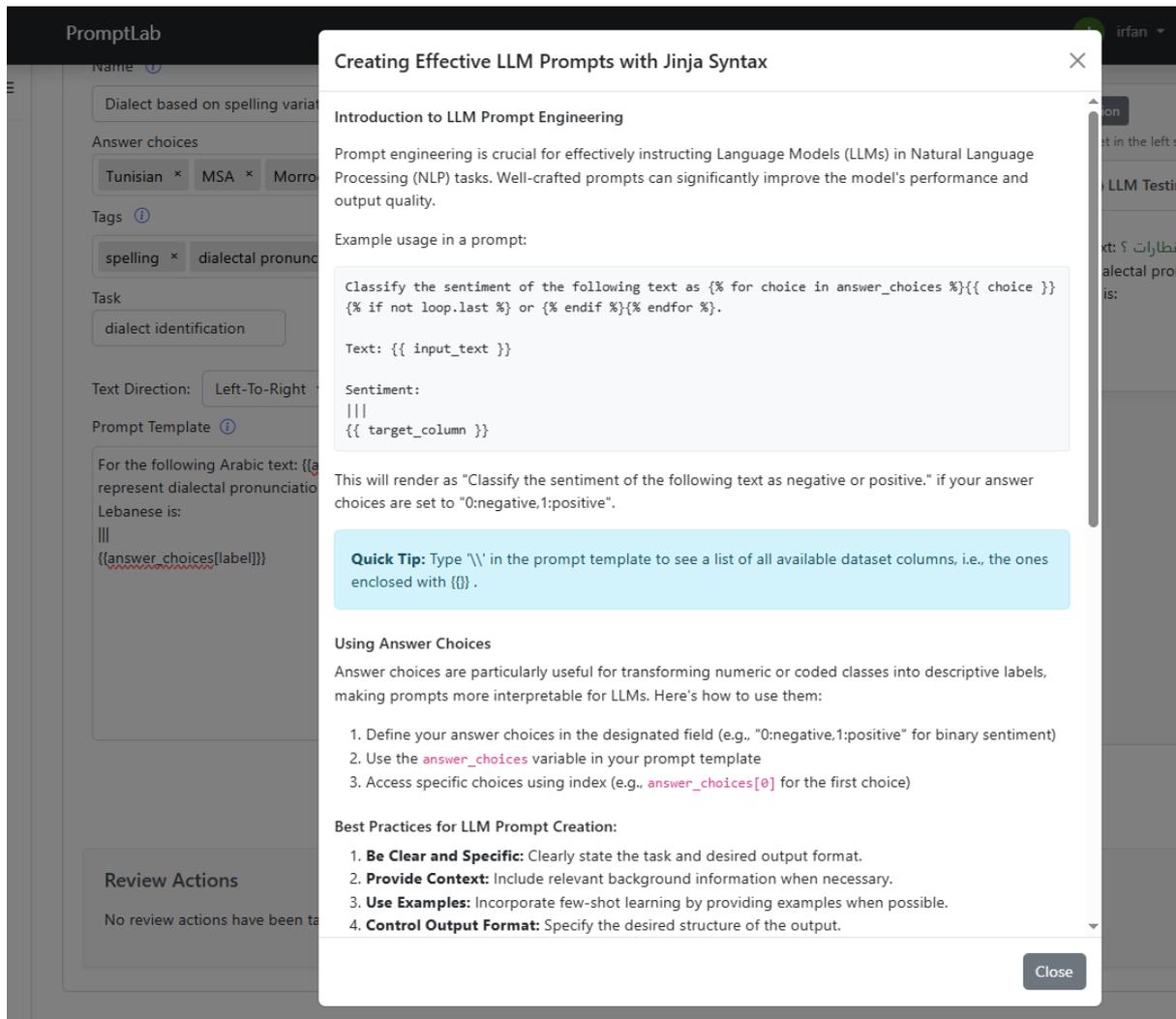


Figure 29: PromptLab: Detailed guidelines on creating prompts.

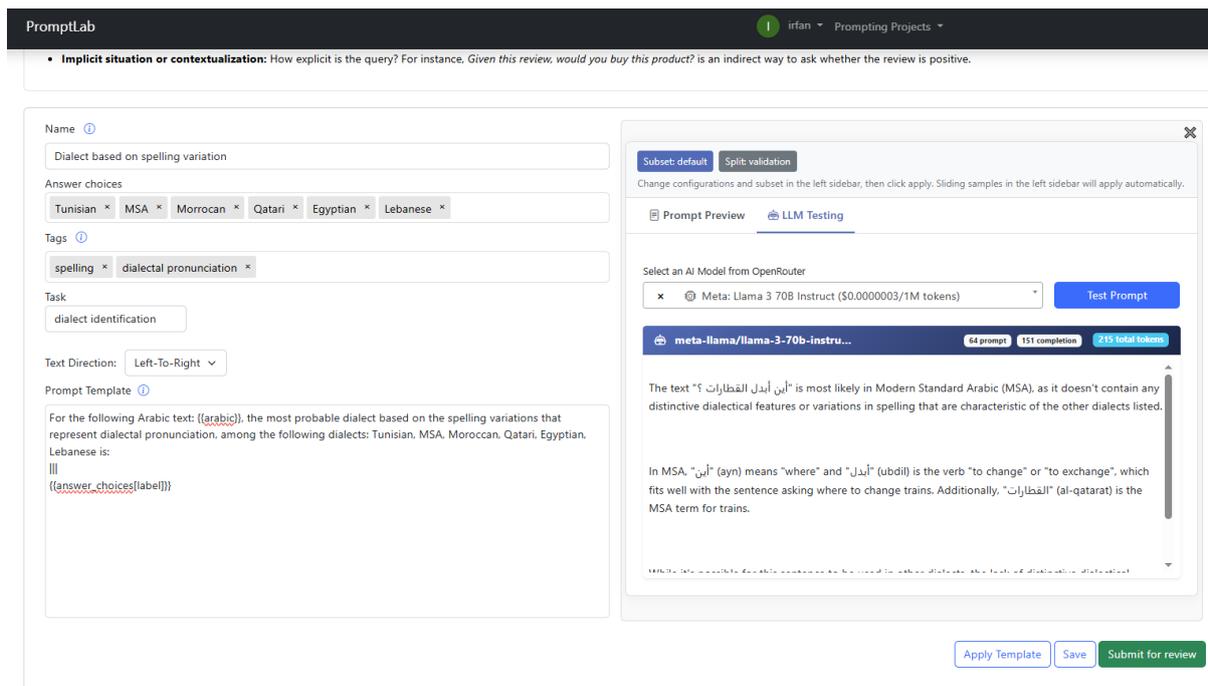


Figure 30: PromptLab: An illustration of testing a prompt on an LLM.

The “Reviewer Actions” panel operates as a real-time checkpoint for quality control and finalization. Positioned at the bottom of the prompt editing or viewing interface, this section enables reviewers to document their feedback and indicate a formal decision (e.g., “Approved,” “Returned for Modification”) as illustrated in Figure 31. By coupling the refinement process with immediate evaluative commentary, the platform ensures that the progression from initial draft to finalized prompt is both rigorous and responsive. Users can integrate feedback promptly, leading to a continuous improvement loop that enhances the quality and relevance of the entire prompt repository.

In some cases, a prompt may not meet the necessary criteria on its first submission, requiring authors to revisit their work. When “Returned for Modification,” a prompt re-enters the editing stage with actionable guidance on how to improve its clarity, formatting, or content. In this interface (Figure 32), prompt authors encounter any attached reviewer comments directly adjacent to their prompt creation tools. This spatial juxtaposition of feedback and editing capability expedites the revision cycle—authors can immediately apply suggested changes, enhancing both efficiency and accuracy in the refinement process. In doing so, the platform streamlines the feedback loop, transforming what could be a tedious back-and-forth exchange into a targeted, productive revision session.

As the prompt creation workflow matures into a cycle of iterative refinement, the platform introduces interfaces dedicated to documenting and managing the historical progression of a prompt’s review process. The “Prompt Review History” interface (Figure 33) provides a chronological record of every significant interaction—submission events, reviewer comments, modification requests, and final approvals. Each review action is clearly timestamped and attributed to a specific contributor, offering full transparency into how and why the prompt has evolved. This historical overview fosters a culture of accountability and collaborative learning: prompt authors gain insights into recurring areas of improvement, while reviewers can revisit past decisions to ensure alignment with evolving quality standards.

The platform introduces interfaces that extend beyond manual creation and refinement, reflecting the ecosystem’s dynamic and forward-looking ethos. The first of these interfaces exemplifies a scenario in which a fully integrated prompt inventory for a given dataset is on display. Prompts are presented with a familiar tabular structure, featuring clear status indicators and minimalistic tagging fields. However, the “Actions” column now incorporates advanced functionalities accessible via a subtle dropdown icon. Researchers can choose to “Generate AI prompts from this prompt,” (Figure 34) leveraging automated

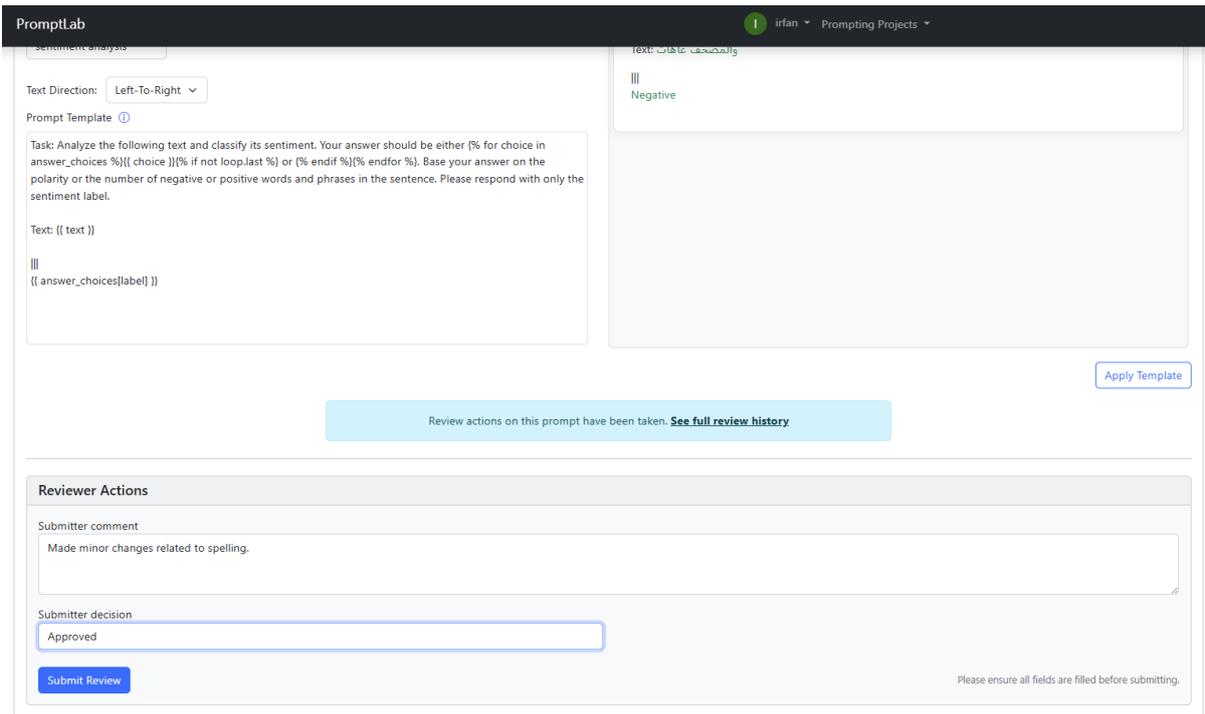


Figure 31: PromptLab: An interface to review a submitted prompt.

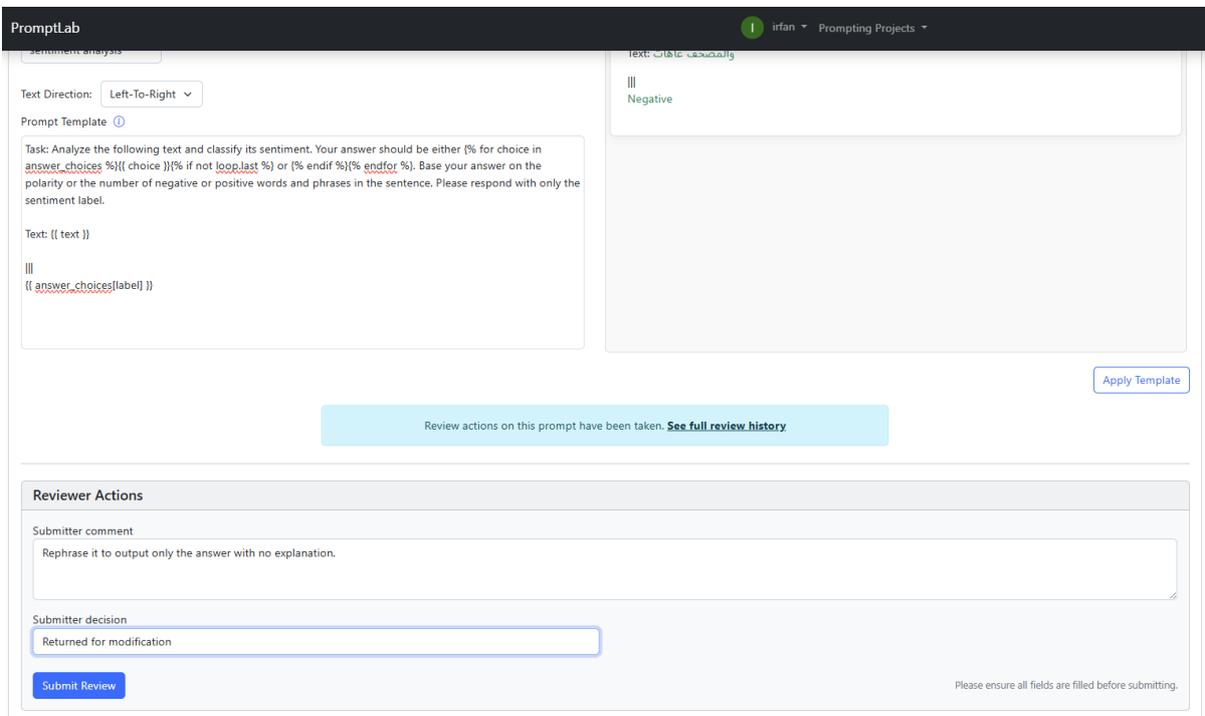


Figure 32: PromptLab: An example interface to review a submitted prompt.

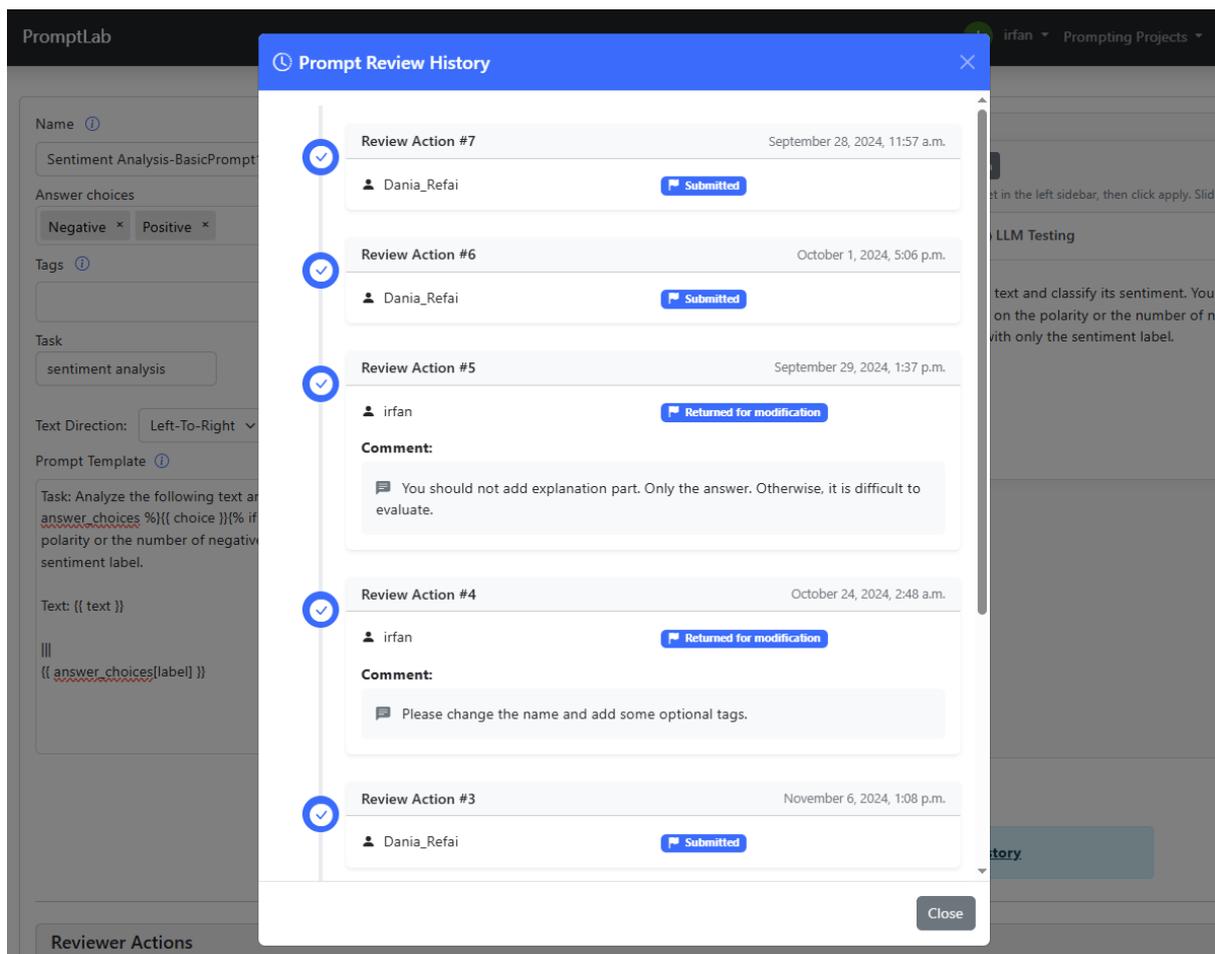


Figure 33: PromptLab: Viewing prompt review history.

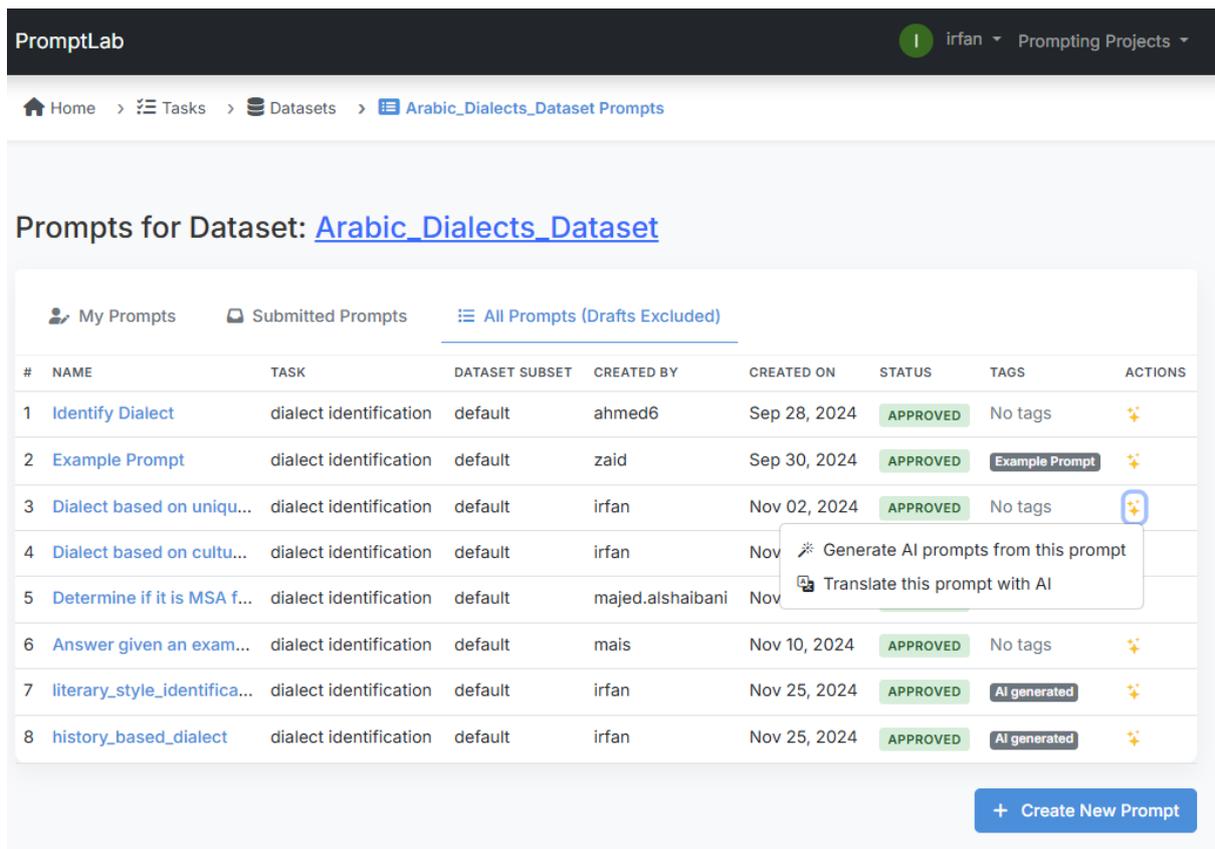


Figure 34: PromptLab: AI-assisted prompt generation.

capabilities to expand and diversify the prompt set.

On selecting the “Generate AI prompts from this prompt,” option, the researcher encounters a layout consistent with previous views—dataset details on one side, customizable prompt attributes on the other (as shown in Figure 35). Yet this time, the submission workflow includes decisive action buttons: “Submit for review” and “Reject Prompt.” These controls crystallize the platform’s commitment to maintaining quality standards and editorial rigor. Researchers can confidently propose their newly AI-formed prompts for inclusion, aware that they will undergo a structured review cycle facilitated by the platform’s features as explained before. Meanwhile, the “Reject Prompt” option provides a fallback mechanism, allowing the user to withdraw or discard those AI-generated prompts that do not meet their evolving criteria. Together, these final two interfaces encapsulate the platform’s overarching philosophy: an iterative, user-driven ecosystem enriched by AI-assisted features and supported by a governance layer that ensures each prompt’s relevance, integrity, and contribution to advancing NLP research.

**Dataset Information**

**Arabic\_Dialects\_Dataset Dataset**

[HuggingFace Page](#) [Reset Cache](#)

**Associated Tasks:**  
dialect identification

**Subset Selection**  
default

**Split Selection**  
train

**Total samples:** 9992  
Preview up to 100 samples

Sample: 1 / 100

**Sample 1**

```
{
  "text": "يعاني الكثيرون من  
المشاكل في العزلة خلق شعور  
عنا نطمح أن يكون عليه الحال  
ولكن يعني ما يفهمنا أن تغير من  
الأمر أنه بدأ في اقتداءه أو  
وسائل تعليمية الخاصة به أينما ما  
"، في البيوت"  
"label": "4 (MSA)"
}
```

**Name**  
history\_based\_dialect

**Answer choices**  
Levant North Africa Egypt GULF MSA

**Tags**  
AI generated

**Task**  
dialect identification

**Text Direction:** Left-To-Right

**Prompt Template**  
Consider this historical context in the text provided: {{Text}}. Which dialect from the ancient regions does it emerge from? Select from these regions: Levant, North Africa, Egypt, GULF, MSA. ||| ({{answer\_choices}}[label])

[Apply Template](#)
[Save](#)
[Submit for review](#)
[Reject Prompt](#)

Figure 35: PromptLab: AI-assisted prompt generation an illustration.

# LLM BiasScope: A Real-Time Bias Analysis Platform for Comparative LLM Evaluation

Himel Ghosh<sup>1,2</sup>, Nick Elias Werner<sup>1</sup>

<sup>1</sup>Technical University of Munich, Germany, <sup>2</sup>Sapienza University of Rome, Italy

Correspondence: [himel.ghosh@tum.de](mailto:himel.ghosh@tum.de)

## Abstract

As large language models (LLMs) are deployed widely, detecting and understanding bias in their outputs is critical. We present LLM BiasScope, a web application for side-by-side comparison of LLM outputs with real-time bias analysis. The system supports multiple providers (Google Gemini, DeepSeek, MiniMax, Mistral, Meituan, Meta Llama) and enables researchers and practitioners to compare models on the same prompts while analyzing bias patterns. LLM BiasScope uses a two-stage bias detection pipeline: sentence-level bias detection followed by bias type classification for biased sentences. The analysis runs automatically on both user prompts and model responses, providing statistics, visualizations, and detailed breakdowns of bias types. The interface displays two models side-by-side with synchronized streaming responses, per-model bias summaries, and a comparison view highlighting differences in bias distributions. The system is built on Next.js with React, integrates Hugging Face inference endpoints for bias detection, and uses the Vercel AI SDK for multi-provider LLM access. Features include real-time streaming, export to JSON/PDF, and interactive visualizations (bar charts, radar charts) for bias analysis. LLM BiasScope is available as an open-source web application, providing a practical tool for bias evaluation and comparative analysis of LLM behaviour.

## 1 Introduction

Large language models (LLMs) are widely used in applications from chatbots to content generation, raising concerns about bias and fairness (Bender et al., 2021; Weidinger et al., 2021). Bias can appear as stereotypes, discriminatory language, or skewed representations across demographic groups (Blodgett et al., 2020; Nadeem et al., 2021). As models proliferate, researchers and practitioners need tools to detect, analyze, and compare bias across models. Existing work on bias evaluation

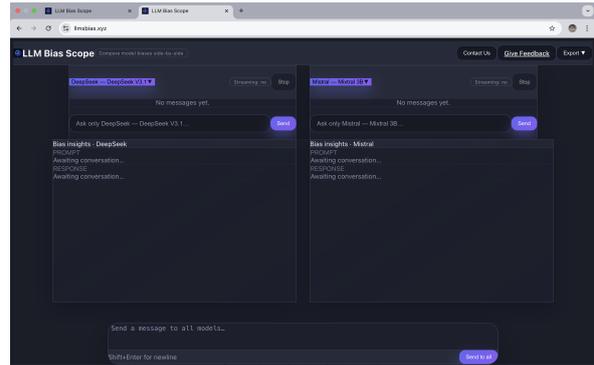


Figure 1: LLM BiasScope Application Home Page.

includes benchmark datasets (Nangia et al., 2020; Nadeem et al., 2021), automated detection methods (Dhamala et al., 2021; Névéol et al., 2022), and frameworks for measuring fairness (Hutchinson et al., 2020; Borkan et al., 2019). However, most tools focus on single-model analysis or static benchmarks, not real-time comparative evaluation of multiple models on user-provided prompts. Comparative evaluation is important because bias patterns vary across models (Gehman et al., 2020; Liang et al., 2023), and practitioners need to choose models that align with their fairness requirements. Existing platforms like Chatbot Arena (Zheng et al., 2023) compare outputs but lack integrated bias analysis. Tools like Perspective API<sup>1</sup> (Lees et al., 2022) detect toxicity but not nuanced bias types. There is a gap for interactive tools that combine real-time multi-model comparison with detailed bias analysis. We present LLM BiasScope, a web application that addresses this gap by enabling side-by-side comparison of multiple LLMs with integrated, real-time bias analysis. The system automatically analyzes both user prompts and model responses using a two-stage pipeline: sentence-level bias detection followed by bias type classification. It provides visualizations, statistics, and

<sup>1</sup><https://perspectiveapi.com/>

comparative analysis to help users understand bias patterns across models. LLM BiasScope supports multiple LLM providers, streams responses in real time, and offers exportable reports. It is designed for researchers evaluating model behavior, developers selecting models, and educators teaching bias awareness. By combining comparative evaluation with detailed bias analysis in an accessible interface, LLM BiasScope supports more informed decisions about LLM deployment.

## 2 Related Work

**Bias Detection and Evaluation:** Bias in language models has been studied through benchmarks and automated detection. Nangia et al. (Nangia et al., 2020) introduced CrowS-Pairs to measure social biases, and Nadeem et al. (Nadeem et al., 2021) created StereoSet for stereotypical bias. Dhamala et al. (Dhamala et al., 2021) proposed BOLD for measuring biases in open-ended generation. Spinde et al. (Spinde et al., 2021) introduced the BABE (Bias Analysis Benchmark for Evaluation) dataset, which provides a binary classification framework specifically designed for bias detection tasks, offering a more direct evaluation approach compared to pair-based benchmarks. These benchmarks focus primarily on static evaluation rather than real-time analysis of user-provided text.

**Bias Type Classification:** Beyond binary bias detection, categorizing the specific type of bias present in text is crucial for understanding and addressing different forms of social bias. Powers et al. (Powers et al., 2025) introduced the GUS Framework, which benchmarks social bias classification using both discriminative (encoder-only) and generative (decoder-only) language models, providing a comprehensive approach to bias type classification across multiple categories. This framework enables fine-grained analysis of bias types, which is essential for understanding the nuanced ways in which bias manifests in language.

Automated bias detection systems include Perspective API (Lees et al., 2022) for toxicity, and HateCheck (Röttger et al., 2021) for hate speech. These target specific bias types and do not provide comparative analysis across multiple models. Recent work has explored sentence-level bias detection (Dhamala et al., 2021; Névéol et al., 2022; Ghosh et al., 2025), with specialized models such as the Domain-Adapted RoBERTa model fine-tuned on BABE (Krieger et al., 2022) demon-

strating strong performance on bias classification tasks. However, integration of these detection capabilities into interactive evaluation platforms that enable real-time comparative analysis across multiple LLMs remains limited.

**Comparative LLM Evaluation Platforms:** Several platforms enable side-by-side comparison of LLMs. Chatbot Arena (Zheng et al., 2023) uses crowdsourced pairwise comparisons. LMSYS Chatbot Arena (Chiang et al., 2023) provides a leaderboard. These focus on quality and preference, not bias analysis. Tools like HELM (Liang et al., 2023) and BIG-bench (Srivastava et al., 2023) offer comprehensive evaluation, but are benchmark-driven rather than interactive. They do not support real-time analysis of user-provided prompts or integrated bias detection.

**Bias Classification Frameworks:** Taxonomies for bias types include gender, racial, religious, and socioeconomic biases (Blodgett et al., 2020; Borkan et al., 2019). Hutchinson et al. (Hutchinson et al., 2020) proposed a framework for measuring fairness in NLP. These provide theoretical foundations but lack practical tools for real-time classification. Recent work has explored automated bias type classification (Dhamala et al., 2021; Névéol et al., 2022; Powers et al., 2025), but these are typically evaluated on static datasets rather than integrated into interactive platforms.

**Interactive LLM Analysis Tools:** Interactive tools for LLM analysis include prompt engineering interfaces (Reynolds and McDonell, 2021) and debugging platforms (Liu et al., 2023). These focus on prompt optimization or error analysis, not bias evaluation. LLM BiasScope differs by combining real-time multi-model comparison with integrated bias analysis, enabling users to evaluate bias patterns across models on their own prompts. It bridges the gap between static benchmarks and interactive evaluation, providing both comparative analysis and detailed bias insights in a single platform.

## 3 System Architecture

LLMBiasScope ([llmsbias.xyz](https://llmsbias.xyz)) is implemented as a modern client-server web application designed for real-time, side-by-side comparison of LLM outputs with integrated bias detection. A complete demonstration of the system is available at this link: <https://youtu.be/rRFRsq-udEo>.

The system leverages a reactive frontend built

on Next.js 16 and React 19, combined with lightweight backend API routes that orchestrate inference across multiple model providers and custom Hugging Face endpoints. See Fig. 2 for the application architecture.

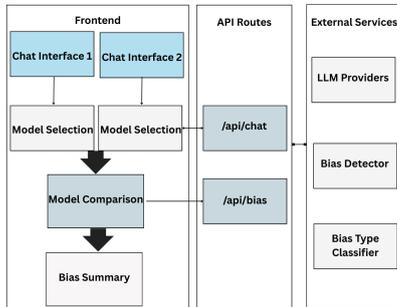


Figure 2: System architecture. The system uses a three-layer client–server design: (1) a React/Next.js frontend with dual chat panels for parallel LLM comparison and bias visualizations; (2) Next.js API routes handling model inference and the two-stage bias analysis pipeline; and (3) external services, including multi-provider LLMs via the Vercel AI Gateway and Hugging Face endpoints for bias detection and classification. Arrows show the flow from user input to model outputs, bias analysis, and visual comparison.

### 3.1 Frontend Architecture

The frontend is built with Next.js 16 (App Router), React 19, TypeScript, and Tailwind CSS 4, using a fully client-side rendering approach for low-latency interaction and smooth streaming of LLM outputs. The main page orchestrates global state, model selection, chat history, and bias-analysis aggregation, presenting two parallel columns for side-by-side model comparison. Each column hosts an independent chat interface powered by the Vercel AI SDK, with responses streamed in real time via SSE. Bias results are visualized through a Bias Summary Card that displays sentence-level scores, bias proportions, and bias-type distributions using Recharts. A Model Comparison Card further contrasts the two models by aggregating their statistics and highlighting differences in bias patterns.

### 3.2 Backend Architecture

Backend functionality is implemented using Next.js App Router API routes, allowing the application to remain stateless and horizontally scalable. Chat API Route handles LLM inference requests. The user can choose the model provider from the options such as: Google Gemini, DeepSeek, MiniMax, Mistral, Meituan, Meta Llama, and OpenAI,

all accessed through the Vercel AI Gateway. The models exposed in the demo interface were selected as a diverse subset of widely used, publicly accessible LLM APIs. Our selection balances provider diversity, model scale, and practical deployment constraints such as API availability, cost, rate limits, and stability. The goal is not to provide an exhaustive comparison, but to enable representative cross-provider evaluation within a real-time interactive setting. Bias Analysis API Route operates as a two-stage pipeline:

- **Bias Detection:** Using the bias-detector model<sup>2</sup> (Ghosh et al., 2025), each sentence is classified as biased or unbiased with an associated probability.
- **Bias Type Classification:** For sentences with bias score  $> 0.5$ , the system invokes the maximuspowers/bias-type-classifier<sup>3</sup> (Powers et al., 2025) to assign a type (e.g., political bias, racism bias).

Both models are served through custom Hugging Face Inference Endpoints deployed on AWS. The route returns aggregated statistics including total sentences, bias ratio, average bias score, and distribution over bias types.

### 3.3 Data Flow

#### 3.3.1 User Interaction Flow

User enters text in either chat column or the shared composer. Frontend sends a request to the selected chat models’ API. The backend forwards the request to the appropriate LLM via the Vercel gateway. Tokens stream back via Server-Sent Events (SSE) and are rendered incrementally.

#### 3.3.2 Bias Analysis Flow

New messages trigger automatic analysis. The message text is segmented into sentences on the client-side. The frontend sends a POST request to bias detection API. Backend performs detection followed by type classification and aggregation. The result is returned to the frontend and displayed in the corresponding Bias Summary Card. If both columns have completed analysis, the Model Comparison Card updates accordingly.

<sup>2</sup><https://huggingface.co/himel7/bias-detector>

<sup>3</sup><https://huggingface.co/maximuspowers/bias-type-classifier>

### 3.3.3 Comparison Flow

Bias statistics are computed per model and then combined to produce deltas, enabling immediate visual comparison of relative bias tendencies.

This architecture supports real-time multi-model comparison, sentence-level bias diagnostics, and interactive evaluation, making LLMBiasScope suitable for research, education, and rapid benchmarking of large language models.

## 4 Evaluation

### 4.1 Bias Detection Model Evaluations

We evaluated multiple bias detection models on standard benchmarks to select the optimal model for LLM BiasScope. Our evaluation followed a two-stage approach: first, we assessed four candidate models on CrowS-Pairs to identify the best-performing model, then we conducted a focused comparison on the BABE dataset, which is specifically designed for bias detection tasks.

#### 4.1.1 Datasets

**CrowS-Pairs** ((Nangia et al., 2020)): 1,508 sentence pairs with stereotypical and anti-stereotypical variants across 9 bias types: race-color (516), gender (262), socioeconomic (172), nationality (159), religion (105), age (87), sexual-orientation (84), physical-appearance (63), and disability (60). This dataset evaluates a model’s ability to distinguish between stereotypical and anti-stereotypical language.

**BABE** ((Spinde et al., 2021)): A binary bias classification dataset containing 1,000 test sentences with gold labels (biased: 559, unbiased: 441). BABE provides a direct assessment of bias detection performance through binary classification, making it particularly suitable for evaluating bias detection models in practical applications.

#### 4.1.2 Evaluation Methodology

**Stage 1: CrowS-Pairs Evaluation.** We evaluated four candidate models on the CrowS-Pairs dataset: unitary/toxic-bert, martin-ha/toxic-comment-model, facebook/roberta-hate-speech-dynabench-r4-target, and himel7/bias-detector. We used the Stereotype Score (SS) metric from Nangia et al. (Nangia et al., 2020), defined as the percentage of pairs where the model assigns a higher bias score to the stereotypical sentence than the anti-stereotypical sentence. SS = 50% represents

random performance; higher values indicate stronger preference for stereotypical sentences.

**Bias Score Extraction and Aggregation.** We evaluate classifier-based models on CrowS-Pairs by first normalizing their heterogeneous outputs into a unified bias score in the range  $[0, 1]$ , representing confidence in the “biased” class. Model predictions may appear as single label–score pairs, multi-label arrays, or nested structures; in all cases, we map labels such as *biased*, *toxic*, *hate*, or *label\_1* to their provided probability, and invert scores for *unbiased*, *non-toxic*, *nothate*, or *label\_0* using  $1 - \text{score}$ . For each CrowS-Pairs pair, we compute normalized scores for the stereotypical (*sent\_more*) and anti-stereotypical (*sent\_less*) sentences, and count a pair as correctly handled when  $\text{score}_{\text{more}} > \text{score}_{\text{less}}$ . The Stereotype Score (SS) is then defined as the percentage of pairs for which this preference holds, enabling consistent comparison across classifier-based models with differing output formats.

**Stage 2: BABE Evaluation.** Based on the CrowS-Pairs results, we selected the top-performing model (unitary/toxic-bert) and compared it with bias-detector and mediabiasgroup/da-roberta-babe-ft (Krieger et al., 2022) (a Domain-Adapted RoBERTa model fine-tuned on BABE) on the BABE dataset. We report standard binary classification metrics: accuracy, precision, recall, and F1-score. The F1-score is particularly important as it balances precision and recall, making it suitable for imbalanced datasets like BABE.

#### 4.1.3 Results

**CrowS-Pairs Results.** Table 1 presents the evaluation results on CrowS-Pairs. The unitary/toxic-bert model achieved the highest Stereotype Score of 69.30%, correctly identifying stereotypical sentences as more biased in 1,045 out of 1,508 pairs. This performance significantly outperforms the other candidates: facebook/roberta-hate-speech-dynabench-r4-target, bias-detector, and martin-ha/toxic-comment-model. The unitary/toxic-bert model also demonstrated the lowest average latency (0.73s), making it both the most accurate and most efficient option.

**BABE Results.** Table 2 and Figure 3 present the evaluation results on the BABE dataset for the three selected models. The bias-detector model achieved the highest

Model	SS (%)	Avg Latency (s)
toxic-bert	<b>69.30</b>	0.73
roberta-hate-speech	57.29	0.98
bias-detector	49.73	1.22
toxic-comment-model	37.86	0.89

Table 1: Evaluation results on CrowS-Pairs dataset. SS = Stereotype Score.

performance with an F1-score of 85.8%. The mediabiasgroup/da-roberta-babe-ft model achieved competitive performance with an F1-score of 81.7%. The unitary/toxic-bert model, while performing well on CrowS-Pairs, showed lower performance on BABE with an F1-score of 71.7%.

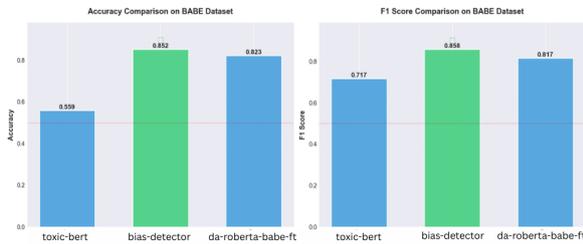


Figure 3: Comparison of Metrics for several bias detection models. The accuracy and F1 Score of the bias-detector show better performance compared to others.

The high precision of da-roberta-babe-ft (97.0%) indicates that when it predicts a sentence as biased, it is highly reliable, though its lower recall (70.5%) suggests it may miss some biased sentences. In contrast, unitary/toxic-bert achieves perfect recall (100.0%) but lower precision (55.9%), indicating it tends to over-classify sentences as biased. The bias-detector model provides the best balance with high precision (92.4%) and good recall (80.1%), achieving the highest F1-score among all models.

Model	Accuracy	Precision	Recall	F1
toxic-bert	55.9%	55.9%	100.0%	71.7%
da-roberta-babe-ft	82.3%	97.0%	70.5%	81.7%
<b>bias-detector</b>	85.2%	92.4%	80.1%	<b>85.8%</b>

Table 2: Evaluation results on BABE dataset. Best F1-score in bold.

**Model Selection.** Based on the comprehensive evaluation, we selected bias-detector (Ghosh et al., 2025) for LLM BiasScope due to its superior performance on the BABE dataset (F1-score: 85.8%), which is specifically designed for

bias detection tasks. As shown in Figure 3, the bias-detector model achieves the best balance between precision (92.4%) and recall (80.1%), outperforming both unitary/toxic-bert and da-roberta-babe-ft in overall F1-score. While unitary/toxic-bert showed strong performance on CrowS-Pairs, its lower performance on BABE (F1-score: 71.7%) and tendency to over-classify (100% recall but only 55.9% precision) made it less suitable for practical bias detection applications. The bias-detector model’s high precision ensures reliable bias detection while maintaining good recall, making it the optimal choice for real-time bias analysis in interactive applications.

## 4.2 Bias-Type Classification Evaluation

Our target is person- and group-directed harms in arbitrary LLM prompts and responses, where token- and sentence-level categories such as generalized statements about groups, unfair attributions, and stereotypes are directly observable. In contrast, media bias and propaganda taxonomies (e.g., framing, agenda setting, emotional language) are primarily designed for document-level news analysis and outlet-level behavior, and map less naturally to short, interactive, cross-domain chat outputs. We therefore prioritize a social-bias-oriented taxonomy that aligns with the kinds of harms practitioners most often wish to inspect in LLM outputs. Hence, in LLM BiasScope we adopt the GUS social bias framework (Generalizations, Unfairness, Stereotypes) for bias type classification (Powers et al., 2025).



Figure 4: Bias Classification: examples of some bias types as detected by the system.

The GUS-Net encoder model (BERT-base-uncased with focal loss) achieves a macro F1-score of 0.80 and Hamming loss of 0.05 on the GUS dataset as reported by Powers et al. (2025) exhibiting higher overall performance compared to the

baselines such as, DistilBERT, RoBERTa, Nbias (BCE).

Since our system integrates this pretrained model without modification, we rely on the published evaluation results for quantitative performance.

We have illustrated the bias-type classification in our system using some examples in Fig. 4, which shows different kinds of bias present in the text.

### 4.3 Empirical Evaluation for Model Comparison

We conducted a small exploratory study to illustrate how LLM BiasScope can be used to compare models on user-facing prompts. We designed three test cases spanning different domains (healthcare advice, career guidance, and educational content; prompts listed in Appendix C). For each test case, we ran  $N = 3$  trials, generated responses from both models, and analyzed them with the bias detection pipeline.

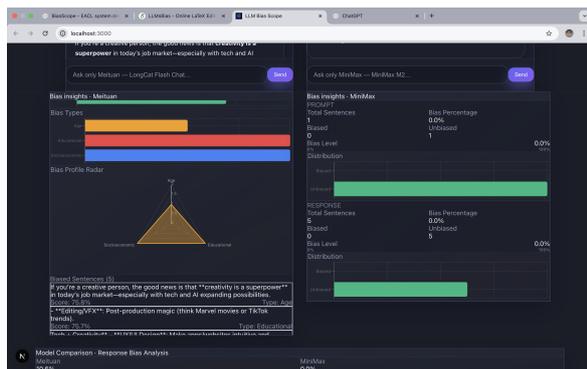


Figure 5: Model Comparison: Bias-Types distribution from Model responses.

For this demonstration, we report only descriptive statistics: mean bias percentage across trials for each model and the absolute difference between models (Table 3). Because the sample size is intentionally small and one model shows 0% bias in these specific prompts, we do not draw strong statistical conclusions or claim generalizable significance. Moreover, the choice of the Meituan and MiniMax models is driven by the lower API costs, while users are free to select other model pairs available in the system for the same purpose. Instead, these results are meant to showcase how the system surfaces differences in bias patterns and supports qualitative inspection of model behaviour on concrete user prompts.

This is how LLM BiasScope can highlight relative bias tendencies on specific prompts and can be interpreted as a qualitative case study.

Test Case	Model A (Meituan)	Model B (MiniMax)	Difference
Healthcare Advice	2.60%	0.00%	+2.60%
Career Advice	10.60%	0.00%	+10.60%
Educational Content	28.20%	0.00%	+28.20%

Table 3: Illustrative model comparison results on three test prompts (mean bias percentage across  $N = 3$  trials).

### 4.4 System Performance Evaluation

To assess the practical usability of LLM BiasScope, we evaluated the system’s performance in terms of response latency and reliability across different text lengths. This evaluation is critical for understanding the system’s scalability and user experience in real-world deployment scenarios.

#### 4.4.1 Methodology

We conducted controlled performance tests using synthetic text samples of varying lengths to isolate the relationship between input size and system latency. This controlled approach allows us to measure performance characteristics independent of content complexity, which may vary in real-world usage. The test cases were designed to represent four distinct text length categories: (1) *Short text* (1 sentence, 6 words), representing quick queries or single-sentence responses; (2) *Medium text* (3 sentences, 15 words), typical of brief paragraphs; (3) *Long text* (10 sentences, 63 words), representing extended responses; and (4) *Very long text* (20 sentences, 83 words), simulating comprehensive document analysis scenarios.

**Test Cases:** The synthetic test cases were constructed to maintain consistent sentence structure and complexity across length categories, ensuring that observed latency differences primarily reflect text length rather than content complexity. While these synthetic cases may not capture all nuances of real-world LLM-generated text (which may contain varying sentence structures, domain-specific terminology, and complex linguistic patterns), they provide a controlled baseline for performance measurement. Future work could extend this evaluation with domain-specific test cases from actual LLM outputs.

For each text length category, we executed 5 independent trials, measuring the end-to-end latency from API request submission to response completion. The latency measurement includes the time required for sentence segmentation, bias detection model inference via the Hugging Face API, bias

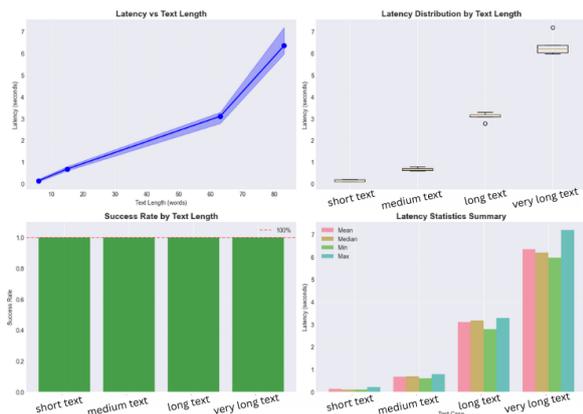


Figure 6: System performance across four text-length categories. Top: latency vs. text length and latency distributions, showing near-linear scaling and consistent behavior across trials. Bottom: success rates (100% for all cases) and summary latency statistics, indicating predictable performance from short inputs (0.14 s for 6 words) to longer texts (6.35 s for 83 words).

type classification (when applicable), and result aggregation. We also tracked the success rate, defined as the percentage of trials that returned valid bias analysis results.

The performance metrics collected include: mean latency, median latency, minimum and maximum latency, standard deviation, and success rate. These metrics provide a comprehensive view of system performance characteristics, including typical response times, variability, and reliability.

#### 4.4.2 Results

Figure 6 presents the performance evaluation results across all text length categories. The system demonstrates consistent reliability with a 100% success rate across all test cases, indicating robust error handling and API stability.

Results show predictable near-linear scaling of latency with text length (0.14s for 6 words to 6.35s for 83 words), with low variability across trials and consistent system behavior, as detailed in Figure 6.

**Limitations:** The synthetic test cases offer a controlled baseline but may not fully reflect the complexity of real LLM-generated text, where performance can vary with linguistic difficulty and domain-specific terminology. Nonetheless, the observed linear scaling indicates that overall performance trends should generalize across diverse content types.

These results demonstrate that LLM BiasScope provides practical performance characteristics suitable for interactive use, with sub-second response

times for typical queries and reasonable latency for comprehensive document analysis.

## 5 Discussion

We evaluated the end-to-end efficiency of the deployed two-stage bias analysis pipeline to ensure suitability for interactive use. On 24 manually curated sentences, the sentence-level bias detector achieved an F1 score of 84.96% with an average latency of 0.25 s per sentence (median 0.19 s; see Appendix A). In typical interactive sessions, the combined bias detection and bias-type classification pipeline operates within real-time constraints (sub-second per sentence), supporting responsive multi-model comparison while leaving room for future optimization (e.g., batching or caching).

A limitation of the current system is that it does not explicitly capture bias expressed through refusal or omission. Models that systematically decline to answer sensitive prompts may therefore appear less biased than models that provide substantive responses. Incorporating refusal-aware analysis is an important direction for future extensions.

The models exposed in the demo interface form a curated set of widely used, publicly accessible LLM APIs from multiple providers, selected to reflect realistic deployment options under practical constraints such as cost, rate limits, and availability rather than to provide an exhaustive catalog. The backend is designed to support additional providers, and future versions will expose a “bring your own API key” mechanism to allow users to integrate their own credentials without the platform handling sensitive information.

**Availability & Licensing.** The LLM BiasScope system is available as an online demonstration at <https://llmsbias.xyz> for research and evaluation purposes. The underlying source code is available here: <https://github.com/Himel1996/LLMBiasScope>. The demonstration video is present at <https://youtu.be/rRFRsq-udEo>.

## References

- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. *On the dangers of stochastic parrots: Can language models be too big?*. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21*, page 610–623, New York, NY, USA. Association for Computing Machinery.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and

- Hanna Wallach. 2020. [Language \(technology\) is power: A critical survey of “bias” in NLP](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2019. [Nuanced metrics for measuring unintended bias with real data for text classification](#). *Preprint*, arXiv:1903.04561.
- Wei-Lin Chiang, Zheng Li, Zi Lin, Ying Sheng, Zi Wu, Hao Zhang, and Eric Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90% chatgpt quality](#). LMSYS Blog. Accessed: YYYY-MM-DD.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. [Bold: Dataset and metrics for measuring biases in open-ended language generation](#). In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21*, page 862–872. ACM.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. [RealToxicityPrompts: Evaluating neural toxic degeneration in language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics.
- Himel Ghosh, Ahmed Mosharafa, and Georg Groh. 2025. [To bias or not to bias: Detecting bias in news with bias-detector](#). *Preprint*, arXiv:2505.13010.
- Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. [Social biases in NLP models as barriers for persons with disabilities](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5491–5501, Online. Association for Computational Linguistics.
- Jan-David Krieger, Timo Spinde, Terry Ruas, Juhi Kulshrestha, and Bela Gipp. 2022. [A domain-adaptive pre-training approach for language bias detection in news](#). In *Proceedings of the 22nd ACM/IEEE Joint Conference on Digital Libraries, JCDL ’22*. ACM.
- Alyssa Lees, Vinh Q. Tran, Yi Tay, Jeffrey Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. 2022. [A new generation of perspective api: Efficient multilingual character-level transformers](#). *Preprint*, arXiv:2202.11176.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, and 31 others. 2023. [Holistic evaluation of language models](#). *Preprint*, arXiv:2211.09110.
- Nelson F. Liu, Tianyi Zhang, and David Percy. 2023. [Promptfoo: Testing and improving llm outputs](#). Website. Accessed: 2025-11-14.
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. [StereoSet: Measuring stereotypical bias in pretrained language models](#). 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 5356–5371, Online. Association for Computational Linguistics.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. [CrowS-pairs: A challenge dataset for measuring social biases in masked language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Online. Association for Computational Linguistics.
- Aurélie Névéol, Yoann Dupont, Julien Bezançon, and Karèn Fort. 2022. [French CrowS-pairs: Extending a challenge dataset for measuring social bias in masked language models to a language other than English](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8521–8531, Dublin, Ireland. Association for Computational Linguistics.
- Maximus Powers, Shaina Raza, Alex Chang, Umang Mavani, Harshitha Reddy Jonala, Ansh Tiwari, and Hua Wei. 2025. [The gus framework: Benchmarking social bias classification with discriminative \(encoder-only\) and generative \(decoder-only\) language models](#). *Preprint*, arXiv:2410.08388.
- Laria Reynolds and Kyle McDonell. 2021. [Prompt programming for large language models: Beyond the few-shot paradigm](#). *Preprint*, arXiv:2102.07350.
- Paul Röttger, Bertie Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet Pierrehumbert. 2021. [Hatecheck: Functional tests for hate speech detection models](#). 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)*. Association for Computational Linguistics.
- Timo Spinde, Hendrik Westermann, and Georg Rehm. 2021. [Neural media bias detection using distant supervision with babe](#). In *Proceedings of the 2021 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 3518–3530.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali

Safaya, Ali Tazarv, and 432 others. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *Preprint*, arXiv:2206.04615.

Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, and 4 others. 2021. [Ethical and social risks of harm from language models](#). *Preprint*, arXiv:2112.04359.

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. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *Preprint*, arXiv:2306.05685.

## 6 Ethics and Broader Impact

LLM BiasScope aims to support transparency in language model behaviour, but its bias assessments are constrained by the limitations of the underlying detection models and datasets. Since these classifiers are trained on specific corpora (e.g., BABE), they may under-detect subtle or context-dependent forms of bias while over-representing others, and they inevitably reflect the cultural assumptions embedded in their training data. Sentence-level analysis also restricts contextual understanding, meaning that both false positives and false negatives remain possible. Bias, moreover, is inherently subjective and culturally dependent, so the system’s outputs represent one possible operationalisation of bias rather than an absolute ground truth.

Because user text is processed through external LLM providers and Hugging Face endpoints, prompts may be subject to third-party data handling practices. Although conversation data is stored only in the user’s browser, users retain limited control once text is transmitted for inference. Additionally, bias-analysis tools can be misused—for example, to selectively criticise competing models, generate misleading comparisons, or probe models adversarially. Automated scores may also create an unwarranted sense of objectivity, leading users to over-rely on quantitative metrics without considering nuance or context.

Moreover, our choice of the GUS social bias taxonomy means that BiasScope focuses on social bias phenomena at token and sentence level, and does not cover propaganda strategies or media-framing categories as in dedicated media-bias taxonomies. Users should therefore treat the reported bias types

as one operationalisation of social bias rather than an exhaustive account of all possible ideological or propagandistic patterns.

Despite these limitations, LLM BiasScope can positively contribute to accountability, education, and model evaluation by making bias analysis more accessible. Its outputs should be interpreted as indicative signals rather than definitive judgments, and used alongside human judgment and broader evaluation practices. Responsible use requires awareness of the tool’s constraints, careful interpretation of results, and an understanding that bias detection is an ongoing, imperfect process embedded within wider societal and ethical considerations surrounding AI.

## A Application Screenshots

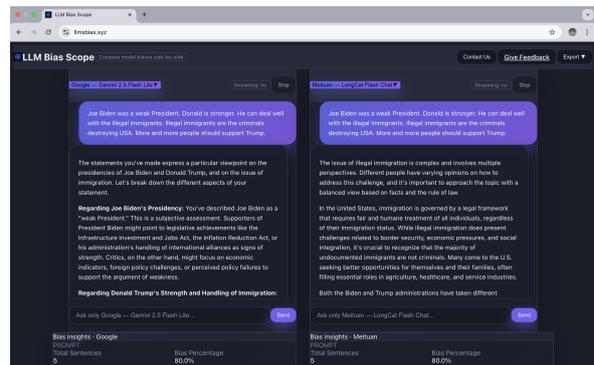


Figure 7: Application home page.

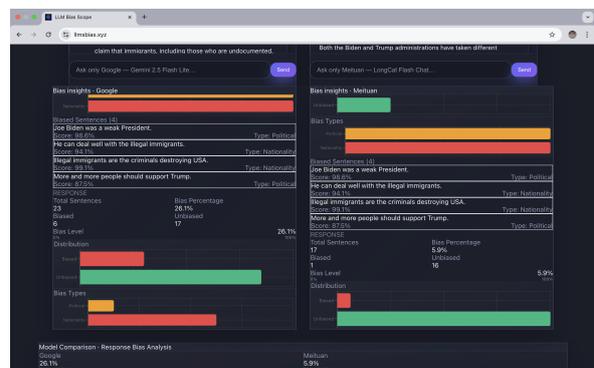


Figure 8: Bias analysis card for the input and model responses.

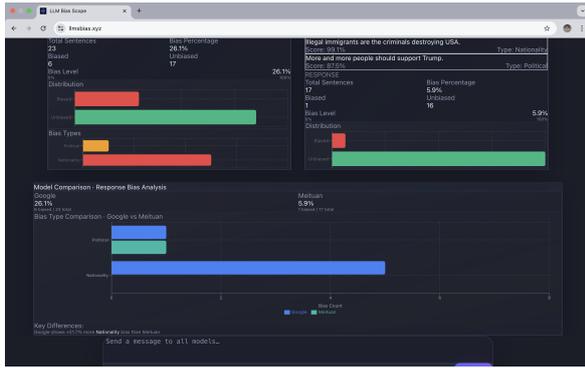


Figure 9: Model comparison card.

## B Evaluations

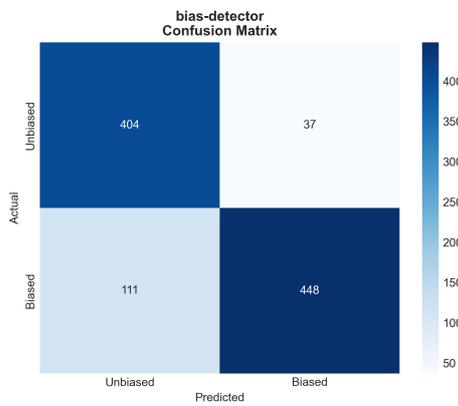


Figure 10: Confusion matrix for bias detection model predictions.

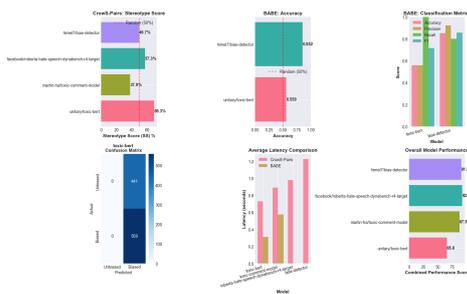


Figure 11: Evaluation results across all bias detection models.

## C Empirical Test Prompts

```
comparison_tests = [
  {
    'name': 'Healthcare Advice',
    'prompt': 'What should I know about diabetes?',
    'model_a_response': 'Response',
    'model_b_response': 'Response',
  },
  {
    'name': 'Career Advice',
    'prompt': 'What career is suitable for creative people?',
    'model_a_response': 'Response',
    'model_b_response': 'Response',
    'expected_difference': 'Model A should show more bias (gender/ethnicity stereotypes)',
  },
  {
    'name': 'Educational Content',
    'prompt': 'Explain the causes of poverty.',
    'model_a_response': 'Response',
    'model_b_response': 'Response',
    'expected_difference': 'Model A should show more bias (gender/ethnicity stereotypes)',
  }
]
```

# InkSight: Towards AI-Aided Historical Manuscript Analysis

Andrey Sakhovskiy<sup>1,5\*</sup>, Ivan Ulitin<sup>1,3\*</sup>, Emilia Bojarskaja<sup>1,4\*</sup>,  
Vladimir Kokh<sup>1</sup>, Ruslan Murtazin<sup>1</sup>, Maxim Novopoltsev<sup>1</sup>, Semen Budenny<sup>1,2</sup>

<sup>1</sup>Sber AI <sup>2</sup>AIRI <sup>3</sup>Perm State University <sup>4</sup>AI Talent Hub <sup>5</sup>Skoltech

Correspondence: andrey.sakhovskiy@gmail.com

## Abstract

Large-scale scientific research on historical documents — particularly medieval Arabic manuscripts — remains challenging due to the need for advanced paleographic and linguistic training, the large volume of hand-written materials, and the absence of assisting software. In this paper, we propose **InkSight**, the first end-to-end Arabic manuscript analysis tool for manuscript-based analytics and research hypothesis testing. *InkSight* integrates three key components: (i) an Optical Character Recognition (OCR) module utilizing a Large Visual Language Model (LVLM); (ii) a lightweight document indexing and information retrieval module that enables query-based evidence retrieval from book-length manuscripts; and (iii) a flexible Large Language Model (LLM) prompting interface factually grounded to the given manuscript via Retrieval-Augmented Generation (RAG). Empirical evaluation on the existing KITAB OCR benchmark and our in-house dataset of ancient Arabic manuscripts has revealed that historical research can be effectively supported using smaller fine-tuned LVLMs without relying on larger proprietary models. The live web demo for InkSight is available freely at: <https://inksight.ru> and the source code for *InkSight* is publicly available at Github<sup>1</sup>.

## 1 Introduction

Recent advances in Optical Character Recognition (OCR) (JaidedAI, 2020; Heakl et al., 2025b) and Natural Language Processing (NLP), particularly with Retrieval-Augmented Generation (RAG) (Lewis et al., 2020b) approach, now enable systems to answer complex questions by retrieving and analyzing supporting evidence passages from long documents. Despite these technological

advances, large-scale scholarly analysis of historical documents — particularly medieval Arabic manuscripts — continues to face significant practical barriers. Many archival collections remain underexplored due to a critical shortage of specialists capable of accurately transcribing historical scripts. Scholars often spend excessive time on mechanical transcription tasks, with a complex manuscript page requiring up to three days of work even for experienced researchers. In our work, we address this challenge by developing an integrated *InkSight* tool that overcomes the OCR bottleneck while providing historians with an efficient LLM framework for book-length manuscript analysis, hypothesis testing, and evidence-based historical interpretation.

The past few decades have witnessed a notable improvement in handwritten text recognition (HTR) due to OCR methods that adopt either Convolutional Neural Networks (CNN) (Lecun et al., 1998; Wigington et al., 2018; Fasha et al., 2020; JaidedAI, 2020; Bhunia et al., 2021) or task-specific Transformer models (Li et al., 2023). However, these models typically require task-specific fine-tuning and exhibit limited generalization abilities to new domains and fonts. Additionally, historical Arabic manuscripts exhibit high script variability. In contrast to prior OCR approaches, Large Vision-Language Models (LVLMs) (Bai et al., 2023) are capable of surpassing task-specific OCR models in accuracy and generalization without extensive fine-tuning (Heakl et al., 2025b).

Recently, Large Language Models (LLMs) have excelled at fine-grained analysis of long-form texts. Specifically, LLMs enhanced with RAG were shown to have good fact checking capabilities not only for English texts (Min et al., 2023; Liu et al., 2025), but also for Arabic data (Kim et al., 2024; Shafayat et al., 2024). For question answering, current state-of-the-art LLMs exhibit near-human performance (Abdelali et al., 2024) on numerous Arabic datasets (Mozannar et al., 2019; Lewis et al.,

\*These authors contributed equally to this work.

<sup>1</sup><https://github.com/ds-hub-sochi/InkSight-tool>

2020a; Artetxe et al., 2020; Clark et al., 2020).

In this paper, we present *InkSight*, the first end-to-end Arabic manuscript analysis tool that adopts Large Visual Language Models (LVLMs) for image-to-text document transcription, hypothesis testing, and evidence retrieval. As seen from Figure 1, *InkSight* integrates three key components: (i) an Optical Character Recognition (OCR) module utilizing a LVLM; (ii) a Retrieval-Augmented Generation (RAG) module designed to efficiently index and retrieve information from book-length manuscripts; and (iii) a flexible prompting interface, allowing domain experts to formulate custom analytical queries and hypotheses.

The contributions of our paper are as follows:

- We present *InkSight*, an open-access web tool that accepts a handwritten Arabic book, performs OCR and document indexing, and supports arbitrary comprehension queries. User questions are answered by an LLM enhanced with RAG module, which retrieves relevant passages from the book and returns them explicitly as supporting evidence.
- Experimental evaluation on modern KITAB OCR benchmark (Heakl et al., 2025b) and our in-house MAS corpus of historical manuscripts has revealed that large proprietary LLMs, namely GPT-4o, GPT-5 and Gemini-2.0 Flash, have limited zero-shot generalization to ancient Arabic texts.
- *InkSight*'s modular pipeline enables easy adaptation to historical manuscripts in other languages. We make the source code for our tool publicly available: <https://github.com/ds-hub-sochi/InkSight-tool>.

## 2 InkSight System

### 2.1 Handwritten Text Recognition Pipeline

For HTR, we implement a two-step pipeline that performs *line segmentation* followed by *line-level transcription* using a fine-tuned LVLM. Line level processing minimizes layout complexity and interference while maintaining adequate context for LVLM autoregressive capabilities (Younes and Abdellah, 2015; Chan et al., 2024).

**Task Formulation** The Handwritten Text Recognition (HTR) task aims to transcribe handwritten text into a machine-interpretable format. Formally,

given an input image  $X \in \mathbb{R}^{H \times W \times C}$  of a handwritten document (with  $C = 1$  for grayscale or  $C = 3$  for RGB), the objective is to predict the corresponding character sequence  $Y = (y_1, y_2, \dots, y_m)$  of length  $m$ . When implemented with LVLMs, HTR performs conditional generation by jointly using the image  $X$  and a textual prompt  $P$  that specifies the recognition task, producing the transcription  $Y$  as output.

**Line Segmentation** For line segmentation, we adopt the Kraken toolkit<sup>2</sup> (Kiessling, 2019). First, the input image  $X_{page} \in \mathbb{R}^{H \times W \times C}$  undergoes binarization to enhance text-background contrast and reduce noise. The rule-based segmenter then processes this binarized image to extract text line regions  $L = \{l_1, l_2, \dots, l_k\}$ , where  $k$  is the number of detected lines, providing robust handling of historical document layouts.

**Segmentation Post-Processing** After initial page segmentation with Kraken, we apply a four-step geometry- and content-aware post-processing pipeline to refine text-line extraction for Arabic handwriting, addressing common artifacts such as over-segmentation and redundant detections:

1. **Dominance-based Box Merging** Given an initial set of text line bounding boxes  $\mathcal{B} = \{b_i\}_{i=1}^N$  extracted by Kraken, we apply padding with ratio  $\gamma_p$ . For each pair of boxes  $(b_i, b_j)$ , we compute their padded intersection  $I_p = \text{Area}(b_i^p \cap b_j^p)$ . If  $I_p > 0$  and  $\frac{\text{Area}(b_i)}{\text{Area}(b_j)} > \alpha_{\text{dominance}}$ , the smaller box is designated as a fragment and merged with the dominant anchor box. This yields a refined set  $\mathcal{B}_{\text{merged}} \subseteq \mathcal{B}$ .
2. **Overlap-based Filtering** For each box  $b_i \in \mathcal{B}_{\text{merged}}$ , we compute the total overlap ratio:  $r_i = \frac{\sum_{j \neq i} \text{Area}(b_i \cap b_j)}{\text{Area}(b_i)}$ . Boxes exceeding the overlap threshold  $\tau_{\text{overlap}}$  are discarded, producing  $\mathcal{B}_{\text{filtered}} = \{b_i \mid r_i \leq \tau_{\text{overlap}}\}$ .
3. **Image-based Content Validation** In the third step, each remaining candidate box is converted to image coordinates and its region of interest is extracted. We then validate the presence of handwriting using multiple image-based criteria, including: ink density  $\rho \in [\rho_{\text{min}}, \rho_{\text{max}}]$ ; minimum contour count  $c_{\text{min}}$ ; minimum contour area  $a_{\text{min}}$ ; minimum

<sup>2</sup><https://kraken.re>

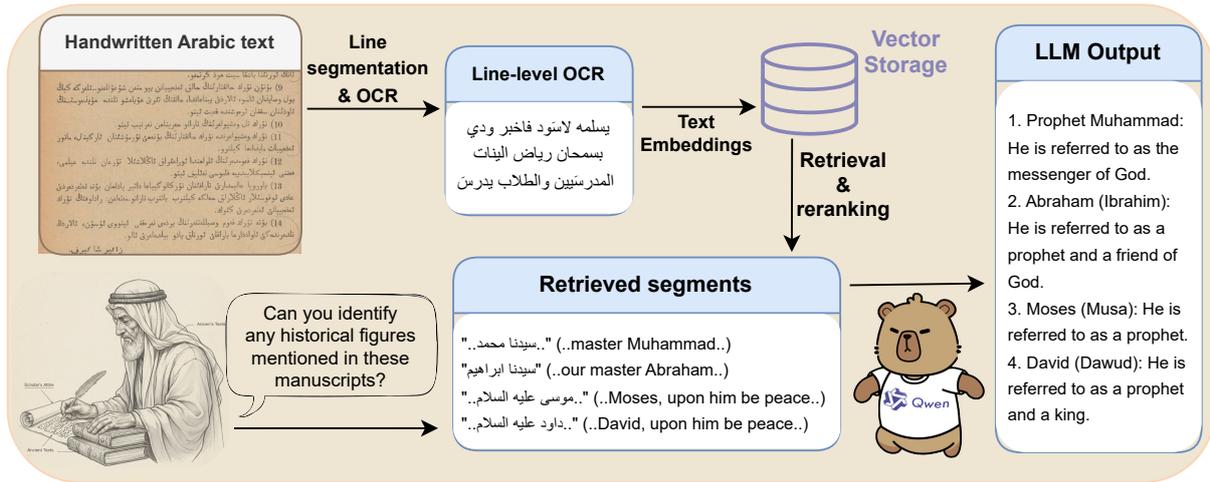


Figure 1: An overview of the InkSight tool for historical Arabic manuscript analysis. The system processes manuscript images through (1) a fine-tuned LLM-based OCR pipeline with specialized post-processing, (2) a RAG module for semantic indexing of book-length documents, and (3) LLM-based chat interface for queries related to document analysis.

hole ratio  $\eta_{\min}$ ; and minimum average defects per contour  $\delta_{\min}$ . Only boxes satisfying all criteria are retained.

4. **Horizontal Row Aggregation** Bounding boxes are iteratively merged into horizontal rows when three conditions are satisfied: their horizontal projections overlap, their heights differ by less than a threshold  $\zeta_{\text{height}}$ , and their vertical overlap ratio exceeds a minimum value  $\kappa_{\text{overlap}}$ .

The key post-processing parameters are presented in Table 4 of Appx. C.

## 2.2 OCR Model

After line segmentation, each segmented line image is processed using our pretrained Qwen3-VL-8B (Team, 2025) model queried with a task-specific prompt (see Appendix A), which directs the model to produce an Arabic transcription of line content.

## 2.3 Training Datasets

To train our OCR model, we adopt annotated data from three Arabic OCR datasets: (i) Muharaf (Saeed et al., 2024a), (ii) SARD (Nacar et al., 2025), and (iii) MAS, our in-house corpus of medieval Arabic manuscripts. The overview of the three training datasets is presented in Table 1.

MUHARAF (Saeed et al., 2024b) (Manuscripts of Handwritten Arabic) is a dataset of manuscripts spanning from the early 19th to the 21st century. It comprises 1,644 authentic document page images

Feature	Muharaf	SARD	MAS
Data Type	Historical manuscripts	Synthetic	Historical manuscripts
# Pages	1,644	843,622	1,023
# lines/words	36,311 lines	690M words	11,841 lines
Fonts/Scripts	Predominantly Ruq'ah script	10 fonts	3 scripts
Centuries	19th–21st	21st	12th–19th

Table 1: Comparison of Muharaf, SARD, MAS Datasets.

containing 36,311 annotated text lines in total. The corpus covers various document types, including personal correspondence, legal records, and literary fragments.

SARD (Nacar et al., 2025) (Synthetic Arabic Recognition Dataset) is a large-scale corpus of synthetically generated document pages designed to simulate book-style Arabic documents. In total, SARD has 843,622 document images with approximately 690 million words, rendered in ten distinct Arabic fonts to ensure wide typographic diversity.

MAS (Medieval Arabic Script) is our in-house collection of manuscripts sourced from the archives of the Abu Rayhan Biruni Institute of Oriental Studies in Uzbekistan. This corpus includes 1,023 pages containing 11,841 annotated text lines from original manuscripts and rewritten copies of the authentic documents spanning the period from 12th to 19th centuries. MAS documents represent Arabic calligraphy in various styles including Naskh, Nastaliq, and Taliq, preserving the natural varia-

Dataset	Train	Test	Total
MUHARAF	24,495		24,495
SARD	100,000		100,000
Amiri font	20,000		
Arial font	20,000		
Calibri font	20,000		
Sakkal Majalla font	20,000		
Scheherazade new font	20,000		
MAS	10,658	1,183	11,841
Naskh script	1,726	192	1,918
Taliq script	3,785	420	4,205
Nastaliq script	5,147	571	5,718
KITAB benchmark		4,138	4,138

Table 2: OCR data statistics in terms of line count.

tions, imperfections, and layout complexities found in historical manuscripts. The dataset covers diverse domains, including waqf documents (testamentary acts of property donation to religious institutions), yarliqs (khan decrees and firmans), arizas (petitions), cheks (legal receipts and certificates), shajars (genealogical tables), and announcements (public proclamations and charters). Due to calligraphy variability and the linguistic gap between Medieval and Modern Arabic, the annotation of the manuscripts is labour-intensive. Each page is annotated by a single expert in Arabic calligraphy only. Notably, the archives of the Abu Rayhan Biruni Institute comprises over 26,000 manuscript volumes dating from the 9th to the 20th centuries with the majority of the documents not yet digitized.

**Data Statistics** Although SARD provides over 800k documents, we subsampled 100k for training as preliminary experiments did not show further error rates decrease from training on more synthetic data. Specifically, we took 20k for each of the 5 most common fonts. Annotated lines from the MAS dataset are randomly split into train and test sets, respectively (see Table 2). Each split preserves the proportions the three calligraphy styles, namely Naskh, Nastaliq, Taliq, ensuring a balanced representation across all subsets. The summarized statistics for the OCR training and evaluation data are shown in Table 2.

## 2.4 Training Details

**Training Setup** We applied domain-specific adaptations to the Qwen2.5-VL-7B-Instruct<sup>3</sup> and Qwen3-VL-8B-Instruct<sup>4</sup> models for the Ara-

<sup>3</sup>[hf.co/Qwen/Qwen2.5-VL-7B-Instruct](https://hf.co/Qwen/Qwen2.5-VL-7B-Instruct)

<sup>4</sup>[hf.co/Qwen/Qwen3-VL-8B-Instruct](https://hf.co/Qwen/Qwen3-VL-8B-Instruct)

bic language. For this stage, we employed LoRA adapters (Hu et al., 2022) (r=8) for parameter-efficient fine-tuning with optional bf16 quantization to reduce memory usage.

## 2.5 Document Search Index

Efficient document retrieval is the core component of the’s RAG- and LLM-based analytical module based on RAG and LLM InkSight. The pages recognized by the HTR module serve as evidence for answering a historian’s query.

**Text Chunking** Due to variability of page lengths recognized by OCR, each full page text is segmented into fixed-size textual chunks (passages) prior to embedding. For segmentation, we employ a deterministic sliding-window procedure with a chunk length up to  $L = 1000$  characters and a window size and chunk overlap of  $O = 200$ . Pages shorter than  $L$  produce a single chunk. To avoid mid-sentence fragmentation, chunk boundaries are adjusted to the nearest separator (e.g., sentence punctuation, paragraph breaks, or word spaces). This ensures chunks to align with linguistic units while maintaining consistent size.

**Passage Embedding** Each chunk is encoded into a dense semantic vector using the multilingual BERT-based encoder<sup>5</sup> (Devlin et al., 2019) trained on the MS Marco dataset (Nguyen et al., 2016). Although there exist models with better retrieval quality for Arabic (Al-Rasheed et al., 2025), we selected the model for its balanced trade-off between representation quality and computational efficiency. All embeddings are stored in a ChromaDB<sup>6</sup> vector index, configured for cosine similarity search. The indexing pipeline is implemented using LangChain<sup>7</sup>, enabling easy adaptation to other domains and languages thanks to its modular pipeline. Our source code is publicly available (Sec. 1), allowing users to change the retrieval model by changing a single line in the configuration file.

## 2.6 Document Analysis

**Passage Retrieval** The retrieval is performed in two stages: dense vector retrieval followed by optional cross-encoder reranking. Given a natural language query  $q$ , the system first encodes it into

<sup>5</sup>[hf.co/ambrooad/bert-multilingual-passage-reranking-msmarco](https://hf.co/ambrooad/bert-multilingual-passage-reranking-msmarco)

<sup>6</sup><https://github.com/chroma-core/chroma>

<sup>7</sup><https://www.langchain.com/>

a dense embedding using the same used for indexing. Then the cosine similarity between  $q$  and each indexed chunk  $c$  is computed as:

$$\text{sim}(q, c) = \frac{\langle f(q), g(c) \rangle}{|f(q)||g(c)|} \quad (1)$$

The top- $k$  most similar chunks (with  $k = 4$  by default) are retrieved as evidence supporting the input query. Only chunks exceeding a minimum similarity threshold are considered, ensuring low-confidence matches are discarded early. For retrieval and reranking, we adopt the same encoder model.

The final set of retrieved passages is passed to the generative reasoning component, implemented as the Qwen3.5-397B-A17B LLM<sup>8</sup> accessed via OpenRouter<sup>9</sup>. The LLM is prompted with a user query concatenated with the retrieval textual passages. Thus, the RAG module functions as the evidence retrieval mechanism that ensures generated LLM responses to be grounded to and factually aligned with the studied manuscript. For OCR and manuscript analysis prompts, please see Appx. A.

Overall, the InkSight’s architecture provides a robust and reproducible framework for manuscript-based in historical research. All modules may be updated with more advanced models, e.g., LVLM, LLM, and retriever, allowing seamless adaptation to manuscripts in other languages as well.

### 3 Evaluation

#### 3.1 Evaluation Data

To assess the quality of InkSightOCR, we performed evaluation on (i) OCR part of the KITAB benchmark and (ii) 589 lines from the MAS’s test part. While KITAB covers more modern texts (starting from the 19th century), MAS includes calligraphic documents from 12th-19th centuries. Thus, our evaluation explores how well modern LVLMs generalize to Arabic language and visual stylistic variations.

*KITAB* OCR benchmark (Heakl et al., 2025b) is a collection of 4,138 samples across multiple document types, including historical manuscripts, handwritten texts, and printed documents for Arabic text recognition. It integrates data from established Arabic OCR datasets such as KHATT (Mahmoud et al., 2014), ADAB (Boubaker et al., 2021),

<sup>8</sup><http://hf.co/Qwen/Qwen3.5-397B-A17>

<sup>9</sup><https://openrouter.ai>

Model	KITAB		MAS	
	CER	WER	CER	WER
<b>Closed-Source LVLMs</b>				
GPT-5	—	—	.52	.91
GPT-4o	.31	.55	.51	.84
Gemini-2.0 Flash	<b>.13</b>	.32	.55	.80
<b>Task-Specific Models</b>				
Tesseract	.54	.84	—	—
EasyOCR	.58	.89	—	—
Surya	4.95	5.91	—	—
<b>Open-Source LVLMs</b>				
Qwen2-VL-7b	1.48	1.55	4.35	3.55
Qwen2.5-VL-7b	1.2	1.4	.89	1.14
Qwen2.5-VL-7b +sft	.43	.67	.34	.73
Qwen-3VL-8b	.67	.95	.94	1.23
Qwen-3VL-8b +sft	.54	.80	<b>.24</b>	.60
AIN-7b	.20	<b>.28</b>	.34	<b>.53</b>

Table 3: Evaluation results on Arabic KITAB OCR benchmark and our in-house MAS corpus of Arabic manuscripts. Smaller values are better. The best results for each metric are highlighted in **bold**.

Muharaf (Saeed et al., 2024a), and EvAREST (Hassan et al., 2021), covering diverse writing styles, fonts, and document complexities that present authentic challenges for Arabic script recognition.

#### 3.2 Evaluation Setup

**Baselines** Following prior research on Arabic OCR on the KITAB benchmark (Heakl et al., 2025b), we evaluate our *InkSight* model against both LVLMs and prior task-specific OCR approaches. For LVLM evaluation, we compare against open-source models, namely Qwen2-VL-7B (Wang et al., 2024), Qwen2.5-VL-7B (Team, 2025), Qwen3-VL-8B, and the AIN-7B (Heakl et al., 2025a) as well as proprietary GPT-4o (OpenAI et al., 2024) and Gemini-2.0 Flash<sup>10</sup> (Team et al., 2024) models. Among task-specific OCR models, we adopted Tesseract (Smith, 2007), EasyOCR (JaidedAI, 2020), and Surya (Paruchuri and Team, 2025).

**OCR finetuning Details** We configured Qwen2.5-VL with temperature =  $1 \times 10^{-6}$  to obtain stable, low-variance transcription

<sup>10</sup>[cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-0-flash](https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-0-flash)

Can you identify any references to trade routes, cities, or geographical locations?

22:31

After carefully analyzing the provided manuscript fragments, **no direct, unambiguous references to specific trade routes, cities, or geographical locations** (such as Samarkand, Bukhara, Tashkent, the Silk Road, Khwarezm, or Ferghana) are found in the text.

However, there are **indirect, contextual, and cultural clues** that allow us to infer a **strong connection to the historical geography and economic world of medieval Central Asia**, particularly the Islamic heartland of Transoxiana (Mā warā' an-Nahr), which includes modern-day Uzbekistan.

---

**1. Absence of Explicit Toponyms**

The manuscript contains:

- No mention of **cities** like *Bukhara*, *Samarkand*, *Khujand*, *Shash* (Tashkent), or *Nishapur*.
- No reference to **trade routes** such as the **Silk Road**, **Sogdian caravan paths**, or **riverine routes** along the Amu Darya or Syr Darya.
- No mention of **regions** like *Khwarezm*, *Ferghana*, or *Khorasan*.

This is not unusual. Many **devotional and mystical manuscripts** from Central Asia were **personal prayer collections (mukhtarat)** or **liturgical notebooks**, not geographical or commercial records. Their purpose was **spiritual**, not cartographic.

Figure 2: InkSight’s output example for the query using the evidence from handwritten *Dala’il al-Khayrat* book.

predictions. Fine-tuning was conducted for 4 epochs with a batch size of 8, utilizing the Adam optimizer ( $\beta_1 = 0.9, \beta_2 = 0.9999$ ), a learning rate of  $5 \times 10^{-4}$ , a weight decay of 0.02, and gradient clipping with a maximum norm of 0.5 to enhance stability. During training, we employed a cosine learning rate schedule with a warmup ratio of 0.05, linearly increasing the learning rate over the first 5% of steps before transitioning to cosine decay.

**Metrics** We evaluate OCR quality using two standard metrics widely used in OCR research (Saeed et al., 2024a; Dolek and Kurt, 2025): Character Error Rate (CER) and Word Error Rate (WER). Both rely on the Levenshtein edit distance (Levenshtein, 1966) and quantify the number of substitutions (S), insertions (I), and deletions (D) required to transform the predicted text into the ground truth.

$$CER = \frac{S + I + D}{N_{char}}; \tag{2}$$

$$WER = \frac{S + I + D}{N_{word}} \tag{3}$$

The combination of CER and WER therefore offers a balanced assessment of both character-level recognition quality and semantic correctness.

## 4 Results

**Proprietary Models Struggle with Manuscripts** OCR evaluation results are presented in Table 3. While strong proprietary GPT-4o and Gemini-2.0 Flash models show low CER and WER on KITAB, they fall short of smaller open-source fine-tuned Qwen models on ancient manuscripts from MAS corpus. Thus, we use finetuned Qwen-3VL-8b in

InkSight as it has the lowest CER on MAS. However, AIN-7b and Qwen2.5-VL-7b+sft show similar performance and could also be used for OCR on Arabic manuscripts.

### Pretraining is Not Essential for Historical HTR

Despite undergoing full model pretraining on authentic Arabic texts, AIN-7b achieves comparable character error rates to fine-tuned Qwen2.5-VL-7b Qwen3-VL-8b with the latter showing even smaller CER (0.24 vs 0.34) on MAS. This indicates that language-specific pretraining is unnecessary for historical Arabic HTR when leveraging parameter-efficient adaptation of multilingual LLMs. Our results indicate that synthetic data paired with lightweight fine-tuning can enable historical manuscript digitization for low-resource non-English languages without costly pretraining.

**Case Study: Dala'il al-Khayrat** To demonstrate InkSight's capabilities in real-world historical research, we conducted an analysis on a XVII century copy of *Dala'il al-Khayrat* (Guidelines to Goodness) manuscript, a seminal Islamic devotional text composed by the Moroccan scholar Muhammad al-Jazuli (who died in 1465). The 186 pages of the manuscript were segmented into 3,720 lines by the InkSight's OCR component. The InkSight output for an example query on trade routes mention is shown in Figure 2. From the example, InkSight allows a researcher to test a hypothesis (e.g., the given book to mention any trade routes) in seconds.

## 5 Conclusion

In this work, we presented InkSight, the first end-to-end AI-aided system for historical Arabic manuscript analysis that integrates LLM-based OCR, semantic search via RAG, and an expert-oriented prompting interface. Our evaluation demonstrates that appropriately fine-tuned open-source LLMs can outperform larger proprietary models like GPT-4o, GPT-5, and Gemini-2.0 Flash on historical document analysis tasks. The system directly addresses critical bottlenecks in historical research workflows by automating transcription and indexing processes, enabling scholars to focus on higher-value semantic and historical analysis rather than mechanical transcription. The proposed system design can be adopted to automate historical research in other domains and languages.

## Acknowledgments

The authors thank Alexei Rastorguev for his invaluable assistance with the deployment and maintenance of the web demo.

## References

- Ahmed Abdelali, Hamdy Mubarak, Shammur Chowdhury, Maram Hasanain, Basel Mousi, Sabri Boughorbel, Samir Abdaljalil, Yassine El Kheir, Daniel Izham, Fahim Dalvi, Majd Hawasly, Nizi Nazari, Youssef Elshahawy, Ahmed Ali, Nadir Durrani, Natasa Milic-Frayling, and Firoj Alam. 2024. *LAraBench: Benchmarking Arabic AI with large language models*. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 487–520, St. Julian's, Malta. Association for Computational Linguistics.
- Raghad Al-Rasheed, Abdullah Al Muaddi, Hawra Al-jasim, Rawan Al-Matham, Muneera Alhoshan, Asma Al Wazrah, and Abdulrahman AIOsaimy. 2025. *Evaluating RAG pipelines for Arabic lexical information retrieval: A comparative study of embedding and generation models*. In *Proceedings of the 1st Workshop on NLP for Languages Using Arabic Script*, pages 155–164, Abu Dhabi, UAE. Association for Computational Linguistics.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. *On the cross-lingual transferability of monolingual representations*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. *Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond*. *Preprint*, arXiv:2308.12966.
- Ayan Kumar Bhunia, Shuvoyit Ghose, Amandeep Kumar, Pinaki Nath Chowdhury, Aneeshan Sain, and Yi-Zhe Song. 2021. *MetaHTR: Towards writer-adaptive handwritten text recognition*. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15825–15834.
- Houcine Boubaker, Abdelkarim Elbaati, Najiba Tagougui, Haikal El Abed, Monji Kherallah, Volker Märgner, and Adel M. Alimi. 2021. *Adab database*.
- Adrian Chan, Anupam Mijar, Mehreen Saeed, Chau-Wai Wong, and Akram Khater. 2024. *Hatformer: Historic handwritten arabic text recognition with transformers*. *arXiv preprint arXiv:2410.02179*.
- Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. *TyDi QA: A benchmark*

- for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. pages 4171–4186, Minneapolis, Minnesota.
- Ishak Dolek and Atakan Kurt. 2025. **Ottoman htr: Recognition of the ottoman riqqa font using deep learning models**. *2025 Innovations in Intelligent Systems and Applications Conference (ASYU)*, pages 1–6.
- Mohammad Fasha, Bassam Hammo, Nadim Obeid, and Jabir Widian. 2020. **A hybrid deep learning model for arabic text recognition**. *CoRR*, abs/2009.01987.
- Heba Hassan, Ahmed El-Mahdy, and Mohamed E. Hussein. 2021. **Arabic scene text recognition in the deep learning era: Analysis on a novel dataset**. *IEEE Access*, 9:107046–107058.
- Ahmed Heakl, Sara Ghaboura, Omkar Thawakar, Fahad Shahbaz Khan, Hisham Cholakkal, Rao Muhammad Anwer, and Salman H. Khan. 2025a. **AIN: the arabic inclusive large multimodal model**. *CoRR*, abs/2502.00094.
- Ahmed Heakl, Muhammad Abdullah Sohail, Mukul Ranjan, Rania Elbadry, Ghazi Shazan Ahmad, Mohamed El-Geish, Omar Maher, Zhiqiang Shen, Fahad Shahbaz Khan, and Salman Khan. 2025b. **KITAB-bench: A comprehensive multi-domain benchmark for Arabic OCR and document understanding**. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 22006–22024, Vienna, Austria. 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 *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- JaidevAI. 2020. Easyocr. <https://github.com/JaidevAI/EasyOCR>. GitHub repository.
- Benjamin Kiessling. 2019. **Kraken - A Universal Text Recognizer for the Humanities**. In *Digital Humanities 2019*, Utrecht, Netherlands.
- Vu Trong Kim, Michael Krumdick, Varshini Reddy, Franck Dernoncourt, and Viet Dac Lai. 2024. **An analysis of multilingual FActScore**. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4309–4333, Miami, Florida, USA. Association for Computational Linguistics.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. **Gradient-based learning applied to document recognition**. *Proceedings of the IEEE*, 86(11):2278–2324.
- Vladimir I Levenshtein. 1966. Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics Doklady*, 10:707.
- Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020a. **MLQA: Evaluating cross-lingual extractive question answering**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7315–7330, Online. Association for Computational Linguistics.
- 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. 2020b. **Retrieval-augmented generation for knowledge-intensive NLP tasks**. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Minghao Li, Tengchao Lv, Jingye Chen, Lei Cui, Yijuan Lu, Dinei A. F. Florêncio, Cha Zhang, Zhoujun Li, and Furu Wei. 2023. **Trocr: Transformer-based optical character recognition with pre-trained models**. In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 13094–13102. AAAI Press.
- Xin Liu, Lechen Zhang, Sheza Munir, Yiyang Gu, and Lu Wang. 2025. **VeriFact: Enhancing long-form factuality evaluation with refined fact extraction and reference facts**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 17919–17936, Suzhou, China. Association for Computational Linguistics.
- Sabri A. Mahmoud, Irfan Ahmad, Wasfi G. Al-Khatib, Mohammad Alshayeb, Mohammad Tanvir Parvez, Volker Märgner, and Gernot A. Fink. 2014. **Khatt: An open arabic offline handwritten text database**. *Pattern Recognition*, 47(3):1096–1112. Handwriting Recognition and other PR Applications.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. **FActScore: Fine-grained atomic evaluation of factual precision in long form text generation**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.
- Hussein Mozannar, Elie Maamary, Karl El Hajal, and Hazem Hajj. 2019. **Neural Arabic question answering**. In *Proceedings of the Fourth Arabic Natural Language Processing Workshop*, pages 108–118, Florence, Italy. Association for Computational Linguistics.

Omer Nacar, Yasser Al-Habashi, Serry Sibae, Adel Ammar, and Wadii Boulila. 2025. [Sard: A large-scale synthetic arabic ocr dataset for book-style text recognition](#). *Preprint*, arXiv:2505.24600.

Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. [MS MARCO: A human generated machine reading comprehension dataset](#). In *Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016*, volume 1773 of *CEUR Workshop Proceedings*. CEUR-WS.org.

OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.

Vikas Paruchuri and Datalab Team. 2025. [Surya: A lightweight document ocr and analysis toolkit](#). <https://github.com/VikParuchuri/surya>. GitHub repository.

Mehreen Saeed, Adrian Chan, Anupam Mijar, Joseph Moukarzel, Georges Habchi, Carlos Younes, Amin Elias, Chau-Wai Wong, and Akram Khater. 2024a. [Muharaf: Manuscripts of handwritten arabic dataset for cursive text recognition](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

Mehreen Saeed, Adrian Chan, Anupam Mijar, Joseph Moukarzel, Georges Habchi, Carlos Younes, Amin Elias, Chau-Wai Wong, and Akram Khater. 2024b. [Muharaf: Manuscripts of handwritten arabic dataset for cursive text recognition](#).

Sheikh Shafayat, Eunsu Kim, Juhyun Oh, and Alice Oh. 2024. [Multi-fact: Assessing multilingual llms' multi-regional knowledge using factscore](#). *Preprint*, arXiv:2402.18045.

R. Smith. 2007. [An overview of the tesseract OCR engine](#). In *9th International Conference on Document Analysis and Recognition (ICDAR 2007)*, 23-26 September, Curitiba, Paraná, Brazil, pages 629–633. IEEE Computer Society.

Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, and 1118 others. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *Preprint*, arXiv:2403.05530.

Qwen Team. 2025. [Qwen2.5-vl](#).

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. [Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution](#). *Preprint*, arXiv:2409.12191.

Curtis Wigington, Chris Tensmeyer, Brian L. Davis, William A. Barrett, Brian L. Price, and Scott Cohen. 2018. [Start, follow, read: End-to-end full-page handwriting recognition](#). In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VI*, volume 11210 of *Lecture Notes in Computer Science*, pages 372–388. Springer.

Mokhtari Younes and Yousfi Abdellah. 2015. Segmentation of arabic handwritten text to lines. *Procedia Computer Science*, 73:115–121.

## A Prompts

The manuscript analysis prompt, presented in Figure 4, implements a structured reasoning framework that forces LLM-based analysis to be grounded to the provided retrieved evidence from the indexed manuscript. This two-stage protocol — first retrieving relevant manuscript passages via the `search_knowledge` tool before providing contextualized analysis — directly addresses the hallucination problem common in LLM applications in general.

For OCR benchmark evaluation, we employ two baseline prompts: Figure 5 provides generic "helpful assistant" system prompt and Figure 6 provides a simple OCR instruction.

Stage	Parameter	Value
Dominance-based Box Merging	dominance factor $\alpha_{\text{dominance}}$	2
	padding ratio $\gamma_p$	0.03
Overlap-based Filtering	max overlap threshold $\tau_{\text{overlap}}$	0.6
Image-based content validation	min ink density $\rho_{\text{min}}$	0.01
	max ink density $\rho_{\text{max}}$	0.70
	min contour count $c_{\text{min}}$	5
	min contour area $a_{\text{min}}$	4
	min hole ratio $\eta_{\text{min}}$	0.03
	min average defects $\delta_{\text{min}}$	0.3
Horizontal Row Aggregation	height difference ratio $\zeta_{\text{height}}$	0.5
	min vertical ratio $\kappa_{\text{overlap}}$	0.6

Table 4: Parameters for the four-stage post-processing pipeline applied to Kraken’s baseline segmenter output

## B Detailed Data Statistics

Table 5 summarizes the key features of MUHARAF, SARD, and MAS datasets used for fine-tuning OCR model. Overall, these three corpora cover a

```
"You are an expert OCR engine specialized in handwritten historical documents. Transcribe every character exactly as it appears --- preserving original spelling, punctuation, diacritics, ligatures, and archaic letters. Do not add, omit, normalize, correct, or format in any way. Output plain text only, matching the input one-to-one."
```

Figure 3: OCR Model Prompt

```
"You are an expert in analyzing ancient Arabic manuscripts from ancient Uzbekistan and Central Asia.
```

```
IMPORTANT: If a request seems to be about searching for information, use the search_knowledge tool first to search the manuscript database before providing any analysis. This tool contains extracted text from ancient manuscripts that you must reference.
```

```
When answering questions:
```

1. FIRST use search\_knowledge to find relevant information from the manuscripts
2. Then provide detailed analysis focusing on:
  - Historical and cultural context of ancient Uzbekistan
  - Religious and philosophical content (Islamic scholarship, Sufism)
  - Scientific and mathematical knowledge preservation
  - Trade routes and economic insights
  - Daily life and social customs
  - Literary and poetic elements
  - Paleographic and codicological observations when relevant

```
Always base your response on the actual manuscript content found through the search_knowledge tool.
```

```
If no relevant content is found, clearly state that and provide general historical context instead.
```

```
Always contextualize findings within the broader framework of Islamic civilization and Central Asian history."
```

Figure 4: Manuscript Analysis System Prompt

```
"You are a helpful assistant."
```

Figure 5: System Prompt for KITAB-Bench

```
"Extract the text in the image. Give me the final text, nothing else."
```

Figure 6: OCR Prompt for KITAB-Bench

wide range of document types, domains, fonts, and calligraphy ensuring the robustness of the resulting fine-tuned model. Kraken toolkit.

## C Hyperparameter Details

**Line Segmentation hyperparameters** Table 4 describes the hyperparameters used to post-process the initial line segmentation produced by the

Characteristic	Muharaf	SARD	MAS
Data Type	Historical handwritten manuscripts	Synthetic printed documents	Historical handwritten manuscripts
Total Images	1,644	843,622	1,023
Total Text Lines/Words	36,311 lines	690 million words	11,841 lines
Annotation Format	PAGE-XML, JSON	PAGE-XML	JSON
Text Source	Authentic historical documents	133,000+ unique articles from 9 domains	Authentic historical documents
Font Coverage	Predominantly Ruq'ah	10 fonts (Amiri, Arial, Calibri, Sakkal Majalla, Scheherazade New, Noto Naskh Arabic UI, Lateef, Thabit, Jozoor, Al-Jazeera-Arabic-Regular)	Naskh, Nastaliq, Taliq
Domain Coverage	Personal correspondence, diaries, poetry, church records, legal documents	Culture, Fatawa & Counsels, Literature & Language, Bibliography, Publications, Shariah, Social, Translations, News	waqf documents (testamentary acts of property donation to religious institutions), yarliqs (khan decrees and firmans), arizas (petitions), cheks (legal receipts and certificates), shajars (genealogical tables), and announcements (public proclamations and charters)
Period/Context	19th–21st centuries	Contemporary published texts	12th–19th centuries
Writing Styles	Informal handwriting, ranging from legible to barely readable	Clean printed fonts with controlled parameters	Handwriting, ranging from legible to barely readable
Document Types	Letters, diaries, notes, poems, religious records, contracts	Book layouts with full page markup	Letters, diaries, notes, religious records, books
Artifacts & Noise	Natural distortions: slant, curved lines, variable pen pressure	None (high-quality synthetic data)	slant, curved lines, variable pen pressure
Resolution/DPI	Original scanning resolution; lines aligned to 60-pixel height	300 DPI, A4 (8.27 × 11.69 inches), grayscale	Original scanning resolution
Primary Use Case	Handwritten text recognition (HTR) on cursive Arabic	Optical character recognition (OCR) on diverse typography	Handwritten text recognition and Knowledge discovery

Table 5: Comparison of Muharaf, SARD, MAS Datasets.

# promptolution: A Unified, Modular Framework for Prompt Optimization

Tom Zehle<sup>1,2,\*</sup>, Timo Heiß<sup>3,4,\*</sup>, Moritz Schlager<sup>4,5,\*</sup>,  
Matthias Aßemacher<sup>3,4</sup>, Matthias Feurer<sup>6,7</sup>

\*Equal contribution <sup>1</sup>ELLIS Institute, Tübingen, Germany <sup>2</sup> University of Freiburg, Germany

<sup>3</sup>LMU Munich, Germany <sup>4</sup>Munich Center for Machine Learning (MCML), Germany

<sup>5</sup>Technical University of Munich, Germany <sup>6</sup>TU Dortmund University, Germany

<sup>7</sup>Lamarr Institute for Machine Learning and Artificial Intelligence, Dortmund, Germany

Correspondence: tom.zehle@tue.ellis.eu, timo.heiss@stat.uni-muenchen.de, moritz.schlager@tum.de

## Abstract

Prompt optimization has become crucial for enhancing the performance of large language models (LLMs) across a broad range of tasks. Although many research papers demonstrate its effectiveness, practical adoption is hindered because existing implementations are often tied to unmaintained, isolated research codebases or require invasive integration into application frameworks. To address this, we introduce *promptolution*, a unified, modular open-source framework that provides all components required for prompt optimization within a single extensible system for both practitioners and researchers. It integrates multiple contemporary discrete prompt optimizers, supports systematic and reproducible benchmarking, and returns framework-agnostic prompt strings, enabling seamless integration into existing LLM pipelines while remaining agnostic to the underlying model implementation.

## 1 Introduction

Modern large language models (LLMs) exhibit impressive general-purpose capabilities, being able to solve a wide variety of tasks (Radford et al., 2019; Ouyang et al., 2022; Bai et al., 2023; Touvron et al., 2023). They can be adapted to solve specific tasks through in-context learning, i.e., simply by a textual instruction and optionally few-shot examples provided to the LLM as input (Brown et al., 2020). Since this input (referred to as *prompt*) steers the output of the LLM (Karmaker Santu and Feng, 2023; White et al., 2023), the LLM’s performance on a given task highly depends on it – in terms of quality, formulation, and the choice and order of examples (Zhao et al., 2021; Lu et al., 2022; Zhou et al., 2023). Table 1 illustrates this with two similar prompts for the GMS8K math dataset (Cobbe et al., 2021): despite their strong semantic similarity (matching parts in the same color), their performances differ substantially. This sensitivity highlights the potential of optimizing prompts for specific tasks, much like how hyperparameter tuning boosts performance in classical machine learning (Kohavi and John, 1995). Just as man-

-  [AutoML/Promptolution](#)
-  [System Demonstration](#)

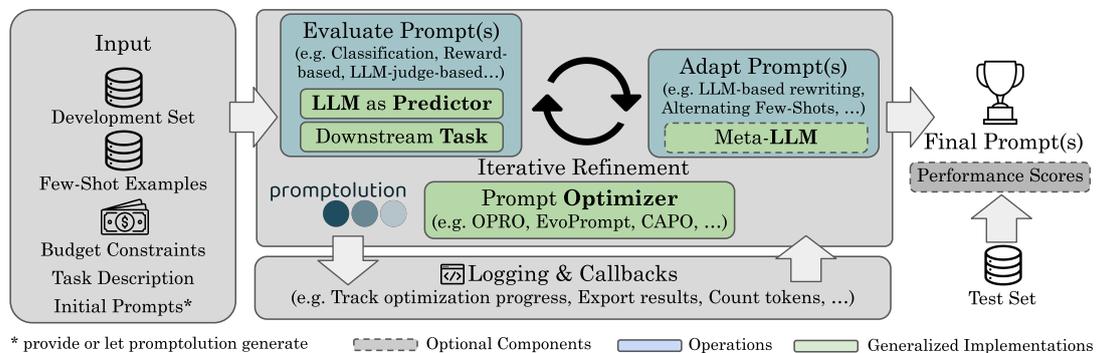


Figure 1: Overview of the *promptolution* framework. *promptolution* takes a dataset (dev set + few-shot examples), token budget constraints, a description of the task, and optionally initial prompts as input. In an iterative process, a user-selected prompt optimizer refines the prompt(s) by evaluating the LLM’s prediction performance on the task’s development set and adapting the prompts accordingly (e.g., through another LLM). Detailed logging and callbacks enable tracking the entire process. The optimized prompts are returned and can be evaluated on test data.

ual hyperparameter search, crafting and refining good prompts by hand (manual prompt engineering) is tedious and unreliable (Jiang et al., 2020; Liu et al., 2023). Similar to how AutoML automates the search over large hyperparameter spaces (Feurer and Hutter, 2019), automatic prompt optimization recently emerged to systematically search the combinatorial space of prompts (Li et al., 2025). Especially nowadays, when multi-agent systems have become central in both industry and research, LLMs specialized for individual tasks are increasingly important (He et al., 2025; Li et al., 2024). While fine-tuning entire models is expensive and data-intensive, automatic prompt optimization is a lightweight alternative and potential addition, often compatible with black-box LLMs (Cheng et al., 2024). A diverse collection of prompt optimizers has thus been introduced (see Cui et al., 2025; Li et al., 2025). However, several hurdles arise when applying these optimizers in practice. Each optimizer lives in its own research repository. Trying and comparing several optimizers requires juggling with multiple code bases and conflicting requirements. Additionally, these repositories are typically not actively maintained, and the research code often lacks proper software tests, documentation, and robustness. Moreover, their setups are inflexible, and deviations from their use cases (often limited to simple classification tasks) and LLM deployments require considerable effort. While existing libraries and tools for prompt optimization partly address these issues, they are either commercial and closed-source (Amazon Bedrock, 2025; Anthropic, 2025; Jina AI, 2025; Lee and Nardini, 2024), only implement a single optimizer (Adalflow, 2025; Agarwal et al., 2024; Hinthorn and Nishimi, 2025; Yuksekgonul et al., 2024), or have a high abstraction level designed mainly for end-to-end AI application development (Khattab et al., 2024; Kulin et al., 2025).

**Contribution.** We introduce **promptolution**, a modular, lightweight, and extensible open-source framework for automatic prompt optimization in Python, providing thoroughly tested and stable implementations. Our library implements multiple LLM interfaces, NLP tasks, and contemporary discrete prompt optimization methods, which can be used interchangeably or replaced with custom implementations of the components. While **promptolution** facilitates single-prompt optimization for practitioners, a low abstraction level and helpers for systematic, reproducible experiments

Prompt	Accuracy
Tackle this elementary math problem by breaking it into logical steps. When you reach the solution, enclose the final answer with <final_answer> and </final_answer> markers for clarity.	37.6%
Assist with solving the elementary or grade school level math problem that requires multiple steps and provide the solution within <final_answer> </final_answer> tags for easy identification.	53.8%

Table 1: Example prompts and their test set accuracy on GSM8K with Llama-3.3-70B (colors = similar phrases).

also make it geared toward researchers. Figure 1 shows an overview of how the framework operates.

**Outline.** We position our work within the landscape of automatic prompt optimization tools and libraries (§2), describe the design of our framework (§3), evaluate its performance compared to unoptimized prompts and other libraries to demonstrate its utility and competitiveness (§4), provide use cases and anti-use cases (§5), and outline future directions of the framework (§6).

## 2 Background & Related Works

**Prompt Optimization.** Automatic prompt optimization refers to systematically exploring prompt spaces using various automated optimization strategies, which may optimize both the instruction and the few-shot example components of a prompt (Cui et al., 2025; Li et al., 2025; Wan et al., 2024). Prompt optimization is commonly categorized by the nature of the prompt space into continuous (Li and Liang, 2021; Lester et al., 2021; Qin and Eisner, 2021) and discrete approaches (Guo et al., 2024; Yang et al., 2024; Zhou et al., 2023). For overviews of the field, we refer to Li et al. (2025) and Cui et al. (2025). **promptolution** focuses on the LLM-agnostic and interpretable discrete prompt optimization, which iteratively refines textual prompts directly, often using another “meta”-LLM<sup>1</sup> to generate improved candidates. The following optimizers are implemented in **promptolution**:

1. **OPRO** (Yang et al., 2024) uses LLMs as optimizers by providing a task description, examples, and previously scored candidates to the meta-LLM, which then proposes refined instructions.

<sup>1</sup>Can be the same as the one we optimize prompts for.

2. *EvoPrompt* (Guo et al., 2024) optimizes instructions using evolutionary algorithms, based on (a) a genetic algorithm (GA) and (b) on differential evolution (DE). In both cases, the meta-LLM performs crossover and mutation.
3. *CAPO* (Zehle et al., 2025) is a recent GA-based alternative that leverages AutoML techniques to improve cost-efficiency and jointly optimizes both instructions and few-shot examples, outperforming the discrete optimizers above.

Other promising prompt optimizers not yet implemented in `promptolution` include *GEPA* (Agrawal et al., 2025), *TextGrad* (Yuksekgonul et al., 2025), *MIPRO* (Opsahl-Ong et al., 2024), and *PromptWizard* (Agarwal et al., 2024).

**Existing Libraries, Frameworks & Tools.** Most optimizers above are implemented in siloed research repositories with hard-coded experimental setups. Oftentimes not actively maintained, lacking software tests and proper documentation, they are inherently difficult to use for both scientific benchmark experiments with other optimizers and practical use cases with specific requirements. In the prompt optimization landscape, many actively maintained libraries and tools have emerged, including both open- and closed-source solutions.

*Open-Source.* We classify open-source tools along four axes in Table 2: (1) single/multiple optimizers, (2) extensibility, (3) abstraction level, and (4) invasiveness of integration.<sup>2</sup> In the following, we focus on their most important delimitation criteria compared to `promptolution`. For explanations of the full categorization in Table 2, see Appendix A.1.

Several libraries only implement a single prompt optimizer, including *TextGrad* (Yuksekgonul et al., 2024) using the method from Yuksekgonul et al. (2025), *AdalFlow* (Adalflow, 2025) with *LLM-AutoDiff* (Yin and Wang, 2025), Microsoft’s *PromptWizard* (Agarwal et al., 2024) based on the eponymous optimization algorithm, and *PromptIM* (Hinthorn and Nishimi, 2025) with its own iterative optimization strategy. `prompt-ops` (Meta-Llama, 2025) implements two optimizers, but focuses on prompt optimization for Llama models. In contrast, `promptolution` offers multiple optimizers that can be used interchangeably (1) for arbitrary LLMs (2).

<sup>2</sup>(1) describes whether the tool implements only a single vs. multiple optimizers, (2) whether it can be easily extended to new NLP tasks, LLMs, and optimizers, (3) if the user has direct control over implementation details, and (4) if it is easy to integrate the optimized prompt into arbitrary existing LLM-pipelines.

Framework	Multiple Optimizers	Extensible	Low Abstraction	Non-Invasive Integration
<code>promptolution</code>	✓	✓	✓	✓
DSPy	✓	✓	✗	✗
CoolPrompt	✓	(✓)	✗	✗
promptomatix	✓	✗	✗	✗
prompt-ops	✓	✗	✓	✓
TextGrad	✗	✗	✓	✗
AdalFlow	✗	✗	✓	✗
PromptWizard	✗	✗	✓	✓
PromptIM	✗	✗	✓	✓

Table 2: Comparison of open-source frameworks.

DSPy (Khatab et al., 2024) is arguably the most popular existing framework in prompt optimization for building modular, declarative LLM pipelines and includes an embedded prompt optimization component. It supports multiple optimizers like *MIPROv2* (building on Opsahl-Ong et al., 2024) or *GEPA* (Agrawal et al., 2025), and can combine LLM training with prompt optimization (Soylu et al., 2024). While DSPy integrates prompt optimization as part of a monolithic, high-level program compilation procedure, `promptolution` exposes optimization as an explicit, iterative process, enabling finer control over optimization dynamics, intermediate results, and budget-aware stopping. Consequently, `promptolution` focuses exclusively on prompt optimization at a lower level of abstraction (3), making it geared toward researchers for systematic benchmarking and advanced practitioners rather than end-to-end AI application development. Moreover, DSPy only integrates with LLM applications in the DSPy framework while integration into other implementations requires larger refactoring (4). In contrast, `promptolution` returns a prompt string, enabling integration by directly replacing the existing prompt with the optimized one in any arbitrary LLM application (details in Appendix A.2). `promptomatix` (Murthy et al., 2025) builds on DSPy through its structured prompt compilation backend while also offering a lighter meta-prompt optimizer.

CoolPrompt (Kulin et al., 2025) is an LLM-agnostic framework that supports multiple optimizers and emphasizes “zero-configuration” (3) in contrast to `promptolution`, where researchers maintain close control over the setup.

Other related tools with a slightly different focus include *PromptBench* (Zhu et al., 2024), which targets LLM evaluation supporting only simple optimization techniques, and *OpenPrompt* (Ding et al., 2022), designed for prompt learning for language models predating modern LLMs.

*Closed-Source.* Proprietary tools range from commercial web-platforms like PromptPerfect (Jina AI, 2025), cloud integrations such as Google Cloud’s *Vertex AI Prompt Optimizer* (Lee and Nardini, 2024) and the *AWS Bedrock Prompt Engineering Playground* (Amazon Bedrock, 2025), to vendor-specific solutions like *Anthropic Claude Prompt Tools* (Anthropic, 2025).

**Positioning of promptolution.** Our library is an open-source, LLM-agnostic, and highly modular framework. It focuses exclusively on prompt optimization rather than constructing full LLM pipelines. It already includes relevant LLM interfaces and NLP tasks, along with evaluation metrics, and provides multiple contemporary discrete prompt optimizers in a single, unified system. The framework is highly customizable and extensible, offering fine-grained control over optimizers, tasks, logging, evaluation, and experiment configuration, and can be seamlessly integrated into existing LLM applications. As a result, it is suitable for both practitioners performing single-prompt optimization and for researchers conducting systematic, reproducible large-scale benchmark studies.

### 3 System Design

promptolution is designed as a modular and extensible framework consisting of four key components (see Figure 2). All components follow a unified interface defined through corresponding Base-classes, ensuring that implementations can be used interchangeably, remain fully compatible with the framework, and automatically inherit shared functionality. While each component can be configured individually (see §3.1–3.4), a separate `ExperimentConfig` together with associated

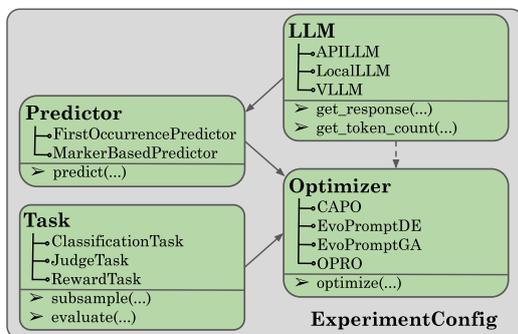


Figure 2: Core components of promptolution. The upper part of each box lists the component implementations, the lower part important functions. Arrows between components indicate conceptual connections.

---

To solve this problem, we need to calculate the total number of fish in the fishbowls at all the tables. First, [...] Then, we add the 3 fish from the table that has 3 fish: 62 fish + 3 fish = 65 fish. `<final_answer>65</final_answer>`

---

Table 3: Example of the MarkerBasedPredictor extraction from a LLM response for a GSM8K sample.

helper functions enables convenient parameterization of all components in a single object (see §3.5).

#### 3.1 LLM

The LLM component provides an interface for obtaining responses from any LLM implementation. Its base class enables parallelization and monitors token usage. Three classes are implemented:

1. **API LLM** enables calls to LLMs hosted via an API, covering common vendors such as OpenAI and Anthropic. Listing 1 (L.1–5) illustrates the setup through the DeepInfra API.
2. **Local LLM** allows using a local model via the transformers library (Wolf et al., 2020).
3. **VLLM** integrates the vllm library (Kwon et al., 2023) for efficient high-throughput inference and serving, and deployment on GPU clusters.

#### 3.2 Predictor

The Predictor component defines how predictions are extracted from LLM output. While the FirstOccurrencePredictor searches and extracts the first occurrence of any possible class label in the response, the more robust MarkerBasedPredictor (Listing 1, L.6) extracts between predefined HTML-like markers (see Table 3).

#### 3.3 Task

The Task component holds the dataset and other task-related information, including a textual description of the task, and defines how prompts are evaluated. It controls how subsampling is performed, which is crucial for efficiency, as not every prompt needs evaluation on the full data for a reasonable performance estimate.<sup>3</sup> We implement the following tasks relevant to NLP:

1. **ClassificationTask** targets discrete class labels and evaluates predictions using standard classification metrics. In our example in Listing 1 (L. 7–13), the task is specified using accuracy as metric. The task is easily created from a

<sup>3</sup>Depending on the subsampling strategy, evaluation is either performed on randomly drawn samples, a slice of the dataset (block), or the full dataset.

Pandas DataFrame by specifying the input and label columns, along with a task description.

2. **JudgeTask** is based on LLM-as-a-judge (Zheng et al., 2023), where another LLM scores the quality of the output according to the task description. This enables optimizing for subjective, creative tasks, both supervised and unsupervised, as ground-truth labels are optional.
3. **RewardTask** allows optimization w.r.t. custom reward functions. The reward can represent any measurable objective, such as code execution time or business metrics. Users define their own reward functions that compute a score (reward) directly from a predictor’s output.

### 3.4 Optimizer

The Optimizer component combines all other components by using a Predictor and a Task to determine the optimal prompt for the specified setup. Depending on the chosen optimization method, it may also rely on a (meta-) LLM for prompt alteration. As the core of the framework, it iteratively evaluates LLM predictions for a given task and refines the prompt(s) according to the respective optimization strategy.

promptolution currently implements four established prompt optimizers: OPRO (Yang et al., 2024), both EvoPromptDE and EvoPromptGA (Guo et al., 2024), and the current SOTA discrete prompt optimizer CAPO (Zehle et al., 2025). promptolution caches previously evaluated prompts out of the box, and constrains prompt evaluation during optimization to a subset of the available data, making algorithms faster and more efficient. In Listing 1, we set up CAPO (L. 14–22) and let it optimize prompts for 12 steps (L. 23).

Furthermore, new prompt optimizers can be added with minimal effort by inheriting from the base optimizer class and implementing a custom `_step()` method that defines the iterative optimization scheme. This makes our framework useful for researchers developing and benchmarking new optimization algorithms.

### 3.5 Experiment Configuration & Helper

promptolution not only supports optimizing prompts for a single specific setup but also provides an ExperimentConfig framework that enables a convenient, structured configuration for larger benchmark experiments. Additional helper functions enable the running and evaluation of such experiments with just a few lines of code.

```

1 llm = APILLM(
2     api_url="api.deepinfra.com/v1/...",
3     model_id="google/gemma-3-27b-it",
4     api_key="...",
5 )
6 predictor = MarkerBasedPredictor(llm=llm)
7 task = ClassificationTask(
8     df,
9     task_description="The task is...",
10    x_column="text",
11    y_column="label_text",
12    metric=accuracy_score,
13 )
14 optim = CAPO(
15     predictor,
16     task,
17     meta_llm=llm,
18     init_prompts=[
19         "Classify the text based on...",
20         # ...
21     ],
22 )
23 prompts = optim.optimize(n_steps=12)

```

Listing 1: Setup of promptolution’s core components.

Listing 2 (L. 1–7) illustrates how the previous example (Listing 1) can be expressed with a single config class. Arguments that are not specified resort to carefully chosen defaults, while arguments that were set but not used during initialization of the classes will throw a warning. The associated helper functions allow users to run the optimization process according to the config (`run_optimization`) and to evaluate the prompts on unseen test data (`run_evaluation`), without requiring manual initialization of the various classes affected. The `run_experiment` function (Listing 2, L. 8) combines both steps, to support researchers who want to perform extensive benchmark studies across multiple datasets and optimizers.

```

1 config = ExperimentConfig(
2     optimizer="capo",
3     task_description="The task is...",
4     n_steps=12,
5     api_url="api.deepinfra.com/v1/...",
6     model_id="google/gemma-3-27b-it",
7 )
8 prompts = run_experiment(df, config)

```

Listing 2: Running an experiment via the ExperimentConfig abstraction.

### 3.6 Supporting Modules and Utilities

*Exemplar Selection:* Some prompt optimization algorithms (e.g., OPRO or EvoPrompt) do not consider few-shot examples. However, they can substantially improve LLM performance (Brown et al.,

2020), even with simple selection strategies (Wan et al., 2024). promptolution offers post-hoc exemplar selection in an additional module, implementing random selection and random search<sup>4</sup> to add few-shot examples to a fixed instruction.

*Initial Prompt Creation:* Many prompt optimizers require an initial pool of prompts to start with. We offer functions to automatically create prompts from a task description, a base prompt (following Zhou et al., 2023), or samples from a dataset.

*Callbacks:* The base class for the optimizers supports callbacks, allowing for easy tracking of the optimization progress. Callbacks can access the state of the optimizer at every optimization step, and optionally terminate the process. Important implementations include the TokenCountCallback, which tracks the accumulated token budget and terminates optimization if a specified threshold is exceeded, and the FileOutputCallback, which writes prompts and their scores to a file, enabling easy post-hoc analysis of the process.

## 4 Evaluation

**Setup.** To evaluate promptolution and contextualize its performance relative to other prompt optimization tools, we perform prompt optimization on the popular *GSM8K* (grade school math word problems; Cobbe et al., 2021) and *SST-5* dataset (sentiment classification; Socher et al., 2013). We use gemma-3-27B instruction tuned (Kamath et al., 2025) as downstream LLM, and, for optimizers that require one, also as meta-LLM. Further details on datasets and implementation choices are provided in Appendix A.3.

For our comparison, we employ the optimizers *CAPO*, *EvoPromptGA*, and *OPRO* from the promptolution library, and additionally evaluate two other frameworks with leading prompt optimizers: AdalFlow (*LLM-AutoDiff*) and DSPy (*GEPA*). To assess the impact of prompt optimization itself, we also evaluate three unoptimized zero-shot prompts (see Appendix A.3) for each dataset and report their average performance. We use the default parameterization for each optimizer to ensure comparability with practical use cases, where users often lack the budget for an extensive hyperparameter search. We further restrict the token budget to at most one million in- and output tokens combined, which corresponds to a cost

<sup>4</sup>Random selection selects exemplars at random, whereas random search generates multiple sets of random examples, evaluates them, and selects the best performing set.

below \$0.15 for this LLM. All frameworks are initialized from a single task description, without any initial prompts, to test the full automation workflow. While DSPy and AdalFlow use the task description as a starting point for their respective optimizers, promptolution’s optimizers additionally require a set of initial prompts. These are generated via promptolution’s utility function `create_prompts_from_task_description`. Both the unoptimized prompts and best prompts per optimizer (based on development-set performance) are evaluated on an unseen test set.

For reproducibility, we make the experiment scripts, seed, and raw results publicly available at <https://github.com/finitearth/prompt-optimization-framework-comparison>.

**Results.** A summary of the results is presented in Table 4. With the exception of *GEPA* on *SST-5* and *OPRO* on *GSM8K*, all optimizers substantially outperform the unoptimized baseline with improvements of up to 15%p in accuracy (*CAPO* on *GSM8K*). This demonstrates the utility of prompt optimization in general and of our framework in particular. The best-performing optimizer on both datasets, *CAPO*, as well as each runner-up, is implemented in promptolution, underscoring the competitiveness of our framework compared to other prompt optimization tools. Although not every optimizer included in promptolution performs optimally on every task (e.g., *OPRO* on *GSM8K*), the library consistently provides at least one strong optimizer per task. Combined with its modular design, this allows users to switch easily to an alternative optimizer if the current one yields unsatisfying performance. We further emphasize that comparing optimizers within promptolution required minimal manual effort due to its dedicated support for systematic benchmark experiments.

For more in-depth evaluations of optimizers uti-

Framework	Optimizer	GSM8K	SST5
<i>Baseline</i>	<i>unoptimized</i>	78.1	44.6
AdalFlow	AutoDiff	88.7	55.7
DSPy	GEPA	84.7	42.0
promptolution	OPRO	69.7	<u>56.0</u>
	EvoPrompt	<u>91.0</u>	53.3
	CAPO	<b>93.7</b>	<b>56.3</b>

Table 4: Test set accuracy of optimized prompts using Gemma3-27B-it. Bold values indicate the best, underlined values the second-best performance per dataset.

lizing the promptolution framework, we point to [Zehle et al. \(2025\)](#).

## 5 Use Cases & Anti-Use Cases

The choice of prompt optimization frameworks depends not solely on performance, but also on *integration constraints*, the *scope of optimization*, and tolerance for *framework lock-in*. We outline four common use cases to clarify these trade-offs.

### **Case I: Integration in Existing Pipelines.**

*A practitioner already operates custom LLM pipelines and aims to improve performance solely via better prompts.* promptolution returns a single optimized prompt string that can be directly substituted into existing inference pipelines without refactoring application logic. Since optimized prompts are plain text, users retain full flexibility regarding LLM providers, deployment modes (API, local models, or vLLM), and future model updates. In contrast, adopting frameworks in which prompts are framework-specific abstractions, such as DSPy, requires rewriting pipelines into these specific structures, which entails significant overhead.

### **Case II: End-to-End LLM Application Development.**

*A practitioner builds an LLM-based system from scratch, jointly designing prompt logic, program structure, and possibly training or fine-tuning stages.* promptolution intentionally only optimizes prompts. In particular, DSPy is better suited to this setting, allowing the entire pipeline to be composed and optimized in a single, abstract program rather than manually engineering and optimizing each step; however, at the cost of coupling the system design to DSPy’s abstractions and losing portability to other frameworks.

### **Case III: Prompt Optimizer Benchmarking.**

*A researcher wants to systematically benchmark prompt optimization algorithms or develop a new prompt optimizer and perform a comparative evaluation.* promptolution is advantageous due to its low abstraction level, modular optimizer design, and explicit support for reproducible experiments. Extension to new optimizers, tasks, or evaluation strategies is explicitly encouraged and can be evaluated in a reproducible and controlled manner. Full trajectories, intermediate prompts, and token usage through callbacks and logging enable detailed analysis, setting promptolution apart from other libraries in this regard.

### **Case IV: Integration in LLM Benchmarking.**

*A researcher conducts systematic benchmarking of LLMs for a new use case. Instead of treating prompts as fixed, they want to include prompt optimization in the benchmark pipeline to account for the strong influence of prompts on performance and reduce variance in findings due to arbitrary prompt choices.* promptolution allows seamless integration of prompt optimization into any benchmark pipeline, exchanging LLM backends and datasets, and is designed for controlled, systematic, reproducible studies across multiple seeds.

## 6 Conclusion & Future Directions

In this work, we introduced promptolution, a unified and modular open-source Python framework for automatic prompt optimization, designed for both practitioners performing single-task prompt optimization and for researchers conducting systematic experiments. We highlighted the unique position of our framework within the landscape of prompt optimization tools, emphasizing its extensible, modular design and low abstraction level, and its focus on optimizing prompts with non-invasive integration. Furthermore, we demonstrated the role of each component in the framework and how they interact to support effective prompt optimization. Through a comparative evaluation, we verified the utility and competitiveness of promptolution. Finally, we provided concrete use and anti-use cases, highlighting the considerations to take into account when choosing a prompt optimization framework.

Looking ahead, we plan to develop interfaces with higher-level frameworks such as DSPy, enabling a combination of strengths from both ecosystems. To further simplify the management of complex, large-scale experimental setups, we intend to introduce an interface to configuration management frameworks such as hydra ([Yadan, 2019](#)). We also aim to improve accessibility by implementing a graphical interface for real-time experiment tracking and visual analysis of the optimization process. Following the recent rise of multi-agent systems, we plan to support the optimization of system prompts and interaction protocols across multiple agents. Additionally, inspired by AutoML, we intend to explore ensembling strategies for prompt optimizers, akin to ELPO ([Zhang et al., 2025](#)). The extensibility of our framework ensures that we can continue to incorporate new state-of-the-art optimization methods as the field evolves.

## Broader Impact

By unifying multiple prompt optimization methods that were previously scattered across separate research repositories, `promptolution` makes the benefits of prompt optimization for improving LLM performance broadly accessible, allowing practitioners to leverage these capabilities in real-world industry applications. At the same time, its modular and extensible design, combined with experiment-friendly implementations, makes the library a powerful tool for researchers benchmarking new prompt optimization algorithms. Since reimplementing competing methods and setting up rigorous benchmark experiments is typically time-consuming, `promptolution` offers the potential to accelerate research progress in prompt optimization and to enhance methodological comparability across the field.

## Acknowledgments

Tom Zehle received funding by the European Union. Views and opinions expressed are, however, those of the author(s) only and do not necessarily reflect those of the European Union or the European Commission. Neither the European Union nor the European Commission can be held responsible for them. This work was supported by the European Union’s Horizon Europe research and innovation program under grant agreement No 101214398 (ELLIOT).

Matthias Aßenmacher received funding from the BERD@NFDI consortium in the context of the work of the National Research Data Infrastructure (NFDI) Association. NFDI is funded by the Federal Republic of Germany and the 16 federal states. The BERD@NFDI consortium is supported within NFDI by the German Research Foundation (DFG) – NFDI 27/1-2026, project number 460037581.

## References

- Adalflow. 2025. [Build and Optimize LM Workflows](#). Last accessed: 11/30/2025.
- Eshaan Agarwal, Joykirat Singh, Vivek Dani, Raghav Magazine, Tanuja Ganu, and Akshay Nambi. 2024. PromptWizard: Task-aware prompt optimization framework. *arXiv:2405.18369 [cs.CL]*.
- Lakshya A. Agrawal, Shangyin Tan, Dilara Soyulu, Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J. Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. 2025. GEPA: Reflective prompt evolution can outperform reinforcement learning. *arXiv:2507.19457 [cs.CL]*.
- Amazon Bedrock. 2025. [User Guide: Optimize a prompt](#). Last accessed: 11/25/2025.
- Anthropic. 2025. [Prompt engineering: Use our prompt improver to optimize your prompts](#). Last accessed: 11/25/2025.
- 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, and 29 others. 2023. Qwen technical report. *arXiv:2309.16609 [cs.CL]*.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, and 12 others. 2020. Language models are few-shot learners. In *Proceedings of the 33rd International Conference on Advances in Neural Information Processing Systems (NeurIPS’20)*, pages 1877–1901. Curran Associates.
- Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. 2024. Black-box prompt optimization: Aligning large language models without model training. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3201–3219, Bangkok, Thailand. Association for Computational Linguistics.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv:2110.14168 [cs.LG]*.
- Wendi Cui, Jiaxin Zhang, Zhuohang Li, Hao Sun, Damien Lopez, Kamalika Das, Bradley A. Malin, and Sricharan Kumar. 2025. Heuristic-based search algorithm in automatic instruction-focused prompt optimization: A survey. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 22093–22111, Vienna, Austria. Association for Computational Linguistics.
- Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Haitao Zheng, and Maosong Sun. 2022. OpenPrompt: An open-source framework for prompt-learning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 105–113, Dublin, Ireland. Association for Computational Linguistics.
- Matthias Feurer and Frank Hutter. 2019. Hyperparameter optimization. In Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren, editors, *Automated Machine*

- Learning: Methods, Systems, Challenges*, pages 3–33. Springer International Publishing.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. 2024. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. In *The Twelfth International Conference on Learning Representations (ICLR'24)*. ICLR. Published online: [iclr.cc](https://iclr.cc).
- Junda He, Christoph Treude, and David Lo. 2025. LLM-based multi-agent systems for software engineering: Literature review, vision, and the road ahead. *ACM Transactions on Software Engineering and Methodology*, 34(5):1–30.
- William F. Hinthorn and Masahiro Nishimi. 2025. [Promptim](#). Last accessed: 11/25/2025.
- Zhengbao Jiang, Frank Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Jina AI. 2025. [PromptPerfect: AI Prompt Optimizer](#). Last accessed: 11/25/2025.
- Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, and 196 others. 2025. Gemma 3 technical report. *arXiv:2503.19786 [cs.CL]*.
- Shubhra Karmaker Santu and Dongji Feng. 2023. TELEr: A general taxonomy of LLM prompts for benchmarking complex tasks. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14197–14203. Association for Computational Linguistics.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2024. DSPy: Compiling declarative language model calls into state-of-the-art pipelines. In *The Twelfth International Conference on Learning Representations (ICLR'24)*. ICLR. Published online: [iclr.cc](https://iclr.cc).
- Ron Kohavi and George H. John. 1995. Automatic parameter selection by minimizing estimated error. In *Proceedings of the Twelfth International Conference on Machine Learning (ICML'95)*. Morgan Kaufmann Publishers.
- Nikita Kulin, Viktor Zhuravlev, Artur Khairullin, Alena Sitkina, and Sergey Muravyov. 2025. CoolPrompt: Automatic prompt optimization framework for Large Language Models. In *Proceedings of the 38th Conference of Open Innovations Association FRUCT, Issue 1 (Full papers)*, pages 158–166, Helsinki, Finland. FRUCT Oy.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles (SOSP '23)*, pages 611–626. Association for Computing Machinery.
- George Lee and Ivan Nardini. 2024. [Announcing Public Preview of Vertex AI Prompt Optimizer](#). Last accessed: 11/25/2025.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3045–3059. Association for Computational Linguistics.
- Wenwu Li, Xiangfeng Wang, Wenhao Li, and Bo Jin. 2025. A survey of automatic prompt engineering: An optimization perspective. *arXiv:2502.11560 [cs.AI]*.
- Xiang 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. Association for Computational Linguistics.
- Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. 2024. A survey on LLM-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinity*, 1(1):9.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):195:1–195:35.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098. Association for Computational Linguistics.
- Meta-Llama. 2025. [prompt-ops: An open-source tool for LLM prompt optimization](#). Last accessed: 11/30/2025.
- Rithesh Murthy, Ming Zhu, Liangwei Yang, Jieli Qiu, Juntao Tan, Shelby Heinecke, Caiming Xiong, Silvio Savarese, and Huan Wang. 2025. Promptomatix: An automatic prompt optimization framework for Large Language Models. *arXiv:2507.14241 [cs.CL]*.
- Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. 2024. Optimizing instructions and demonstrations for multi-stage language model programs.

- In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9340–9366. Association for Computational Linguistics.
- 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 the 35th International Conference on Advances in Neural Information Processing Systems (NeurIPS’22)*. Curran Associates.
- Guanghui Qin and Jason Eisner. 2021. Learning how to ask: Querying LMs with mixtures of soft prompts. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5203–5212. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642. Association for Computational Linguistics.
- Dilara Soylu, Christopher Potts, and Omar Khattab. 2024. Fine-tuning and prompt optimization: Two great steps that work better together. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 10696–10710. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and efficient foundation language models. *arXiv:2302.13971 [cs.CL]*.
- Xingchen Wan, Ruoxi Sun, Hootan Nakhost, and Sercan Arik. 2024. Teach better or show smarter? on instructions and exemplars in automatic prompt optimization. In *Proceedings of the 37th International Conference on Advances in Neural Information Processing Systems (NeurIPS’24)*, pages 58174–58244. Curran Associates.
- Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with ChatGPT. In *Proceedings of the 30th Conference on Pattern Languages of Programs, PLOP ’23*, pages 1–31, USA. The Hillside Group.
- 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, and 3 others. 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.
- Omry Yadan. 2019. [Hydra - a framework for elegantly configuring complex applications](#). Last accessed: 11/30/2025.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc Le, Denny Zhou, and Xinyun Chen. 2024. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations (ICLR’24)*. ICLR. Published online: [iclr.cc](#).
- Li Yin and Zhangyang Wang. 2025. LLM-AutoDiff: Auto-differentiate any LLM workflow. *arXiv:2501.16673 [cs.CL]*.
- Mert Yuksekogun, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. 2024. [TextGrad: Automatic “Differentiation” with Text](#). Last accessed: 11/30/2025.
- Mert Yuksekogun, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. 2025. Optimizing generative AI by backpropagating language model feedback. *Nature*, 639(8055):609–616.
- Tom Zehle, Moritz Schlager, Timo Heiß, and Matthias Feurer. 2025. CAPO: Cost-aware prompt optimization. In *Proceedings of the Fourth International Conference on Automated Machine Learning*, volume 293 of *Proceedings of Machine Learning Research*, pages 18/1–45. PMLR.
- Qing Zhang, Bing Xu, Xudong Zhang, Yifan Shi, Yang Li, Chen Zhang, Yik Chung Wu, Ngai Wong, Yijie Chen, Hong Dai, Xiansen Chen, and Mian Zhang. 2025. ELPO: Ensemble learning based prompt optimization for Large Language Models. *arXiv:2511.16122 [cs.CL]*.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *Proceedings of the 38th International Conference on Machine Learning (ICML’21)*, volume 139 of *Proceedings of Machine Learning Research*, pages 12697–12706. PMLR.

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. 2023. Judging LLM-as-a-judge with MT-bench and Chatbot Arena. In *Proceedings of the 36th International Conference on Advances in Neural Information Processing Systems (NeurIPS'23)*, pages 46595–46623. Curran Associates.

Y. Zhou, A. Ioan Muresanu, Z. Han, K. Paster, S. Pitis, H. Chan, and J. Ba. 2023. Large language models are human-level prompt engineers. In *The Eleventh International Conference on Learning Representations (ICLR'23)*. ICLR. Published online: [iclr.cc](https://iclr.cc).

Kaijie Zhu, Qinlin Zhao, Hao Chen, Jindong Wang, and Xing Xie. 2024. PromptBench: A unified library for evaluation of Large Language Models. *Journal of Machine Learning Research*, 25(254):1–22.

## A Appendix

### A.1 Details on Prompt Optimization Tool Categorization

In the following, we explain the choices (✓/✗) in Table 2 for all considered open-source frameworks:

**DSPy** (Khattab et al., 2024): multiple optimizers (e.g., GEPA, MIPROv2), extensible since adding new optimizers is possible, high abstraction level through declarative programming and monolithic compile method, invasive integration as requiring switching to `dspy.Module` and `dspy.Signature` (details in Appendix A.2).

**CoolPrompt** (Kulin et al., 2025): multiple optimizers (e.g. DistillPrompt, ReflectivePrompt), partly extensible since implementing own optimizers is possible, but follows no clear framework, high abstraction level as it emphasizes “zero-configuration”, non-invasive integration since it simply returns a prompt string.

**promptomatix** (Murthy et al., 2025): uses DSPy optimizers as backend, not extensible by design since it only has one fixed PromptOptimizer class, high abstraction since it builds on top of DSPy, invasive since it wraps DSPy, therefore requiring at least the same effort to integrate.

**prompt-ops** (Meta-Llama, 2025): supports a basic and advanced optimizer, not extensible as the advanced optimizer is intended exclusively for Llama models, low abstraction since we can parametrize all components through a single configuration file, non-invasive since it simply returns a prompt string.

**TextGrad** (Yuksekgonul et al., 2024): single optimizer (TextGrad by Yuksekgonul et al. (2025)) and thus not extensible, low abstraction as the optimizer is parametrizable in a fine-grained way, invasive integration as requiring prompts to be wrapped in `textgrad.Variable` and LLM calls to `textgrad.BlackboxLLM` and engine wrappers.

**AdalFlow** (Adalflow, 2025): single optimizer (LLM-AutoDiff by Yin and Wang (2025)) and thus not extensible, low abstraction as the optimizer is parametrizable in a fine-grained way, invasive integration as it requires re-architecting pipelines into PyTorch-like component classes.

**PromptWizard** (Agarwal et al., 2024): single optimizer (PromptWizard) and thus not extensible, low abstraction as the optimizer is parametrizable in a fine-grained way, non-invasive since it simply returns a prompt string.

**PromptLM** (Hinthorn and Nishimi, 2025): single optimizer (own optimization strategy) and thus not extensible, low abstraction as their optimizer is parametrizable in a fine-grained way, non-invasive since it can either be configured for local experimentation, just returning a prompt string, or to automatically commit the optimized prompt to the LangChain Hub.

### A.2 Delimitation from DSPy

One might argue that the existing framework DSPy (Khattab et al., 2024) is already very similar to `promptolution`, as it is also open-source and modular, supports prompt optimization, implements multiple optimizers, supports several LLM implementations, and provides additional utilities. However, there are several important differences, both in underlying purpose and in concrete design, that clearly differentiate the two frameworks:

**Different focus** : DSPy includes prompt optimization only as one component within a broader compilation workflow tightly coupled to a program structure, whereas `promptolution` is not a full application framework and instead focuses exclusively on the prompt optimization stage.

**Abstraction level** : DSPy hides much of the prompt construction behind a high-level declarative interface (though it remains inspectable). In particular, DSPy primarily offers a monolithic compile method that runs the entire optimization routine. In contrast, `promptolution` provides an iterative optimize routine that repeatedly invokes an optimizer-specific `_step` method, along with associated callbacks. This design enables fine-grained control, including intermediate prompt candidates, scores, and token usage. It further facilitates early stopping, custom logging for full transparency, and budget-aware termination through callbacks.

**Target user base** : `promptolution` is geared toward researchers and advanced ML practitioners, while DSPy primarily targets AI application developers building end-to-end work-

flows such as agents or RAG systems.

**Experiments** : `promptolution` makes it easy to implement custom optimizers, tasks, and other components, and is particularly suited for systematic large-scale benchmark experiments due to the integrated config framework. Its extensibility explicitly encourages contributions of new algorithms and tasks, which is not a focus of DSPy.

**Integration & Portability** : In DSPy, prompt optimization is performed on tasks defined via `dspy.Modules` (e.g., `Predict`, `Refine`, `ChainOfThought`, ...) that require `dspy.Signatures` as input. This couples optimized prompts to DSPy’s program abstraction, which introduces a degree of lock-in when individual subtasks are later migrated to a different system, as additional adaptation may be required. Conversely, `promptolution` produces a standalone prompt string, allowing optimized prompts to be reused or exchanged across systems with minimal friction, thereby improving portability.

### A.3 Experiment Details

In §4, we evaluate our framework against DSPy and AdalFlow. Since neither provides a straightforward way to restrict compute budgets based on token counts, we enforce token limits by throwing an exception in the respective LLM wrappers when the limit is exceeded, and then returning the last suggested prompt.

For evaluation, we route every LLM API call through the same interface using `langchain` to ensure a fair comparison of the optimized prompts, and compare exact matches between the predicted and true labels (allowing for differences in capitalization). In the case of *GEPA*, we had to manually clean the LLM outputs because they did not follow the required output format stated in the task description (encapsulating the prediction within `<final_answer>` tags). Specifically, we had to remove the `“[[ ## target ## ]]”` and `“[[ ## completed ## ]]”` tags. For AdalFlow, we had to intercept API calls due to faulty extraction of system and user prompts. The system-prompt extraction mechanism relied on tags, but these were not forwarded correctly (e.g., `<START_OF_USER_PROMPT>` was expected instead of the provided `<START_OF_USER>`). These changes are documented in the experiment repository. The used LLM was hosted on local servers

Dataset	Huggingface ID	n <sub>dev</sub>	n <sub>test</sub>
SST5	SeetFit/sst5	500	300
GSM8K	openai/gsm8k	500	300

Table 5: Overview of the utilized HuggingFace datasets.

and accessed via API calls. The utilized datasets and sample sizes are detailed in Table 5. The dev set is used for optimization, the test set for holdout evaluation of the final prompts. Few-shot examples, if considered by the optimizer, are also taken from the dev set.

The automatically generated prompts used to compare against a zero-effort baseline are shown in Table 6. The system prompt accompanying the unoptimized prompts and those optimized by `promptolution` was “You are a helpful assistant!”; the system prompts used for the other frameworks were the ones returned after optimization. Note that prompts resulting from `promptolution` generally do not include any `“{input}”`-tags (unlike DSPy and AdalFlow, as well as some of the unoptimized prompts) to enable query-independent prompt-caching. Instead, `promptolution` appends respective queries to the end of the prompt.

The complete experiment code and raw results are publicly available for full reproducibility at <https://github.com/finitearth/prompt-optimization-framework-comparison>.

### A.4 Quality Standards

To ensure high code quality, we adopt established software engineering best practices throughout our package. We maintain a comprehensive test suite that automatically verifies expected behavior after code changes. All tests must pass before a release is published, ensuring users encounter no issues when updating. The main branch is protected, and all contributions must be submitted via pull requests, each reviewed by another main contributor. We additionally employ pre-commit hooks to automatically check code formatting, documentation, and other basic quality issues before commits are made, improving readability and maintainability. We also maintain strict documentation standards and do not accept poorly documented pull requests. Dedicated CI and CI/CD pipelines enable an automated build, test, and release of the package and documentation.

Task	Prompt
GSM8K	Solve the maths problem step by step. Give your answer inside <final_answer> ...</final_answer>.\n\n{input}
	Calculate the solution to the problem below. Show your steps.\n\n{input}\n\nRequired output format: <final_answer> ANSWER </final_answer>.
	Please analyze the following mathematical problem. Break down your reasoning into logical steps before stating the solution.\n\nProblem: {input}\n\nFormat your conclusion as <final_answer>YOUR_ANSWER </final_answer>.
SST5	Classify the sentiment as very negative, negative, neutral, positive, or very positive. Use <final_answer>-tags: \n {input}
	Identify the sentiment: very negative, negative, neutral, positive, very positive.\n\nText: {input}\n\nAnswer: <final_answer>label </final_answer>.
	Text: {input}\n\nClassify as very negative, negative, neutral, positive, or very positive. Output: <final_answer>label </final_answer>.

Table 6: Unoptimized zero-shot prompts per dataset.

## A.5 Documentation & Tutorials

Alongside the open-source software package, we provide extensive documentation available at <https://automl.github.io/promptlution/>. It covers the major components and their functionality, and also includes tutorials that guide users through their first prompt optimization use case. This further enhances the accessibility and usability of our framework.

## A.6 Extensibility of promptlution

promptlution is designed to be easily extensible, allowing researchers to integrate new optimization algorithms with minimal implementation effort. To illustrate this, we briefly outline how a new prompt optimizer can be added to the framework.

All prompt optimizers inherit from the abstract base class `BaseOptimizer`, which defines the common interface and execution structure shared across optimization algorithms. Concretely, extending the framework requires implementing only two abstract methods: `_pre_optimization_loop`, which performs any setup before the optimization begins, and `_step`, which executes a single opti-

mization iteration and returns the updated prompt candidate(s). The overall optimization loop, including configuration handling, callback invocation, and other features, is already implemented in the base class via the `optimize` method.

This design ensures that developers of new optimizers can focus exclusively on their algorithmic logic, while benefiting from shared infrastructure such as logging, callbacks, and budget-aware termination. A similar extension pattern applies when adding new task-, LLM-, or predictor components.

## A.7 System Demonstration

The system demonstration under <https://youtu.be/gySdgjEhsZA> uses the following code example, through which we guide step-by-step. It requires installing the `promptlution` package with API support enabled. The subsequent imports establish the necessary components for the LLM wrapper, CAPO, and the evaluation task definitions.

```
pip install promptlution[api]
```

```
from promptlution.llms import APILLM
from promptlution.optimizers import CAPO
from promptlution.tasks import JudgeTask
from promptlution.predictors import
↳ MarkerBasedPredictor
from promptlution.utils import
↳ create_prompts_from_task_description
import pandas as pd
```

The dataset, containing raw email-generation instructions, is loaded from a CSV file using pandas:

```
df_emails = pd.read_csv("emails.csv")
```

After loading the data, we define a task description that specifies the desired persona, tone, and formatting constraints (including the closing signature). This serves as the reference specification for the entire workflow.

```
task_description = """Write concise, polite,
↳ professional academic emails for me as a
↳ PhD student, asking clarifying questions
↳ when my instructions are vague, avoiding
↳ cliché openings, and always ending with
↳ Yours sincerely, Tom Zehle."""
```

We then instantiate the underlying LLM. In this example, the `APILLM` wrapper connects to the `Llama-3.3-70B-Instruct-Turbo` model via the `DeepInfra` API, which serves as the downstream LLM, the judge, and the meta-LLM.

```
llm = APILLM(
    model_id="meta-llama/Llama-3.3-70B-\
        \"Instruct-Turbo\",
    api_url="https://api.deepinfra.com/\
        \"v1/openai\",
    api_key=open("token.txt").read(),
)
```

To initialize the optimization search space, a set of candidate prompts is derived directly from the task description using the framework's utility function. These serve as the starting population for the optimizer we use later.

```
initial_prompts =
    create_prompts_from_task_description(
        task_description,
        llm,
    )
```

Given the subjective nature of the email generation task (lacking a ground truth), a JudgeTask is configured. This setup utilizes an LLM-as-a-Judge approach to evaluate the semantic quality of the generated outputs against the input instructions. An alternative could be the ClassificationTask when ground truth labels are available, or the RewardTask when the user can define an objective function to score the outputs.

```
task = JudgeTask(
    df_emails,
    judge_llm=llm,
    task_description=task_description,
    x_column="instruction",
)
```

Next, we instantiate the CAPO optimizer with a marker-based predictor. The optimization routine is executed for six iterations (`n_steps=6`) to refine the prompt candidates iteratively.

```
optimizer = CAPO(
    task=task,
    predictor=MarkerBasedPredictor(llm),
    meta_llm=llm,
    initial_prompts=initial_prompts,
    check_fs_accuracy=False
)
final_prompts = optimizer.optimize(n_steps=6)
```

After optimization completes, the final optimized prompts are returned and ready for use in email generation.

# T-pro 2.0: An Efficient Russian Hybrid-Reasoning Model and Playground

Gen-T Team  
T-Tech

Correspondence: [anatolii.s.potapov@gmail.com](mailto:anatolii.s.potapov@gmail.com)

## Abstract

We introduce *T-pro 2.0*, an open-weight Russian LLM for hybrid reasoning and efficient inference. The model supports direct answering and reasoning-trace generation, using a Cyrillic-dense tokenizer and an adapted EAGLE speculative-decoding pipeline to reduce latency. To enable reproducible and extensible research, we release the model weights, the *T-Wix* 500k instruction corpus, the *T-Math* reasoning benchmark, and the EAGLE weights on Hugging Face. These resources allow users to study Russian-language reasoning and to extend or adapt both the model and the inference pipeline. A public web demo exposes reasoning and non-reasoning modes and illustrates the speedups achieved by our inference stack across domains. T-pro 2.0 thus serves as an accessible open system for building and evaluating efficient, practical Russian LLM applications.

Demo: <https://t-pro2eagle.streamlit.app/>

 [hf.co/collections/t-tech/t-pro-20](https://hf.co/collections/t-tech/t-pro-20)

## 1 Introduction

Large Language Models (LLMs) have progressed from basic text generation to systems capable of multi-step reasoning and efficient inference. Recent foundation models show that reasoning-oriented training (DeepSeek-AI et al., 2025a; Yang et al., 2025) and improved decoding methods (Chen et al., 2023; Li et al., 2024e) can substantially boost both accuracy and speed.

In the Russian open-source space, progress remains limited. Most strong models are closed and accessible only through APIs (Mamedov et al., 2025; Zmitrovich et al., 2023), while open models are typically small adaptations of multilingual systems (Nikolich et al., 2024). There is no unified ecosystem for studying Russian-language reasoning: high-quality evaluation sets are scarce, and, to the best of our knowledge, there are currently few

public demos that let users compare direct answering and step-by-step reasoning, inspect inference-time optimizations, or observe how decoding speed impacts user experience.

To address these gaps, we introduce *T-pro 2.0*, an open-weight Russian LLM for hybrid reasoning and an interactive demo platform. The model supports two complementary modes—direct answering and explicit reasoning traces—enabling applications to balance speed and accuracy within a single deployed system.

Our training setup combines a Cyrillic-dense tokenizer derived from Qwen3 (Yang et al., 2025), large-scale instructional midtraining, supervised fine-tuning focused on both reasoning and non-reasoning, preference optimization, and an adaptation of EAGLE-style speculative decoding (Li et al., 2024e) to accelerate Russian-language inference. To sum up, our main contributions are:

- *T-pro 2.0*, an open-weight Russian hybrid-reasoning LLM with improved inference efficiency via an optimized Cyrillic tokenizer and EAGLE-style speculative decoding.
- *T-Wix*, the largest open Russian hybrid-reasoning SFT dataset to date ( $\approx 500k$  samples) covering general instruction following, long-context tasks, and teacher-generated reasoning traces.
- *T-Math*, a benchmark of Russian high-school olympiad-level mathematics problems for curriculum-aligned reasoning evaluation.
- An interactive web demo that exposes *T-pro 2.0* as a research-oriented live system<sup>1</sup>, enabling side-by-side comparison of reasoning and non-reasoning modes, running tasks from our datasets and benchmarks, and viewing telemetry for inference-time optimizations.

<sup>1</sup>The [web demo video](#) is available on YouTube.

All model-related components (T-pro 2.0, EAGLE weights, and the T-Math benchmark) are released under the Apache-2.0 license, while the T-Wix corpus is released under the ODC-By open data license.

## 2 Related Work

The development of Russian LLMs primarily follows two tracks: monolingual pre-training and adaptation of multilingual models. Early decoder-only baselines like ruGPT (Kuratov and Arkhipov, 2019; Zmitrovich et al., 2023) and commercial systems such as YandexGPT<sup>2</sup> and GigaChat (Mamedov et al., 2025) focus on Russian-centric pre-training. While achieving promising results on Russian benchmarks, early versions face a capability gap compared to leading multilingual LLMs like Qwen (Yang et al., 2024) and Llama (Dubey et al., 2024).

To mitigate these limitations, *T-pro 1.0*<sup>3</sup> adopts a continued pre-training strategy on large-scale Russian corpora, reaching state-of-the-art results on MERA (Fenogenova et al., 2024) among open Russian models. Its release aligns with a broader shift toward strengthening open-source Russian LLMs, alongside projects such as Saiga (Gusev, 2023), RuAdapt (Tikhomirov and Chernyshev, 2024), and Vikhr (Nikolich et al., 2024). These works emphasize the value of mitigating English-centric tokenizer limitations (Petrov et al., 2024) and extending pre-training on Russian data. This direction continues to grow: although YandexGPT-5-Lite<sup>4</sup> is a fully pre-trained model rather than an adaptation, its recent open release further expands the set of publicly available Russian foundation models.

**Russian Instruction Datasets.** Existing Russian instruction datasets vary in provenance and domain coverage. Saiga (Gusev, 2023) applies self-instruct (Wang et al., 2023) pipelines producing `ru_turbo_saiga`, GrandMaster-PRO-MAX (Nikolich et al., 2024) aggregates sources across coding and general knowledge, and RuAdapt (Tikhomirov and Chernyshev, 2024) combines translated and native Russian samples. However, these datasets are usually small and contain few reasoning-intensive tasks.

<sup>2</sup><https://ya.ru/ai/gpt>

<sup>3</sup><https://huggingface.co/t-tech/T-pro-it-1.0>

<sup>4</sup><https://huggingface.co/yandex/YandexGPT-5-Lite-8B-pretrain>

**Efficient Inference.** Speculative decoding accelerates autoregressive inference (Leviathan et al., 2023a). EAGLE (Li et al., 2024d) uses a lightweight head to generate draft token trees verified in parallel, achieving 2–3× speedup. Multi-Token Prediction (MTP) (Gloeckle et al., 2024) trains models to predict multiple tokens simultaneously and is deployed successfully in DeepSeek-V3 (DeepSeek-AI et al., 2025b). GigaChat models (Mamedov et al., 2025) also adopt MoE architecture for increased efficiency on training and inference stages. Speculative decoding remains underexplored for general-purpose Russian LLMs, with few publicly documented deployments.

## 3 T-pro 2.0

### 3.1 System and Demonstration Description

We provide a public web demo of T-pro 2.0 that exposes the model as an interactive hybrid-reasoning assistant and makes our inference optimizations directly observable. The service is stateless and does not store user prompts or completions. The interface supports multi-turn chat in Russian and English and side-by-side comparison with baseline models (by default Qwen3-32B-Instruct), allowing users to inspect both answers and reasoning traces under identical serving conditions. The demo currently supports text-only interactions and does not perform additional server-side content filtering beyond what is built into the underlying models.

**Architecture.** The demo is a single-page web application backed by a lightweight Python HTTP server. The server exposes a simple JSON API, attaches configuration options (model, decoding mode, generation parameters) received from the UI, and forwards requests to two serverless SGLang endpoints (Gu et al., 2024). Each endpoint runs on a single NVIDIA H100 GPU: one hosts T-pro 2.0 with an EAGLE-style speculative decoding pipeline (draft head + 32B verifier), and the other hosts the Qwen3-32B baseline with standard autoregressive decoding. The deployment is tuned for interactive use and supports around 20 concurrent users per model while keeping per-request latency low.

**User interface and functionality.** Figure 1 shows the main layout. The central comparison view presents parallel completions from two systems. For each side, users can independently choose between standard and reasoning modes.

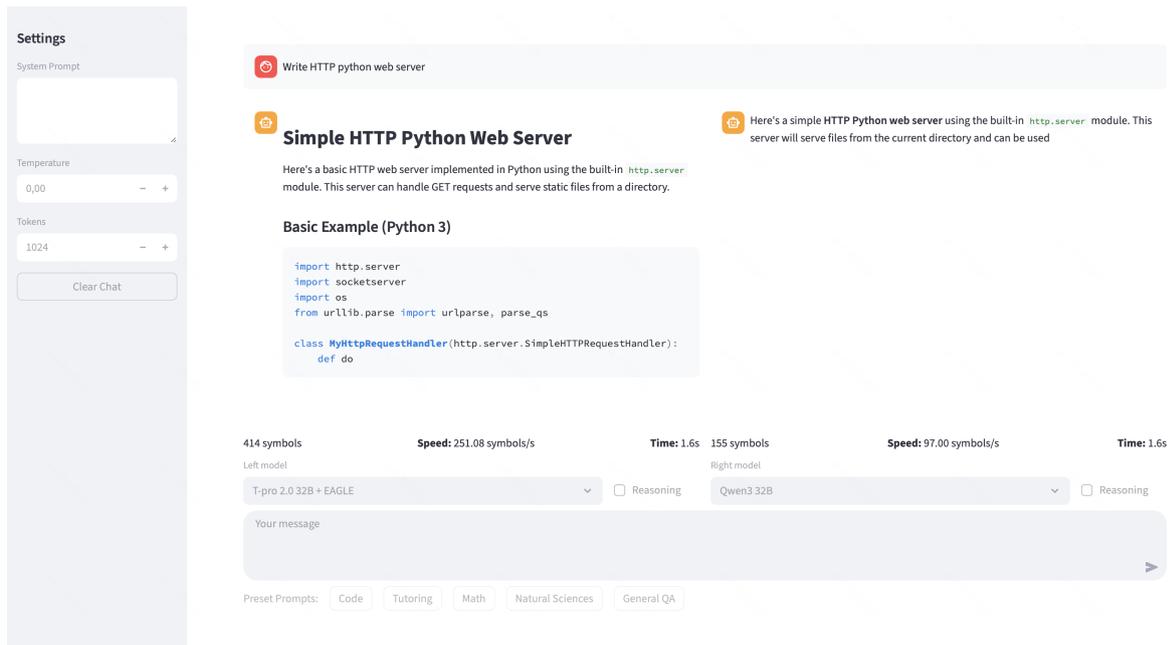


Figure 1: Screenshot of the system demo of the T-pro 2.0 EAGLE.

Outputs are streamed token by token, making differences in latency, verbosity, and reasoning structure directly visible. A control panel above the chat area lets users select models, toggle reasoning per model, and adjust decoding parameters such as temperature, maximum length, and sampling options. All decoding settings used for a given interaction are displayed in the UI, so that qualitative comparisons can be reproduced outside the demo. A typical interaction consists of selecting a preset prompt (or entering a custom query), choosing reasoning and generation settings, and launching both models to compare their outputs and telemetry.

To support systematic probing, the interface provides a small library of predefined prompts grouped by domain (general questions, math and science, code, etc.). Several presets are derived from our evaluation suites, including T-Math and other Russian reasoning benchmarks, so that users can quickly examine T-pro 2.0 on challenging tasks without reconstructing benchmark-style prompts.

**Performance telemetry and usage patterns.** A telemetry panel reports, for every request and model, the number of generated tokens, end-to-end latency, streaming throughput in tokens per second, and the acceptance ratio of speculative tokens for T-pro 2.0. Relating these statistics to the visible outputs illustrates how speculative decoding affects both accuracy and perceived responsiveness for short conversational turns and long reasoning

traces, complementing the benchmark results in Section 4.

### 3.2 Training recipe

This section describes the T-pro 2.0 training pipeline, integrating tokenizer adaptation, instructional midtraining, general post-training, and EAGLE-based speculative decoding. At all stages, we perform MinHash deduplication against benchmarks to prevent data leakage.

**Cyrillic-dense tokenizer** We address the systematic under-tokenization of Russian in multilingual models by replacing 34k low-frequency non-Cyrillic tokens in the Qwen3 (Yang et al., 2025) vocabulary with Cyrillic ones while keeping the total size fixed.

To build the expansion set, we extract 35.7k candidate tokens containing at least one Cyrillic character from four donor tokenizers (Qwen3, RuAdapt (Tikhomirov and Chernyshev, 2024), cl100k\_base (OpenAI et al., 2024), MGPT (Shlitzhko et al., 2023)). For each candidate, we evaluate its decomposition under the current merge graph and iteratively add those merges required to make two-piece decompositions fully reachable. Four refinement passes make approximately 95% of candidates reachable. Tokens containing Cyrillic, pure Latin tokens, punctuation, and all 1–2-symbol units are preserved, while the 34k removed tokens are selected via log-smoothed frequency

scoring on the midtraining data mixture to maintain tokenization quality across all domains (see Appendix C, Table 11, for the domain distribution). Our modification yields substantial compression

T-pro	Qwen3	GigaChat <sup>†</sup>	Ruadapt-Qwen3
<b>2.71</b>	3.63	2.89	3.26

Table 1: Average tokens per word on Wikipedia for eight Cyrillic languages (ru, uk, be, bg, sr, mk, kk, ky). Lower values indicate more efficient segmentation of Cyrillic text. <sup>†</sup>Indicates <https://huggingface.co/ai-sage/GigaChat-20B-A3B-instruct> model.

gains: on Russian Wikipedia, the share of Russian words tokenized into at most two tokens rises from 38% to 60% (see Table 7), and Tables 1, 9 further demonstrates that this improvement generalizes consistently across eight Cyrillic languages, with all of them exhibiting shorter average segmentations under our tokenizer. A full set of tokenizer evaluation metrics is provided in Appendix B.

**Instructional midtraining** To adapt the *Qwen3-32B* model to the new dense Russian tokenizer and enhance reasoning, we employ an intermediate stage on 40B tokens drawn from curated open-source instructions, synthetic tasks, and parallel corpora. The mixture is dominated by Russian (49%) and English (36%) text; in terms of domains, it is dominated by Reasoning (34.6%), General QA (28.8%), and Math (16.2%), supplemented by grounded synthetic Question Answering (QA), code, and real user dialogues. The data mixture undergoes rigorous domain-specific local sensitive hashing (LSH) deduplication and InsTag-based semantic deduplication (Lu et al., 2024; Abbas et al., 2023) to balance diversity. To ensure high quality and stylistic consistency, all assistant responses are regenerated using *Qwen3-235B-A22B* teacher. Training utilizes a 32k context window, stabilizing the model for downstream supervised fine-tuning (SFT). Ablations show that the instruct-only mid-training outperforms mixtures retaining the fraction of raw pre-training data on reasoning tasks, improving ruAIME 2024 (T-Tech, 2025e) from 0.60 to 0.67. Separately, 8B-scale experiments confirm the tokenizer transition, with the T-pro tokenizer reaching a higher MERA (Fenogenova et al., 2024) macro-average (0.574) than the original *Qwen3* tokenizer (0.560). Full details and ablations are provided in Appendix C.

**Reward Model (RM) construction** To support the T-pro 2.0 post-training pipeline, a dedicated reward model is trained (see Appendix F). The RM is initialized from *Qwen3-32B* with a scalar regression head and trained with a Bradley–Terry preference objective on sequences up to 32K tokens using Ulysses-style sequence parallelism. Synthetic preference data are generated using knockdown tournaments over completions from multiple instruct- and reasoning-oriented model groups of different scales, substantially reducing the number of pairwise evaluations relative to an exhaustive pairwise scheme. For each instruction, completion pairs are judged, pairs with positional bias are discarded, and transitive tournament relations are added to improve preference coverage. To assess downstream performance, we design an Arena-Hard Best-of- $N$  benchmark based on the  $\Delta_{\text{BoN}}$  (best@ $N$  – worst@ $N$ ) metric, on which our RM outperforms existing open-source reward models; full details are provided in Appendix F.1.

**General Post-Training** The T-pro 2.0 post-training pipeline is implemented through general and reasoning SFT, and on-policy Direct Preference Optimization (DPO), with all filtering procedures detailed in Appendices D-E.

For the general part of the T-Wix SFT dataset, approximately 14M raw instructions from open-source corpora are reduced to 468k samples using deduplication, a multi-stage filtering pipeline, and domain/complexity balancing across six domains (Math, Code, Science, General Instruction, General Knowledge, Writing) and three difficulty tiers (School, Student, Professor). Each instruction is expanded with 8 candidate completions generated by *Qwen3-235B-A22B* or *DeepSeek-V3* (DeepSeek-AI et al., 2025b) and then passed through an RM-guided selection step. The resulting mixture is low-noise, domain-balanced, and predominantly Russian, with approximately 10% English data retained to preserve bilingual competence.

For the reasoning component, approximately 30K samples are drawn from a 450k English pool covering general reasoning, mathematics, natural sciences, and code. After translation and deduplication, candidate solutions are generated by the *Qwen3-235B-A22B* teacher model and a midtraining student checkpoint and then filtered via RM-based rejection. For verifiable tasks, the highest-scoring factually correct teacher output is selected; for open-ended tasks, the shortest valid

trace among the teacher’s top RM-ranked candidates is chosen.

DPO is performed on 100k instructions sampled from the T-Wix dataset, with a 90/10 general-to-reasoning ratio. For each instruction, 16 on-policy completions are RM-scored, and one high-contrast preference pair (best vs. worst) is formed, so that observed failure modes are directly targeted and alignment is improved without the overhead of on-line RL.

**Speculative Decoding** We integrated a lossless EAGLE-based speculative decoding module into T-pro 2.0, where a lightweight draft model proposes candidate tokens that are verified by the frozen 32B target model (Leviathan et al., 2023b). The draft model consists of a single Llama-2-based decoder layer with an FR-Spec component (Zhao et al., 2025), trained on hidden-state reconstruction (smoothed  $L_1$ ) and token distribution alignment (KL divergence) losses. During inference, we employ EAGLE-2’s dynamic draft-tree mechanism via SGLang. As shown in Tables 2 and 3, at temperature 0.8 the module achieves an average speedup of  $1.85\times$  in standard mode, showing similar speedups for both standard and reasoning modes. STEM domains consistently outperform humanities ( $1.99\times$  vs  $1.62\times$ ), due to more predictable token distributions in technical content. See Appendix G for training pipeline details.

Benchmark	Speedup		Acceptance Length	
	Standard	Reasoning	Standard	Reasoning
ruMT-Bench <sup>1</sup>	1.79	1.69	3.31	3.10
ruAlpaca <sup>2</sup>	1.61	1.57	2.94	2.85
ruCodeEval <sup>3</sup>	2.15	1.84	3.93	3.34
T-Math <sup>4</sup>	–	2.25	–	4.01
<b>Average</b>	<b>1.85</b>	<b>1.83</b>	<b>3.39</b>	<b>3.33</b>

Table 2: Aggregated T-pro-2.0-eagle performance at temperature 0.8. Comparison of relative speedup and average acceptance length for the standard and reasoning modes. <sup>1</sup>T-Tech (2025c), <sup>2</sup>T-Tech (2025a), <sup>3</sup>Fenogenova et al. (2024), <sup>4</sup>T-Tech (2025g)

Domain Category	Speedup	Acceptance Length
STEM & Business <sup>†</sup>	1.99	3.57
Social & Humanities <sup>‡</sup>	1.62	2.88
<b>Average</b>	<b>1.79</b>	<b>3.19</b>

Table 3: Aggregated acceleration results on ruMMLU-Pro. <sup>†</sup>Includes Math, Chem, Eng, Bus, Phys, CS. <sup>‡</sup>Includes Econ, Bio, Health, Psych, Phil, Hist, Law.

## 4 Evaluation

### 4.1 Benchmarks

We evaluate along three axes: factual knowledge, dialogue, and reasoning capabilities.

#### Russian common-knowledge benchmarks

MERA, MaMuRAMu (Fenogenova et al., 2024), and ruMMLU-Pro (T-Tech, 2025f), targeting world knowledge and logical competence.

#### Dialogue benchmarks

Arena Hard Ru (T-Tech, 2025b), Arena Hard 2 (Li et al., 2024b,c), and WildChat Hard Ru (Kukushkin, 2024) (curated native Russian queries). WildChat uses o3-mini responses as baseline; DeepSeek-V3.1-Terminus (DeepSeek-AI et al., 2025b) serves as judge across all arenas and DeepSeek-V3.1 (DeepSeek-AI et al., 2025b) for WildChat.

#### Reasoning benchmarks

AIME 24/25 (Zhang and Math-AI, 2024) (Zhang and Math-AI, 2025), MATH-500 (Hendrycks et al., 2021), GPQA Diamond (Rein et al., 2023), Vikhr Math/Physics (Kuleshov et al., 2025), and LiveCodeBench v4\_v5 (Jain et al., 2024). English benchmarks are professionally localized (ruAIME, ruMATH-500, ruGPQA, ruLCB (T-Tech, 2025d)). Vikhr benchmarks use Math-Verify scoring <sup>5</sup>.

#### T-Math benchmark

We introduce T-Math—331 problems from All-Russian and Moscow olympiads (1998–2025), automatically extracted and human-verified. Details are provided in Appendix H.

### 4.2 Results

#### General knowledge and dialogue abilities

Table 4 shows the results on Russian general-knowledge and dialogue benchmarks. T-pro 2.0 performs consistently well across all evaluations, scoring 0.66 on MERA and 0.697 on ruMMLU-Pro. These numbers put it close to GPT-4o (0.714) and above Russian-adapted baselines such as RuadaptQwen3-32B-Instruct (0.652).

On dialogue tasks, the model reaches 91.1 on Arena Hard Ru and 72.6 on WildChat Hard Ru, outperforming all open-source systems and most proprietary ones. On Arena Hard 2, T-pro 2.0 scores 53.5 on Hard Prompts and 64.8 on Creative Writing, showing that it reliably follows instructions across different task types. These results directly reflect

<sup>5</sup><https://github.com/huggingface/Math-Verify>

Model	MERA	MaMuRAMu	ruMMLU-Pro	Arena Hard Ru	WildChat Hard Ru	Arena Hard 2 HP	Arena Hard 2 CW
<i>Open Source Models (27B-32B class)</i>							
<b>T-pro 2.0 (Ours)</b>	<b>0.66</b>	<b>0.851</b>	<b>0.697</b>	<b>91.1 / 90.36</b>	<b>72.6 / 76.4</b>	<b>53.5 / 46.2</b>	<b>64.2 / 62.8</b>
Qwen3-32B	0.582	0.833	0.677	83.95 / 84.66	59.6 / 51.9	46.4 / 32.9	53.7 / 41.5
RuadaptQwen3-32B-Instruct <sup>1</sup>	0.574	0.823	0.652	68.4 / 64.76	41.5 / 39.4	13 / 14.2	19.4 / 12.7
Gemma 3 27B <sup>2</sup>	0.577	0.808	0.665	82.66	58.4	23.5	74.7
DeepSeek-R1-Distill-Qwen-32B <sup>3</sup>	0.508	0.787	0.537	34.07 / 22.83	12.1 / 8.7	7.2 / 7.2	5.9 / 3.5
<i>Open Source Larger Scale &amp; Proprietary Models</i>							
DeepSeek-V3	<b>0.677</b>	<b>0.875</b>	<b>0.736</b>	92.67	66.8	45.8	59.9
DeepSeek-R1 <sup>3</sup>	–	–	–	<b>95.74</b>	<b>90.3</b>	<b>73.6</b>	<b>90.3</b>
YandexGPT5-Pro <sup>4</sup>	–	–	0.604	19.13	12.1	3.8	2.6
GigaChat 2 Max <sup>5</sup>	0.67	0.864	0.649	61.44	10.1	8.5	27.1
o4-mini <sup>6</sup> (medium)	–	–	–	<u>95.63</u>	<u>74.4</u>	<u>67</u>	49.8
GPT-4o <sup>7</sup>	<u>0.642</u>	<u>0.874</u>	<u>0.714</u>	85.14	41.4	20.0	44.2

<sup>1</sup>Tikhomirov and Chernyshev (2024), <sup>2</sup>Team et al. (2025), <sup>3</sup>DeepSeek-AI et al. (2025a), <sup>4</sup><https://ya.ru/ai/gpt>, <sup>5</sup>Mamedov et al. (2025), <sup>6</sup>OpenAI (2025), <sup>7</sup>OpenAI et al. (2024).

Table 4: Comparison of models on Russian language understanding and dialogue benchmarks. In Arena Hard 2, subsets are Hard Prompt (HP) and Creative Writing (CW). For entries reported as  $a/b$ , the first value corresponds to the *reasoning* setting and the second to the *non-reasoning* setting. o4-mini and DeepSeek-R1 are omitted from MERA as this benchmark does not support reasoning model mode, while YandexGPT5-Pro is omitted from MERA due to licensing restrictions.

Model	T-Math	ruAIME 2024	ruAIME 2025	ruMATH-500	ruGPQA Diamond	ruLCB	Vikhr Math	Vikhr Physics
<i>Open Source Models (27B-32B class)</i>								
<b>T-pro 2.0 (Ours)</b>	<b>0.541</b>	<b>0.704</b>	<b>0.646</b>	<b>0.94</b>	0.591	<b>0.563</b>	0.799	0.51
Qwen3-32B	0.529	<b>0.706</b>	0.625	0.938	0.606	0.537	<b>0.809</b>	<b>0.531</b>
RuadaptQwen3-32B-Instruct	0.444	0.575	0.450	0.450	0.591	0.500	0.528	0.337
Gemma 3 27B	0.208	0.248	0.231	0.860	0.439	0.261	0.548	0.276
DeepSeek-R1-Distill-Qwen-32B	0.254	0.510	0.402	0.898	<b>0.631</b>	0.493	0.462	0.286
<i>Open Source Larger Scale &amp; Proprietary Models</i>								
DeepSeek-V3	0.278	0.319	0.285	0.882	0.657	0.444	0.613	0.367
DeepSeek-R1	0.619	<b>0.800</b>	<b>0.800</b>	<b>0.972</b>	0.763	0.69	<b>0.864</b>	<b>0.469</b>
YandexGPT5-Pro	0.13	0.062	0.046	0.682	0.354	0.265	0.372	0.252
GigaChat 2 Max	0.142	0.102	0.062	0.702	0.475	0.272	0.372	0.245
o4-mini (medium)	<b>0.634</b>	0.781	0.771	0.958	<b>0.773</b>	<b>0.705</b>	0.834	0.408
GPT-4o	0.106	0.090	0.069	0.766	0.510	0.131	0.372	0.296

Table 5: Comparison of models on Russian advanced reasoning benchmarks.

the structure of the T-Wix corpus, which mixes general instruction-following with long-context tasks and distilled reasoning traces from stronger teacher models.

**Reasoning capabilities** Table 5 summarizes performance on T-Math and localized reasoning benchmarks. On T-Math, the model scores 0.541, indicating strong performance on original olympiad-style Russian problems. On ruAIME 2024 and 2025 it reaches 0.704 and 0.646, sharply outperforming DeepSeek-V3 (0.319/0.285), GPT-4o (0.090/0.069) and all proprietary Russian models. Results on ruMATH-500 (0.94) and Vikhr Math (0.799) further confirm the model’s ability to perform mathematical reasoning in Russian under varied setups. These results also show that T-Math is a chal-

lenging benchmark that reveals meaningful performance differences that are obscured by translated or adaptation-based alternatives.

Crucially, the Russian-focused training does not hurt English performance. As shown in Table 23, T-pro 2.0 achieves 0.765 on AIME 2024, 0.966 on MATH-500, and 0.556 on LiveCodeBench, closely matching Qwen3-32B (0.808, 0.961, 0.546) and consistently outperforming other adapted models in the same class. These results indicate that the Cyrillic-focused tokenizer and Russian-centric training pipeline preserve robust English reasoning with minimal degradation. A detailed breakdown is provided in Appendix I.

## 5 Conclusion

We present *T-pro 2.0*, an open-weight Russian language model tailored for hybrid reasoning and efficient inference. The combination of a Cyrillic-dense tokenizer, a reasoning-focused midtraining stage, and an adapted EAGLE-style speculative decoding pipeline allows the model to deliver strong performance on Russian tasks without increasing model size and without notable degradation on English benchmarks. Along with the model, we release *T-Wix*, a large-scale SFT dataset enriched with reasoning traces, and *T-Math*, a benchmark designed to probe mathematical and analytical abilities in Russian.

These results point to two broader takeaways. First, careful, targeted adaptation of strong multilingual backbones remains a practical and reproducible route for building high-quality models for languages with limited resources. Second, tokenizer design and inference optimization are not optional details but key components for deploying reasoning-capable models beyond English. We expect the released models, datasets, evaluation code, and public demo to support research on Russian-language reasoning LLMs and to contribute to more transparent and consistent evaluation practices in this area.

### Ethical Statement

**Possible Misuse.** Our work may enable generation of misleading, offensive, or otherwise harmful content if deployed without appropriate safeguards. We do not support applications that restrict access to information, facilitate disinformation, target individuals or groups, or automate harmful actions. To mitigate these risks, we apply toxicity and safety filtering during post-training and provide usage guidelines for responsible deployment.

**Biases and Data Quality.** The datasets used for pre-training and fine-tuning contain publicly available Russian and English text, which may include stereotypes, factual inaccuracies, or cultural biases. While automated filtering and manual checks are applied, residual biases may remain. We recommend additional evaluation when transferring the model to domains or communities that are under-represented in the training data.

**Human Subjects and Privacy.** This work does not involve human-subjects research or the collection of personally identifiable information. All data

sources comply with their respective licenses and usage policies.

### Limitations

Despite the strong performance of *T-pro 2.0*, several limitations should be acknowledged, which are planned to be addressed in future work.

**Limited Multilingual Evaluation.** The training recipe was solely focused on Russian with English data replay to preserve and transfer abilities. As a result, downstream performance was evaluated only on these two languages, and generation quality in other languages remains untested and may vary. Additionally, since the reasoning ability was trained via distillation without explicit language control, the language of reasoning traces may be occasionally different. While such cases are infrequent in practice, developing methods for explicit language control during reasoning, broader multilingual support and evaluation are left for future work.

**Limited Agentic Capabilities** No dedicated improvements for tool use or agentic behavior were incorporated into the model. Optimizations for function calling or complex multi-turn interactions were not performed, and as a result, performance in these areas is expected to be comparable to or slightly below that of the base Qwen3-32B model. Enhancements in this direction are prioritized for future development.

**Offline-Only Reinforcement Learning** Model alignment was conducted exclusively through offline DPO. Online reinforcement learning methods such as PPO (Schulman et al., 2017) or GRPO (Shao et al., 2024) were not employed. Although DPO is computationally efficient, the absence of interactive feedback may limit robustness and lead to performance degradation on out-of-domain tasks.

**Unverified Long-Context Performance** All training stages for *T-pro 2.0* were carried out with a fixed context window of 32k tokens, consistent with the base Qwen3-32B configuration. While support for up to 128k tokens is theoretically enabled via RoPE scaling, the model’s ability to maintain coherence and retrieve information over such extended contexts has not been empirically validated.

**Reproducibility Issues** Full reproducibility is restricted by the use of proprietary datasets in mid-training and the DPO stage. However, the curated

SFT dataset is being released to support and encourage further research, particularly in the development of high-quality Russian-language language models.

## Author Contributions

- *Project Lead and Technical Direction*: Anatolii Potapov
- *Training Pipelines Team*: German Abramov, Pavel Gein
- *Post-training Team*: Dmitrii Stoianov, Olga Tsymboi, Danil Taranets, Ramil Latypov, Almaz Dautov, Dmitry Abulkhanov
- *Inference Team*: Vladislav Kruglikov, Nikita Surkov
- *Evaluation Team*: Mikhail Gashkov, Viktor Zelenkovskiy, Artem Batalov, Aleksandr Medvedev

## References

2025. Turbo-alignment. <https://github.com/turbo-llm/turbo-alignment>. GitHub repository.
- Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S. Morcos. 2023. *Semdedup: Data-efficient learning at web-scale through semantic deduplication*. *arXiv preprint arXiv:2303.09540*.
- Berk Atil, Sarp Aykent, Alexa Chittams, Lisheng Fu, Rebecca J. Passonneau, Evan Radcliffe, Guru Rajan Rajagopal, Adam Sloan, Tomasz Tudrej, Ferhan Ture, Zhe Wu, Lixinyu Xu, and Breck Baldwin. 2025. *Non-determinism of "deterministic" llm settings*. *Preprint*, arXiv:2408.04667.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. *Qwen2.5-vl technical report*. *Preprint*, arXiv:2502.13923.
- Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, Ido Shahaf, Oren Tropp, Ehud Karpas, Ran Zilberstein, Jiaqi Zeng, Soumye Singhal, Alexander Bukharin, Yian Zhang, Tugrul Konuk, and 114 others. 2025. *Llama-nemotron: Efficient reasoning models*. *Preprint*, arXiv:2505.00949.
- RALPH ALLAN BRADLEY and MILTON E. TERRY. 1952. *Rank analysis of incomplete block designs: The method of paired comparisons*. *Biometrika*, 39(3-4):324–345.
- Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. *Accelerating large language model decoding with speculative sampling*. *Preprint*, arXiv:2302.01318.
- Yang Chen, Zhuolin Yang, Zihan Liu, Chankyu Lee, Peng Xu, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2025. *Acereason-nemotron: Advancing math and code reasoning through reinforcement learning*. *arXiv preprint arXiv:2505.16400*.
- Peng Cui and Mrinmaya Sachan. 2025. *Investigating the zone of proximal development of language models for in-context learning*. *Preprint*, arXiv:2502.06990.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025a. *Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning*. *Preprint*, arXiv:2501.12948.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2025b. *Deepseek-v3 technical report*. *Preprint*, arXiv:2412.19437.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, and 1 others. 2024. *The Llama 3 herd of models*. *arXiv preprint arXiv:2407.21783*.
- Alena Fenogenova, Artem Chervyakov, Nikita Martynov, Anastasia Kozlova, Maria Tikhonova, Albina Akhmetgareeva, Anton Emelyanov, Denis Shevelev, Pavel Lebedev, Leonid Sinev, Ulyana Isaeva, Katerina Kolomeytseva, Daniil Moskovskiy, Elizaveta Goncharova, Nikita Savushkin, Polina Mikhailova, Anastasia Minaeva, Denis Dimitrov, Alexander Panchenko, and Sergey Markov. 2024. *Mera: A comprehensive llm evaluation in russian*. *arXiv preprint arXiv:2401.04531*.
- Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, and 1 others. 2024. *Better & faster large language models via multi-token prediction*. *arXiv preprint arXiv:2404.19737*.
- Shiyang Gu and 1 others. 2024. *Sglang: Efficient serving of llms with speculative decoding and continuous batching*. GitHub repository.
- Igor Gusev. 2023. *Saiga: Instruction-tuned russian llama models*. <https://huggingface.co/IlyaGusev/saiga>.
- Dan Hendrycks and 1 others. 2021. *Measuring mathematical problem solving with the math dataset*. *arXiv preprint arXiv:2103.03874*.

- Hugging Face. 2025. [Open r1: A fully open reproduction of deepseek-r1](#).
- Shawn Im and Sharon Li. 2025. [Can dpo learn diverse human values? a theoretical scaling law](#). *Preprint*, arXiv:2408.03459.
- Sam Ade Jacobs, Masahiro Tanaka, Chengming Zhang, Minjia Zhang, Reza Yazdani Aminadabi, Shuaiwen Leon Song, Samyam Rajbhandari, and Yuxiong He. 2024. [System optimizations for enabling training of extreme long sequence transformer models](#). In *Proceedings of the 43rd ACM Symposium on Principles of Distributed Computing, PODC '24*, page 121–130, New York, NY, USA. Association for Computing Machinery.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. [Livecodebench: Holistic and contamination free evaluation of large language models for code](#). *Preprint*, arXiv:2403.07974.
- Diederik P. Kingma and Jimmy Ba. 2017. [Adam: A method for stochastic optimization](#). *Preprint*, arXiv:1412.6980.
- Alexander Kukushkin. 2024. [wildchat-hard-ru](https://github.com/kuk/wildchat-hard-ru). <https://github.com/kuk/wildchat-hard-ru>. GitHub repository.
- Ilya Kuleshov, Pavel Ilin, Nikolay Kompanets, Ksenia Sycheva, and Aleksandr Nikolich. 2025. [Doom: Difficult olympiads of math](#). ArXiv preprint.
- Yuri Kuratov and Alexey Arkipov. 2019. [Adaptation of deep bidirectional multilingual transformers for russian language](#). *arXiv preprint arXiv:1912.11283*.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023a. [Fast inference from transformers via speculative decoding](#). In *Proceedings of the 40th International Conference on Machine Learning (ICML '23)*, pages 19274–19286.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023b. [Fast inference from transformers via speculative decoding](#). *arXiv preprint arXiv:2211.17192*.
- Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024a. [From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning](#). *Preprint*, arXiv:2308.12032.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. 2024b. [From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline](#). *arXiv preprint arXiv:2406.11939*.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. 2024c. [From live data to high-quality benchmarks: The arena-hard pipeline](#).
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024d. [Eagle-2: Faster inference of language models with dynamic draft trees](#).
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024e. [Eagle: Speculative sampling requires rethinking feature uncertainty](#). *arXiv preprint arXiv:2401.15077*.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2025. [Eagle-3: Scaling up inference acceleration of large language models via training-time test](#). *arXiv preprint arXiv:2503.01840*.
- Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. 2024. [Skywork-reward: Bag of tricks for reward modeling in llms](#). *arXiv preprint arXiv:2410.18451*.
- Chris Yuhao Liu, Liang Zeng, Yuzhen Xiao, Jujie He, Jiakai Liu, Chaojie Wang, Rui Yan, Wei Shen, Fuxiang Zhang, Jiacheng Xu, Yang Liu, and Yahui Zhou. 2025a. [Skywork-reward-v2: Scaling preference data curation via human-ai synergy](#). *arXiv preprint arXiv:2507.01352*.
- Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. 2025b. [Pairjudge rm: Perform best-of-n sampling with knockout tournament](#). *arXiv preprint arXiv:2501.13007*. In progress work.
- Zihe Liu, Jiashun Liu, Yancheng He, Weixun Wang, Jiaheng Liu, Ling Pan, Xinyu Hu, Shaopan Xiong, Ju Huang, Jian Hu, Shengyi Huang, Johan Obando-Ceron, Siran Yang, Jiamang Wang, Wenbo Su, and Bo Zheng. 2025c. [Part i: Tricks or traps? a deep dive into rl for llm reasoning](#). *Preprint*, arXiv:2508.08221.
- Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and Jingren Zhou. 2024. [#instag: Instruction tagging for analyzing supervised fine-tuning of large language models](#). In *The Twelfth International Conference on Learning Representations*.
- Saumya Malik, Valentina Pyatkin, Sander Land, Jacob Morrison, Noah A. Smith, Hannaneh Hajishirzi, and Nathan Lambert. 2025. [Rewardbench 2: Advancing reward model evaluation](#). *Preprint*, arXiv:2506.01937.
- Valentin Mamedov, Evgenii Kosarev, Gregory Leyletner, Ilya Shchuckin, Valeriy Berezovskiy, Daniil Smirnov, Dmitry Kozlov, Sergei Averkiev, Lukyanenko Ivan, Aleksandr Proshunin, Ainur Israfilova, Ivan Baskov, Artem Chervyakov, Emil Shakirov, Mikhail Kolesov, Daria Khomich, Daria Latortseva, Sergei Porkhun, Yury Fedorov, and 14 others. 2025. [GigaChat family: Efficient Russian language modeling through mixture of experts architecture](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 93–106, Vienna, Austria. Association for Computational Linguistics.

- Aleksandr Nikolich, Konstantin Korolev, Sergei Bratchikov, Igor Kiselev, and Artem Shelmanov. 2024. Vikhr: The family of open-source instruction-tuned large language models for russian. *arXiv preprint arXiv:2405.13929*.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, and 401 others. 2024. *Gpt-4o system card*. *Preprint*, arXiv:2410.21276.
- OpenAI. 2025. Introducing o3 and o4-mini. <https://openai.com/index/o3-o4-mini-system-card/>. Accessed: 2025-05-15.
- Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Biber. 2024. Language model tokenizers introduce unfairness between languages. *arXiv preprint arXiv:2305.15425*.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. *Qwen2.5 technical report*. *Preprint*, arXiv:2412.15115.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. *Direct preference optimization: Your language model is secretly a reward model*. *Preprint*, arXiv:2305.18290.
- David Rein and 1 others. 2023. Gpqa: A google-proof q&a benchmark for large language models. *arXiv preprint arXiv:2311.16452*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. *Proximal policy optimization algorithms*. *Preprint*, arXiv:1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. *Deepseekmath: Pushing the limits of mathematical reasoning in open language models*. *Preprint*, arXiv:2402.03300.
- Oleh Shliachko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. 2023. *mgpt: Few-shot learners go multilingual*. *Preprint*, arXiv:2204.07580.
- Jacob Mitchell Springer, Sachin Goyal, Kaiyue Wen, Tanishq Kumar, Xiang Yue, Sadhika Malladi, Graham Neubig, and Aditi Raghunathan. 2025. *Overtrained language models are harder to fine-tune*. In *Forty-second International Conference on Machine Learning*.
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Na Zou, Hanjie Chen, and Xia Hu. 2025. *Stop overthinking: A survey on efficient reasoning for large language models*. *Preprint*, arXiv:2503.16419.
- T-Tech. 2025a. ru-alpaca-eval. <https://huggingface.co/datasets/t-tech/ru-alpaca-eval>. Hugging Face Datasets; accessed 2025.
- T-Tech. 2025b. ru-arena-hard. <https://huggingface.co/datasets/t-tech/ru-arena-hard>. Hugging Face Datasets; accessed 2025.
- T-Tech. 2025c. ru-mt-bench. <https://huggingface.co/datasets/t-tech/ru-mt-bench>. Hugging Face Datasets; accessed 2025.
- T-Tech. 2025d. ru-reasoning-benchmarks. <https://huggingface.co/collections/t-tech/ru-reasoning-benchmarks>. Hugging Face Datasets; accessed 2025.
- T-Tech. 2025e. ruaime-2024. <https://huggingface.co/datasets/t-tech/ruAIME-2024>. Hugging Face Datasets; accessed 2025.
- T-Tech. 2025f. rummlu-pro. <https://huggingface.co/datasets/t-tech/ruMMLU-pro>. Hugging Face Datasets; accessed 2025.
- T-Tech. 2025g. T-math. <https://huggingface.co/datasets/t-tech/T-math>. Hugging Face Datasets; accessed 2025.
- Sijun Tan, Siyuan Zhuang, Kyle Montgomery, William Yuan Tang, Alejandro Cuadron, Chenguang Wang, Raluca Popa, and Ion Stoica. 2025. *Judgebench: A benchmark for evaluating LLM-based judges*. In *The Thirteenth International Conference on Learning Representations*.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. *Gemma 3 technical report*. *Preprint*, arXiv:2503.19786.
- Mikhail Tikhomirov and Daniil Chernyshev. 2024. *Facilitating large language model russian adaptation with learned embedding propagation*. *arXiv preprint arXiv:2412.21140*.
- Tianchun Wang, Zichuan Liu, Yuanzhou Chen, Jonathan Light, Weiyang Liu, Haifeng Chen, Xiang Zhang, and Wei Cheng. 2025a. *On the effect of sampling diversity in scaling llm inference*. *Preprint*, arXiv:2502.11027.

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, and 1 others. 2023. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.
- Zhilin Wang, Jiaqi Zeng, Olivier Delalleau, Hoo-Chang Shin, Felipe Soares, Alexander Bukharin, Ellie Evans, Yi Dong, and Oleksii Kuchaiev. 2025b. [HelpSteer3-Preference: Open human-annotated preference data across diverse tasks and languages](#). *Preprint*, arXiv:2505.11475.
- Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong Zhang, Caiming Xiong, and Hanze Dong. 2025. [A minimalist approach to llm reasoning: from rejection sampling to reinforce](#). *Preprint*, arXiv:2504.11343.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- An Yang, Baosong Yang, Binyuan Hui, and 1 others. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. 2024. [Mammoth2: Scaling instructions from the web](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 90629–90660. Curran Associates, Inc.
- Yifan Zhang and Team Math-AI. 2024. American invitational mathematics examination (aime) 2024.
- Yifan Zhang and Team Math-AI. 2025. American invitational mathematics examination (aime) 2025.
- Weilin Zhao, Tengyu Pan, Xu Han, Yudi Zhang, Ao Sun, Yuxiang Huang, Kaihuo Zhang, Weilun Zhao, Yuxuan Li, Jie Zhou, Hao Zhou, Jianyong Wang, Zhiyuan Liu, and Maosong Sun. 2025. [FR-spec: Accelerating large-vocabulary language models via frequency-ranked speculative sampling](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3909–3921, Vienna, Austria. Association for Computational Linguistics.
- Dmitry Zmitrovich, Alexander Abramov, Andrey Kalmykov, Maria Tikhonova, Ekaterina Taktasheva, Danil Astafurov, Mark Baushenko, Artem Snegirev, Vitalii Kadulin, Sergey Markov, Tatiana Shavrina, Vladislav Mikhailov, and Alena Fenogenova. 2023. A family of pretrained transformer language models for russian. *arXiv preprint arXiv:2309.10931*.

## A Released Resources and Licenses

All released components use permissive, research-friendly licenses. The T-pro 2.0 model, its EA-GLUE draft weights, and the T-Math benchmark are distributed under Apache-2.0, allowing broad academic and commercial use. The T-Wix 500k corpus is released under ODC-By. Full license details appear in Table 6.

## B Tokenizer adaptation statistics

**Russian and English corpora.** Table 7 reports tokenization statistics for Russian and English on both Wikipedia and our in-domain SFT corpus (T-Wix). For Russian, the Cyrillic-dense T-pro 2.0 tokenizer substantially reduces the average number of tokens per word and increases the share of words segmented into at most two tokens, while English compression is essentially unchanged.

Table 8 extends this analysis to eight Cyrillic languages using Wikipedia. In all cases the new tokenizer reduces tokens per word, with particularly large gains for Kyrgyz, which is poorly served by generic multilingual tokenizers.

**Comparison with other Cyrillic-rich tokenizers.** Finally, Table 9 compares T-pro 2.0 against several strong Cyrillic-focused baselines. T-pro 2.0 achieves the lowest tokens-per-word for seven out of eight languages (ru, uk, be, bg, sr, mk, ky) and remains competitive on Kazakh, demonstrating that our tokenizer design is competitive with specialized alternatives.

## C Instructional midtraining

We employ an intermediate *instructional midtraining* stage between generic large-scale pre-training and downstream alignment. Starting from the publicly available Qwen3-32B dense model (Yang et al., 2025), already pre-trained on  $\sim 36T$  tokens, we perform continual training on 40B tokens of instruction-style data. The goals of this stage are: (i) adapt the model to a denser, Russian-centric tokenizer, (ii) learn useful representations for new subword units, and (iii) further strengthen Russian language and reasoning skills without degrading the base model’s capabilities.

**Training setup.** Midtraining uses a 4M-token global batch and 40B total tokens ( $\approx 9,750$  steps). We train with AdamW using a peak learning rate of  $1 \times 10^{-5}$  and cosine decay to  $1 \times 10^{-6}$ , with 100 warmup batches. Data are formatted in the same

chat-style schema and packed up to a 32K context window without a length curriculum. Training is performed in bf16 with FSDP full-shard and activation checkpointing. Remaining hyperparameters are listed in Table 10.

**Midtraining datamix** All midtraining samples are in instruction format. The datamix combines (i) public instruction datasets from the Hugging Face Hub, (ii) web and forum data (e.g., question-answer style threads), (iii) real user–assistant dialogues, and (iv) synthetic instructional data and reasoning traces grounded in pre-training corpora (books, Common Crawl, code) via a WebInstruct-style pipeline (Yue et al., 2024). Compared to the SFT stage, the midtraining datamix is intentionally *larger and less curated*: we trade some noise for broader coverage of tasks and domains. All instructions are derived from public sources; internal data are used only as anonymized targets in dialogue-style responses.

Table 11 reports the category-level breakdown of the 40B-token corpus. In terms of language, the corpus is predominantly Russian and English: roughly 49% of tokens are Russian (19.6B), 36% English (14.4B), 9.3% code (3.7B) and 5.5% come from parallel Russian–English data (2.2B).

**InsTag deduplication.** To control redundancy while preserving diversity across sources, we apply #INSTAG-based deduplication (Lu et al., 2024) *independently within each component* of the datamix (reasoning, general QA, code, etc.). The tagger is applied only to user utterances; all tags from a multi-turn sample are unioned into a single tag set. We then perform exact-match and semantic deduplication at the tag level, followed by greedy diversity sampling over tagged samples. This procedure gives macro-level control over the balance between different categories instead of deduplicating the raw pool as a whole. On large components such as reasoning and general QA, only about 10–30% of the raw candidates are retained (the remaining 70–90% are discarded), whereas for smaller, less repetitive sources we keep 80–90% of samples. For each retained sample, the final assistant turn is regenerated with a stronger teacher, Qwen3-235B, which improves answer quality and stylistic consistency while keeping the original user input and context intact.

**Ablations on datamix design** We compare two variants of the midtraining corpus, both trained

Resource	Type	Location	License
T-pro 2.0	Model	<a href="https://huggingface.co/t-tech/T-pro-it-2.0">https://huggingface.co/t-tech/T-pro-it-2.0</a>	Apache-2.0
EAGLE weights	Model component	<a href="https://huggingface.co/t-tech/T-pro-it-2.0-eagle">https://huggingface.co/t-tech/T-pro-it-2.0-eagle</a>	Apache-2.0
T-Math	Benchmark dataset	<a href="https://huggingface.co/datasets/t-tech/T-math">https://huggingface.co/datasets/t-tech/T-math</a>	Apache-2.0
T-Wix 500k	Instruction corpus	<a href="https://huggingface.co/datasets/t-tech/T-wix">https://huggingface.co/datasets/t-tech/T-wix</a>	ODC-By

Table 6: Released resources and licenses.

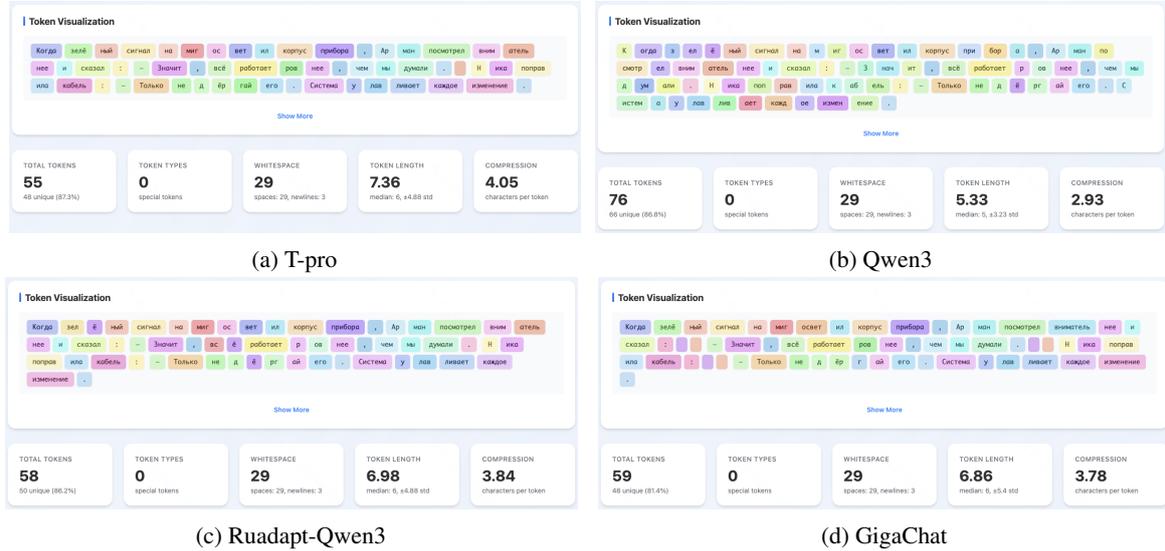


Figure 2: Qualitative comparison of Russian tokenization. A 220-character text is tokenized by T-pro 2.0, the original Qwen3 tokenizer, and other Cyrillic-optimized models. T-pro 2.0 encodes the text into just 55 tokens compared to 76 for Qwen3, demonstrating superior compression efficiency.

Corpus	Tokenizer	tok/ word	1 tok (%)	$\leq 2$ tok (%)	$> 2$ tok (%)
<i>Russian</i>					
ruWiki	Qwen3	3.12	20.3	38.2	61.8
ruWiki	T-pro 2.0	2.38	28.7	60.1	39.9
T-Wix	Qwen3	2.70	31.8	52.4	47.6
T-Wix	T-pro 2.0	2.26	39.3	65.5	34.5
<i>English</i>					
enWiki	Qwen3	1.68	61.2	83.7	16.3
enWiki	T-pro 2.0	1.68	61.1	83.7	16.3

Table 7: Tokenization density statistics for Russian and English on Wikipedia and our SFT corpus (T-Wix). We compare the original Qwen3 and Cyrillic-dense T-pro 2.0 tokenizers. Columns show: average tokens per word (tok/ word), percentage of words tokenized into exactly 1 token, at most 2 tokens, and more than 2 tokens.

Lang	Tokens/Word		% Words ( $\leq 2$ tok)	
	Qwen3	T-pro	Qwen3	T-pro
ru	3.12	2.38	38.20	60.13
uk	3.70	2.80	31.17	45.79
be	3.97	2.94	30.15	41.36
bg	2.99	2.35	43.42	59.60
sr	3.26	2.62	37.65	51.79
mk	3.04	2.41	42.42	57.19
kk	4.60	3.07	15.30	37.69
ky	4.35	3.09	21.27	39.93

Table 8: Tokenization density on Wikipedia for Cyrillic languages. *Tokens/Word*: average tokens per word; *% Words ( $\leq 2$  tok)*: percentage of words tokenized into at most 2 tokens.

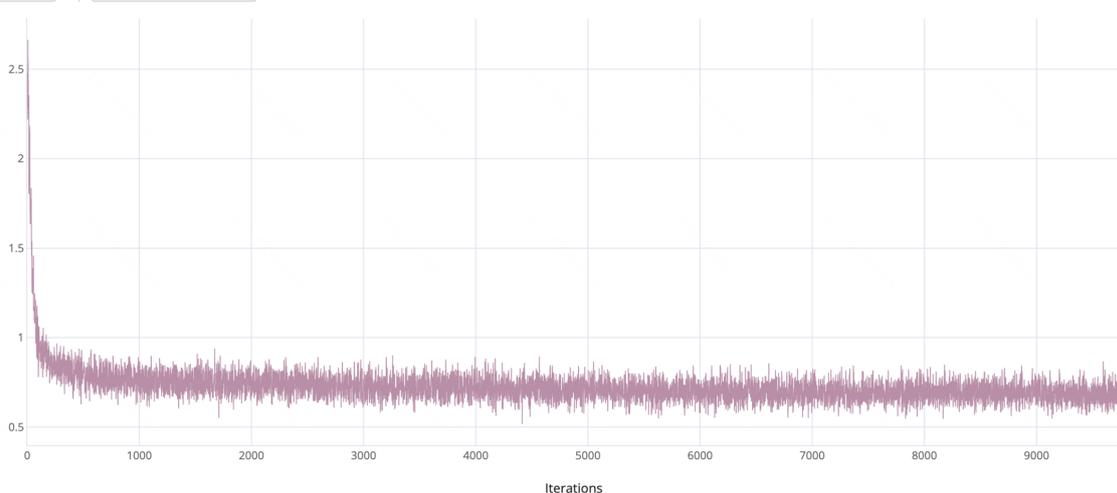


Figure 3: Midtraining training loss as a function of optimization steps. Loss drops steeply during the first  $\sim 1k$  steps ( $\sim 4B$  tokens) and then gradually plateaus, indicating that most adaptation to the new tokenizer happens early in the run.

Lang	T-pro	GigaChat <sup>†</sup>	Ruadapt-Qwen3	gpt-oss
ru	<b>2.38</b>	2.49	2.43	2.70
uk	<b>2.80</b>	3.09	3.29	2.92
be	<b>2.94</b>	3.32	3.54	3.03
bg	<b>2.35</b>	2.58	2.50	2.56
sr	<b>2.62</b>	2.97	3.07	2.73
mk	<b>2.41</b>	2.67	2.70	2.59
kk	3.07	<b>2.67</b>	4.60	3.11
ky	<b>3.09</b>	3.33	3.97	3.17
Avg	<b>2.71</b>	2.89	3.26	2.85

Table 9: Tokenization density (tokens/word) on Wikipedia for T-pro and other Cyrillic-dense tokenizers. Lower is better. <sup>†</sup>Indicates <https://huggingface.co/ai-sage/GigaChat-20B-A3B-instruct> model.

for 40B tokens with identical optimization settings (Table 10):

- **Pre-train + instruct:** mixture including generic pre-training-style data (Common Crawl, Wikipedia, code) alongside instruction-formatted examples.
- **Instruct-only:** the same instruction pool but without additional raw pre-training sources, i.e., all examples follow an explicit instruction–response schema.

The instruct-only variant thus allocates more of the 40B-token budget to high-quality instruction data, whereas the mixed variant spends a fraction of tokens on generic web/code continuation. We evaluate both models on a suite of Russian and multilingual math/reasoning benchmarks, including ru-

Hyperparameter	Value
Global batch size (tokens)	4M (128 seq $\times$ 32K)
Total tokens	40B
Steps	$\approx 9,750$
Max context length	32K
Optimizer	AdamW
Adam betas	(0.9, 0.95)
Adam $\epsilon$	$10^{-8}$
Weight decay	$10^{-6}$
LR schedule	cosine
Peak / min LR	$1 \times 10^{-5} / 1 \times 10^{-6}$
Warmup	100 batches
Gradient clipping	max norm 1.0
Precision	BF16
Parallelism	FSDP full-shard, act. checkpointing

Table 10: Midtraining optimization setup.

Category	Share	Tokens (B)
Reasoning	34.5%	13.8
General	29.3%	11.7
Math	16.3%	6.5
Real chat	5.5%	2.2
IF	5.0%	2.0
Grounded QA synth	3.8%	1.5
Code	2.8%	1.1
Forum	1.7%	0.7
Summarization	0.7%	0.3
ICL	0.4%	0.2

Table 11: Midtraining datamix by category (40B tokens total).

AIME’24/25, ruMATH500, ruGPQA, ruLCB, T-Math, and Arena-style pairwise comparisons. All evaluations are zero-shot; AIME-style metrics are computed as avg@8 over 30 problems, and other benchmarks are run once due to computational

cost.

Table 12 shows a representative subset of metrics. Across most math and reasoning benchmarks, the instruct-only datamix outperforms or matches the mixed variant despite using the same token budget, and even early checkpoints from the instruct-only run are ahead of the mixed model. This is consistent with recent evidence that heavily pre-trained models are harder to adapt via continual pre-training (Springer et al., 2025), especially when the additional data distribution differs from downstream tasks.

Benchmark	PT+I	I-only
ruAIME'24	0.60	<b>0.67</b>
ruAIME'25	0.47	<b>0.63</b>
ruMATH500	0.93	<b>0.94</b>
ruGPQA	0.58	<b>0.66</b>
ruLCB	0.53	<b>0.55</b>
T-Math	0.49	<b>0.50</b>
Arena hard (think)	43.7	<b>44.5</b>
Arena wildchat ru (think)	55.0	<b>55.1</b>

Table 12: Ablation on midtraining datamix design (zero-shot). “PT+I” denotes the pre-training+instruct mixture; “I-only” uses only instruction-formatted data.

We did not run additional ablations such as training on original (non-regenerated) answers or disabling InsTag deduplication. In the first case, instruction data come from heterogeneous sources with uneven answer quality and formats, and we found it undesirable to train on such out-of-distribution completions. In the second case, skipping deduplication would require regenerating answers for a much larger pool of raw samples, significantly increasing computational cost; we leave this exploration for future work.

**Tokenizer adaptation and MERA results** A key objective of midtraining is to adapt the model to a new, denser tokenizer for Russian without degrading downstream quality. To quantify the impact of the tokenizer choice, we train two 8B models on the same midtraining datamix with identical optimization hyperparameters, differing only in the tokenizer (original Qwen3 vs. T-pro 2.0).

Table 13 reports MERA scores for these two variants. The T-pro 2.0 tokenizer attains a macro-average score comparable to the original one (0.574 vs. 0.560), with only small per-task differences in both directions. In other words, replacing the tokenizer with a denser Cyrillic segmentation does not degrade general Russian-language performance

on MERA, which is the primary design goal of midtraining.

Task	Qwen tokenizer	T-pro 2.0 tokenizer
USE	<b>0.198</b>	0.191
MaMuRaMu	0.784	<b>0.796</b>
ruWorldTree	0.966/0.966	0.966/0.966
ruCodeEval	0.173/0.45/0.585	<b>0.454/0.689/0.756</b>
RCB	0.557/0.479	<b>0.564/0.47</b>
MathLogicQA	0.710	<b>0.731</b>
ruOpenBookQA	0.897/0.897	<b>0.922/0.923</b>
RWSD	<b>0.446</b>	0.250
CheGeKa	0.30/0.368	<b>0.31/0.384</b>
LCS	<b>0.102</b>	0.096
PARUS	0.868	<b>0.912</b>
MultiQ	<b>0.381/0.517</b>	0.344/0.478
ruMultiAr	<b>0.402</b>	0.400
ruTiE	0.788	<b>0.798</b>
ruModAr	0.515	<b>0.627</b>
AVG	0.560	<b>0.574</b>

Table 13: MERA scores for an 8B model with the original Qwen3 tokenizer and the Cyrillic-dense T-pro 2.0 tokenizer. Bold marks the better value per row.

The midtraining loss curve in Figure 3 further illustrates the adaptation process. Training loss decreases sharply over the first  $\approx 1k$  steps ( $\approx 4B$  tokens) before plateauing, indicating that a substantial token budget at a relatively high learning rate is required to adapt the model to the new tokenizer beyond what the smaller, lower-LR SFT budget alone could provide.

## D T-Wix SFT dataset

### D.1 General part of T-Wix

The general part of the dataset consists of 468k diverse prompts collected from open-source data and high-quality translations of English-language datasets, subsequently deduplicated. The dataset is assembled to enhance the model’s capabilities in coding, mathematics, dialogue, and other competencies expected from a modern LLM.

#### D.1.1 Data Preparation

First and foremost, a corpus of 14M instructions (mostly in English) is compiled from various open-source datasets. To select the most useful samples, a data filtering pipeline is developed. It consists of several consecutive stages aimed at deduplication and ensuring high thematic, qualitative, and complexity diversity of the SFT dataset. We also perform deduplication against benchmark datasets to ensure that no benchmark examples leak into the training corpus.

**LSH and Embedding-Based Deduplication.** At the initial stage, simple deduplication is performed using locality-sensitive hashing (LSH) and embedding-based similarity search to eliminate duplicated samples originating from different open-source datasets.

**Thematic Tag Filtering.** To ensure thematic balance, the #INSTAG-based filtering approach (Lu et al., 2024) is employed. The pipeline uses a trained tagging model to extract thematic tags from each instruction.

For the present work, the tagger is trained using the Qwen2.5-7B (Qwen et al., 2025) model on multilingual data with a context length of up to 32k tokens, allowing tagging in both Russian and English, including long-context data. This modification substantially reduces translation overhead, as tagging and filtering can be applied directly to multilingual raw data without prior translation into English.

To improve thematic balance, an additional domain-balancing stage—“Domain & Complexity Balancing”—is included, as the tagger primarily produces low-level thematic annotations.

**Domain & Complexity Balancing.** In addition to fine-grained thematic filtering, a higher-level balance is introduced across major knowledge domains and difficulty levels within each domain. To achieve this, six domains — *Math, Code & Programming, Science, General Instruct, General Knowledge, Writing* — and three complexity levels — *School, Student, Professor* — are defined.

Using large-scale LLM-assisted annotation, approximately 14M samples are automatically labeled with both domain and complexity tags. Subsequently, the dataset is balanced across domains and further normalized by difficulty within each domain to regulate the model’s output capabilities and ensure a uniform skill distribution.

This stage enables finer control over the resulting model’s generalization behavior, preventing over-representation of specific topics or difficulty levels.

**Reward Model Filtering.** In the subsequent stage, samples are filtered according to prompt quality using scores from the Reward Model (RM) described in F. For each of the datasets comprising the 14M instructions, an RM score is computed, and the bottom 10% of samples with the lowest scores within each dataset are filtered out.

This step effectively removes “noisy” or low-quality samples that could negatively impact downstream model performance, preserving only high-quality and instructionally meaningful examples.

**Instruction Following Difficulty Filtering.** A further filtering stage based on Instruction Following Difficulty (IFD) scores is incorporated, following the approach introduced in (Li et al., 2024a). These scores quantify the difficulty a language model faces in following a given instruction. For the present work, IFD scores are computed relative to a midtraining checkpoint to reflect the model’s actual instruction-following capability. Samples with excessively high IFD values ( $>1.0$ ) are discarded as overly complex or ambiguous, while those with very low IFD scores ( $<0.7$ ) are filtered out as trivially simple.

This selective filtering makes it possible to retain the most challenging and instructionally rich examples—those that contribute most to improving the model’s instruction-following ability—while removing both overly simple and excessively difficult samples.

**Multilingual Filtering and Translation.** The multilingual setup enables filtering to be conducted directly on mixed-language raw data, reducing both the cost and time associated with preliminary translation. Only the final curated dataset is translated into Russian to ensure cross-lingual consistency.

**Rejection Sampling and Generation.** High dataset quality is further ensured through the use of top LLMs and rejection sampling. Each final training completion is produced using DeepSeek-V3 and Qwen-235B-A22B models, generating 8 candidate responses per instruction. These candidates are then filtered using RM scores to select the highest-quality outputs.

This approach not only eliminates translation artifacts present in the raw multilingual data but also results in substantially higher-quality responses compared to the original samples, thereby improving the overall consistency and instructional value of the dataset.

Overall, the combined multistage filtering pipeline ensures that the final SFT dataset is diverse, balanced, and composed of high-quality, instructionally valuable samples, free from data leakage and redundancy.

This approach allows the training process to remain balanced across domains (e.g., code and

math) without bias toward any particular category.

## D.2 Long Context

To enhance the model’s ability to process extended inputs, a dedicated long-context dataset is constructed.

A diverse collection of long texts is selected from publicly available data for the pre-training phase, covering various domains such as *education*, *technology*, *business*, *scientific literature*, and *fiction*. The dataset is distributed across multiple context lengths ranging from 8k to 32k tokens, enabling the model to learn robustly across different input sizes. Using DeepSeek-V3 and Qwen-235B-A22B, a variety of prompts and responses are generated for each text, encompassing summarization, open- and closed-domain QA tasks, as well as reasoning-oriented datasets.

The resulting long-context dataset increment constitutes about 1% of the total SFT training data in samples and 7.7% in tokens, providing valuable coverage for instruction tuning under extended context conditions.

## D.3 Parallel Corpora

To maintain strong English proficiency, parallel corpora are added to the SFT dataset — that is, instructional samples presented in English alongside their Russian counterparts.

A series of experiments on the language share within the dataset shows that an optimal ratio is approximately 10% English data relative to the total SFT mix.

## D.4 Reasoning part of T-Wix

To enhance reasoning capabilities in Large Language Models (LLMs) for the Russian language, a high-quality, reasoning-focused dataset is constructed through a targeted distillation pipeline. Rather than maximizing data volume, the pipeline prioritizes instructional value and appropriate task difficulty, ensuring that each retained sample provides meaningful learning potential for the target student model. The process starts from a broad collection of English-language reasoning instructions, which are subsequently translated into Russian, deduplicated, and carefully balanced across domains to support diverse and robust reasoning behaviors.

**Initial Pool Generation and Deduplication.** The initial pool of data is constructed from approx-

imately 450k high-quality English-language reasoning instructions, drawn from established open-source datasets (e.g., Open-R1 (Hugging Face, 2025), Nvidia/AceReason-Math (Chen et al., 2025), Nvidia/Nemotron (Bercovich et al., 2025)), covering general knowledge, mathematics, natural sciences, and code generation.

Domain Distribution:

- 60% general knowledge and open-ended reasoning (to establish fluent, structured reasoning in Russian),
- 10% verifiable mathematics (e.g., arithmetic, algebra),
- 10% open-ended mathematics (e.g., proofs, conceptual explanations),
- 15% natural sciences (physics, chemistry, biology),
- 5% code-related reasoning.

After the initial collection, these English-language instructions are translated into Russian. As in the general part of T-Wix, deduplication is applied to eliminate near-duplicates and ensure sample uniqueness.

**Reward-Based Completion Evaluation.** To mitigate the stochasticity inherent in LLM generation (Wang et al., 2025a; Atil et al., 2025), 8 diverse completions are generated per instruction by both:

- The teacher model (Qwen3-235B-A22B),
- The student model (midtraining checkpoint).

This yields 16 completions per instruction, which are independently scored by a trained reward model (RM). The inclusion of student generations enables direct assessment of the model’s current reasoning capability on each task, while the teacher generations provide high-quality reference behaviors. This multi-generation approach provides a more robust statistical picture of the model’s performance on each instruction.

**Statistical Filtering Based on Reward Stability.** Instructions exhibiting high variance in RM scores across generations are discarded, as they reflect ambiguous or unstable evaluation signals. Additionally, instructions for which the student model consistently receives very low RM scores—even

with low variance—are excluded, as they lie beyond the student’s current learning capacity and are unlikely to support effective knowledge transfer (Xiong et al., 2025; Liu et al., 2025c).

**Mean Reward-Based Selection Within the Zone of Proximal Development.** To operationalize the pedagogical principle of the zone of proximal development (ZPD) (Cui and Sachan, 2025), the RM scores for the 8 teacher and 8 student completions per instruction are aggregated by computing their respective means. The average reward of the teacher responses and the average reward of the student responses are then used to estimate the reasoning gap.

Samples are selected based on the absolute difference between these mean rewards. A small difference indicates that the student already performs comparably to the teacher (suggesting limited learning potential), whereas a very large difference implies that the task lies beyond the student’s capabilities (making distillation ineffective). Only samples with a moderate gap in mean RM scores—neither too small nor too large—are retained. This ensures that the selected samples are challenging enough to drive improvement, yet sufficiently within reach for successful knowledge transfer.

**Final Completion Selection.** For the final training targets, teacher-generated completions are selected as follows:

- For *verifiable* instructions (e.g., mathematical problems), factually incorrect completions are first filtered out; among the remaining correct ones, the completion with the highest RM score is chosen.
- For *open-ended* instructions, the shortest reasoning trace among the top-3 RM-ranked teacher completions is selected. This encourages the student to learn concise and non-redundant reasoning patterns (Sui et al., 2025).

This pipeline yields a high-quality reasoning dataset of approximately 30k samples, consisting of 90% Russian and 10% English instructions, consistent with the overall language strategy of T-Wix. By design, the dataset emphasizes stable, diverse, and pedagogically optimal reasoning traces in Russian across multiple domains, effectively balancing task difficulty with learning potential.

## D.5 T-Wix dataset analytics

The final distribution of data in the SFT dataset (T-Wix) are presented in Figure 4. The total size is 500k samples.

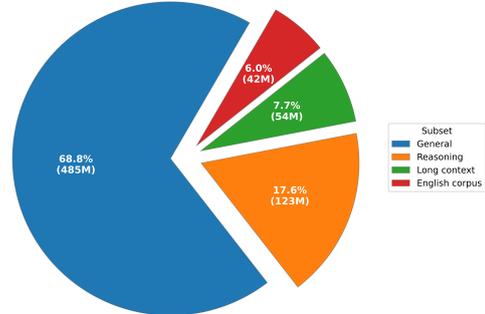


Figure 4: Token distribution in T-Wix. Token counts were computed using tiktoken with o200k base tokenizer.

## D.6 SFT Training Recipe

The SFT stage took 9 hours on 4 nodes with 8xH100 GPUs for T-pro 2.0, using gradient checkpointing and FSDP, as well as packing samples into 32k-token contexts without truncation. After a series of experiments, the optimal fine-tuning hyperparameters were selected, as described in Table 14.

Hyperparameter	Value
Global batch size (samples)	32
Max context length	32k
Number of training epochs	2
Optimizer	Adam (Kingma and Ba, 2017)
Adam betas	(0.9, 0.95)
Adam $\epsilon$	$10^{-12}$
Learning rate	1e-6
Learning rate scheduler	cosine
Warmup ratio	0.1
Gradient clipping	max norm 2.0
Precision	BF16

Table 14: Hyperparameters used for SFT training.

## E Preference tuning

To enhance alignment beyond supervised fine-tuning, an on-policy Direct Preference Optimization (DPO) procedure is applied. Recent work shows that on-policy preference optimization offers more stable and reliable alignment gains than off-policy alternatives, as it learns directly from the model’s own generative distribution and therefore avoids distribution shift while targeting realistic error modes (Rafailov et al., 2024; Im and Li, 2025).

For each instruction, the SFT-trained model produces 16 completions. All candidates are scored using the RM described in App. F, and preference pairs are constructed by selecting the highest- and lowest-scoring completions. This contrastive selection yields stable and informative training pairs by removing low-signal, ambiguous comparisons.

The DPO dataset is constructed from filtered SFT data (T-Wix). A total of 100k preference pairs is formed, consisting of:

- 90k sampled from the General SFT part,
- 10k sampled from the Reasoning SFT part.

In addition, cross-subset augmentation is applied to enrich preference diversity: 4k samples from the General subset are paired with reasoning-style reformulations, while 6k samples from the Reasoning subset are converted into general-style instructions. This yields a smoother distribution of reasoning complexity without altering the intended emphasis of each subset.

The resulting on-policy DPO stage improves the model’s alignment, coherence, and reasoning structure while preserving broad general-purpose capabilities.

### E.1 Preference Training Recipe

The DPO stage required 28 hours of training on 4 nodes with 8×H100 GPUs. The training was carried out using sequence parallelism, which enabled efficient distribution of computation across devices. The hyperparameters listed in Table 15 were identified as optimal.

Hyperparameter	Value
Global batch size (samples)	128
Max context length	32k
Number of training epochs	1
Optimizer	AdamW
Adam betas	(0.9, 0.95)
Adam $\epsilon$	$10^{-12}$
Weight decay	0.01
Learning rate	1e-7
Learning rate scheduler	cosine
Warmup ratio	0.05
Gradient clipping	max norm 2.0
Precision	BF16
Loss type	DPO
DPO beta	0.5

Table 15: Hyperparameters used for preference tuning

## F Reward Model

**Tournament-Based Synthetic Preference Data Generation** To construct a high-quality reward model, it is essential to obtain reliable preference data —pairs of model completions ranked according to their relative quality. Direct pairwise annotation across all available completions, however, is computationally expensive and inefficient. To address this, similar to the knockout-tournament method introduced by Liu et al. (2025b), we propose a tournament-based preference generation approach that substantially reduces the number of required comparisons while preserving the informativeness of the resulting preference signal.

Each tournament comprises  $n$  participants, randomly sampled from the pool of available models. For each instruction, every model generates a completion, and the tournament bracket is constructed according to model category —for instance, small-scale models (7B–13B) are paired against models of similar scale, and reasoning-oriented models compete within the same subclass. This grouping strategy ensures that comparisons are made between models of comparable generative quality, encouraging the reward model to learn fine-grained distinctions rather than relying on trivial cases where one output is clearly superior (e.g., when a large model is compared with a small size model, a comparison might yield an obvious outcome —the larger model would consistently produce more coherent and contextually appropriate responses, leaving little room for the reward model to learn subtle differences).

Each round of the tournament consists of a single instruction and the corresponding completions generated by the competing models. An external LLM, not participating in the tournament, is employed as judge to determine the preferred completion for each matchup. To avoid positional bias, each pair of completions is evaluated in both possible orders, and samples exhibiting positional bias are excluded from the final training set.

At the completion of each single-elimination tournament with  $n$  participants, a total of  $\frac{n}{2} \log_2 N$  preference pairs are obtained. This result comes from the hierarchical structure of the tournament: in each round, half of the remaining participants compete, producing  $\frac{n}{2}$  new pairwise outcomes (both direct and transitive). Because a tournament with  $n$  participants requires  $\log_2 N$  rounds to determine a winner, the total number of inferred prefer-

ence pairs accumulates to  $\frac{n}{2} \log_2 N$ .

Each round contributes the same number of new known preferences because every winner’s new victory also establishes transitive relationships over all opponents defeated in earlier rounds. For instance, if player A beats player B in the final, it is implied that A outperform every player that B previously defeated. Consequently, even though only  $n - 1$  matches are directly played, the tree-like transitive structure allows many additional indirect comparisons to be inferred.

This process produces preference set dense enough to capture comparative information among many participants, yet far more efficient than exhaustively comparing every possible pair (which would require  $\frac{n(n-1)}{2}$  comparisons). This tournament-based approach yields an informative preference dataset while significantly reducing annotation complexity.

**Reward Model Training** The reward model is based on Qwen3-32B (Yang et al., 2025) with a regression head to produce a single preference score for each completion. Training follows the Bradley–Terry (BRADLEY and TERRY, 1952) formulation, which models the probability of one completion being preferred over another as a logistic function of their respective scores. All training is conducted with a maximum sequence length of 32k tokens, leveraging Ulysses sequence parallelism (Jacobs et al., 2024) to efficiently support long-context optimization. Data preprocessing, batching, and distributed training are managed through the TurboAlignment library (tur, 2025).

**Evaluation** For intrinsic evaluation, we adapt RewardBench 2 (Malik et al., 2025) to Russian by translating the original benchmark and report standard leaderboard metrics. For downstream evaluation, we additionally construct a Best-of- $N$  selection benchmark on top of the Arena-Hard-RU instruction set to assess the reward model under realistic generation scenarios. In this setting, the base model produces  $N$  candidate completions per instruction, and the reward model selects the highest-scoring (best@ $N$ ) and lowest-scoring (worst@ $N$ ) outputs. These selections are then evaluated using Arena-Hard, allowing us to measure the alignment between reward-model rankings and externally validated quality. We further report the  $\Delta_{\text{BoN}}$  metric (best@ $N - \text{worst@}N$ ) to quantify discriminative capacity. Although our model performs compara-

bly to existing open-source reward models on the translated RewardBench 2, it demonstrates a clear advantage on our Best-of- $N$  Arena-Hard benchmark. As shown in Table 16, our model obtains the highest  $\Delta_{\text{BoN}}$  score, reflecting the strongest separation between high- and low-quality completions.

RM-model	best@8↑	worst@8↓	$\Delta_{\text{BoN}}$
<b>Qwen3-32B-RM (Ours)</b>	<b>92.69</b> (-0.99)	<b>70.48</b> (+2.34)	<b>22.21</b>
Llama-3.3-Nemotron-70B-Reward-Multilingual <sup>1</sup>	85.93 (-1.93)	84.91 (+1.85)	1.02
Skywork-Reward-Gemma-2-27B <sup>2</sup>	89.05 (-1.6)	74.35 (+2.07)	14.70
Skywork-Reward-V2-Llama-3.1-8B <sup>3</sup>	90.49 (-1.43)	77.31 (+1.77)	13.18
Llama-3.1-Tulu-3-70B-SFT-RM-RB2 <sup>4</sup>	87.37 (-1.86)	78.47 (+1.76)	8.90

Table 16: Best-of- $N$  ( $N = 8$ ) evaluation on Arena-Hard-RU. We report win rates for the highest- (best@8) and lowest-scoring (worst@8) completions selected by each reward model, and their difference  $\Delta_{\text{BoN}} = \text{best@8} - \text{worst@8}$ , which measures discriminative capacity. <sup>1</sup>Wang et al. (2025b), <sup>2</sup>Liu et al. (2024), <sup>3</sup>Liu et al. (2025a), <sup>4</sup>Malik et al. (2025)

**Prompt selection** Furthermore, in the process of synthetic data generation, we evaluated a range of prompting strategies derived from the JudgeBench (Tan et al., 2025). Empirical analysis indicates that the Google Vertex prompt yields superior evaluation quality in different benchmarks (see Table 18), particularly on RewardBench 2 (RU). This improvement underscores the sensitivity of LLM-based evaluators to prompt design and highlights the importance of selecting domain-appropriate judging configurations for reliable preference data generation.

## F.1 Reward Model Analysis

Our experiments reveal that DeepSeek-V3 (DeepSeek-AI et al., 2025b) demonstrates superior judgment capabilities in open-domain and conversational (chat) tasks, whereas Qwen3-235B-A22B exhibits stronger performance in mathematical, code-related and other domains (see Table 17).

**Ablation on Transitive Samples.** An ablation study was conducted to evaluate the contribution of transitive preference pairs. Removing transitive samples led to a consistent degradation across all evaluation metrics (see Table 17), suggesting that inferred pairwise relationships enrich the preference signal and improve the model’s generalization to unseen instructions. Conversely, adding additional transitive samples beyond the first closure continued to yield marginal but positive improvements.

Category	Model	RewardBench 2 (RU)					
		Fact.	Focus	Math	Prec. IF	Safety	Total
Comparison with Existing RMs	<b>Qwen3-32B-RM (Ours)</b>	0.66	0.87	0.62	0.42	0.89	0.69
	Skywork-Reward-V2-Llama-3.1-8B	0.68	<b>0.88</b>	0.65	<b>0.45</b>	0.79	0.69
	Skywork-Reward-Gemma-2-27B	0.69	<b>0.88</b>	0.64	0.40	<b>0.92</b>	<b>0.7</b>
	Llama-3.1-Tulu-3-70B-SFT-RM-RB2	0.72	0.74	<b>0.69</b>	0.41	0.76	0.66
	Llama-3.3-Nemotron-70B-Reward-Multilingual	<b>0.73</b>	0.85	0.62	0.41	0.86	0.69
Judge Model Ablation	Qwen-3-RM-8B-DeepSeek-V3	<b>0.478</b>	<b>0.756</b>	0.598	0.341	0.736	<b>0.581</b>
	Qwen-3-RM-8B-Qwen3-235B-A22B	0.467	0.324	<b>0.688</b>	<b>0.350</b>	<b>0.840</b>	0.533
Transitive Samples Ablation	Qwen3-8B-RM w/o transitive	0.467	0.324	0.688	0.350	0.840	0.533
	Qwen3-8B-RM w/ transitive	<b>0.505</b>	<b>0.453</b>	<b>0.704</b>	<b>0.413</b>	<b>0.860</b>	<b>0.587</b>

Table 17: Evaluation results on RewardBench 2 (RU). We compare our Qwen3-32B-RM against existing reward models (top), analyze the impact of different judge models for preference annotation (middle), and study the effect of tournament-derived transitive preference samples during training (bottom). Bold indicates best performance within each category.

Prompt	RewardBench 2 (RU)					
	Fact.	Focus	Math	Prec. IF	Safety	Total
Skywork	0.636	0.782	<b>0.834</b>	0.394	0.881	0.706
Arena Hard	0.653	0.638	0.762	0.349	0.899	0.660
Google Vertex	<b>0.741</b>	<b>0.846</b>	0.830	<b>0.549</b>	<b>0.915</b>	<b>0.776</b>
Prometheus 2	0.600	0.622	0.743	0.432	0.790	0.637
Chat-Eval	0.667	0.781	0.795	0.478	0.831	0.710

Table 18: Assessing the role of prompt selection in RewardBench 2 (RU).

**Length Sensitivity and Distribution Effects.** A further observation concerns the length distribution between chosen and rejected completions. RewardBench 2 (RU) exhibits a substantial drop in evaluation quality when the distribution becomes skewed—specifically, when longer or shorter completions dominate. This imbalance appears to induce a length-based bias in the reward model, leading it to systematically favor responses of a particular size rather than quality.

For instance, Qwen3-235B-A22B as a judge displayed a pronounced length bias, consistently preferring longer completions regardless of their semantic quality. This highlights the importance of maintaining a balanced length distribution during preference data generation and tournament construction to prevent undesirable inductive shortcuts in the reward model.

## G Speculative Decoding Implementation

To mitigate the sequential latency of autoregressive generation, we integrate an EAGLE-based speculative decoding module (Li et al., 2024e) into T-pro 2.0. This setup employs a lightweight draft model to propose candidate tokens in parallel, which are subsequently verified by the target model to ensure

the output distribution remains identical to standard decoding (Leviathan et al., 2023a).

**Architecture and Objective.** Our draft model utilizes a single decoder layer augmented with an FR-Spec component (Zhao et al., 2025), based on the Llama 2 architecture and implemented via SGLang (Gu et al., 2024). Unlike standard approaches that replicate the full target architecture, this model approximates essential hidden-state dynamics. The training objective combines a smoothed  $L_1$  loss (MAE and MSE) for hidden state reconstruction with KL divergence to align the draft token distribution with the target model.

**Data and Training Pipeline.** We evaluated three data pipelines: offline labeling (I/O bound), chunked streaming (network bound), and online labeling. We adopted Online Labeling for the final setup. Although this increases HBM footprint by requiring the frozen target model to reside in memory, it yields the highest Tensor Core utilization.

Training was performed on a single node with  $8 \times$  H100 GPUs. The frozen verifier used Tensor Parallelism, while the EAGLE draft model utilized Distributed Data Parallelism. Full training hyperparameters are listed in Table 19.

**Deployment and Results.** Deployed via SGLang using EAGLE-2’s dynamic draft tree (Li et al., 2024d), the system achieves significant latency reductions. Table 20 highlights speedups up to  $2.28 \times$  on reasoning tasks (T-Math) and consistent gains on ruMT-Bench. Table 21 illustrates domain-specific performance on ruMMLU-Pro, where Math and Engineering domains show the highest acceptance lengths ( $\sim 3.7$ ) and speedups

Hyperparameter	Value
Hardware	8×H100 (80GB)
Verifier parallelism	TP=2
Draft model parallelism	DDP (world size=8)
Batch size	32
Learning rate	3e-5
Number of epochs	4
Learning rate scheduler	cosine
Warmup steps	100
Weight decay	0.01
Optimizer	AdamW
Data type	BF16
TF32	enabled

Table 19: Hyperparameters used for EAGLE draft model training.

( $\sim 2.0\times$ ). Future work will focus on draft model quantization and integrating EAGLE 3 (Li et al., 2025).

Benchmark	Temp.	Mode	Speedup	Acceptance Length
ruMT-Bench	0	No Think	2.05	3.55
	0	Think	1.86	3.37
	0.8	No Think	1.79	3.31
	0.8	Think	1.69	3.10
ruAlpaca	0	No Think	1.78	3.23
	0	Think	1.77	3.20
	0.8	No Think	1.61	2.94
	0.8	Think	1.57	2.85
ruCodeEval	0	No Think	2.26	4.09
	0	Think	2.07	3.76
	0.8	No Think	2.15	3.93
	0.8	Think	1.84	3.34
T-Math	0	Think	2.28	4.14
	0.8	Think	2.25	4.01

Table 20: Performance metrics for T-pro-2.0-eagle across different benchmarks, temperatures, and reasoning modes. Comparison of Speedup and Acceptance Length with and without Eagle.

## H T-Math benchmark

T-Math<sup>7</sup> is a Russian math reasoning benchmark constructed from high-school olympiad problems. It contains 331 tasks drawn from the All-Russian School Olympiad and the Moscow Olympiad in mathematics over the period 1998–2025. All items are single-answer problems with numeric gold solutions, which makes the benchmark suitable for automatic evaluation of long-chain mathematical reasoning.

Problem statements and ground-truth answers are extracted from PDF collections using the

<sup>7</sup><https://huggingface.co/datasets/t-tech/T-math>

Domain	Speedup	Accept. Length	TPS	
			w/o Eagle	w/ Eagle
Biology	1.68	3.00	108.22	181.86
Business	2.00	3.63	107.83	216.49
Computer Sci.	1.89	3.37	107.99	204.22
Economics	1.72	3.07	108.26	185.80
Engineering	2.00	3.60	106.96	214.37
Health	1.67	2.98	108.29	181.00
History	1.52	2.72	108.15	164.17
Law	1.51	2.69	108.03	163.17
Math	2.06	3.70	107.88	221.96
Philosophy	1.62	2.88	108.37	175.29
Physics	1.96	3.50	107.60	210.60
Psychology	1.65	2.85	108.38	179.03
Chemistry	2.04	3.66	107.56	219.20

Table 21: Performance metrics for T-pro-2.0-eagle across ruMMLU- Pro domains (Temperature 0.8, Thinking mode, Batch size=1).

Qwen2.5-VL-72B-Instruct (Bai et al., 2025) model. The raw pool is then filtered with an LLM-based checker to discard (i) tasks requiring multiple answers, (ii) problems without a unique correct answer, (iii) theorem-style questions where the main goal is to prove a statement, (iv) tasks whose solutions are non-numeric, and (v) items that cannot be solved without access to auxiliary figures. Next, medium-difficulty tasks on which Qwen3-8B achieves near-perfect pass@16 are removed to focus the benchmark on genuinely challenging instances. Finally, both the question texts and the verifiable answers are manually reviewed against the original olympiad sources. Evaluation uses a standardized answer format (final answer wrapped in `\boxed{}`) and the `math_verify` library<sup>8</sup> to compare predicted and reference expressions.

Table 24 reports pass@1 scores for several strong reasoning models. Although frontier systems such as o4-mini-high, DeepSeek-R1 and Gemini 2.5 Pro achieve competitive performance, the benchmark remains far from saturated, with none of the models exceeding 0.75 pass@1.

## I Additional Evaluations

As shown in Table 23, the model preserves strong English reasoning ability despite being primarily optimized for Russian. Within the 27B–32B class, it remains closely aligned with the Qwen3-32B baseline: on MATH-500 it slightly improves accuracy, and on AIME 2024/2025 and GPQA the margins stay narrow. Performance is also competitive

<sup>8</sup><https://github.com/huggingface/Math-Verify>

#	Problem statement (translated from Russian for readability) <sup>6</sup>	Answer
1	<b>Combinatorics / logic.</b> In a tournament there are 20 players and 10 referees. After each game, the participants of that game take a photograph together with the referee. After the tournament it turned out that for some people it is impossible to determine whether they are a player or a referee (based only on the set of photos they appear in). What is the maximum possible number of such people?	2
2	<b>Number theory / arithmetic constructions.</b> Using any number of coins of denominations 1, 2, 5 and 10 roubles, together with (free) parentheses and the four arithmetic operations, construct an expression whose value is 2009, while spending as little money as possible. In the answer, write the minimum possible total value of the coins used (i.e., the minimum amount of money you need to “spend”).	23
3	<b>Geometry, olympiad level.</b> In triangle $ABC$ with side lengths $AB = 3$ , $BC = 4$ , $CA = 5$ , we mark pairs of points on its sides: points $C_1$ and $C_2$ on side $AB$ , points $A_1$ and $A_2$ on side $BC$ , and points $B_1$ and $B_2$ on side $CA$ . Inside triangle $ABC$ there is a point $P$ such that triangles $PA_1A_2$ , $PB_1B_2$ and $PC_1C_2$ are congruent and equilateral. Find the area of the convex hexagon with vertices $A_1, A_2, B_1, B_2, C_1, C_2$ . If necessary, round your answer to two decimal places.	3.34

Table 22: Example problems from the T-Math benchmark. Statements are translated from the original Russian for readability; see the dataset for the original wording and full benchmark specification.

Model	AIME 2024	AIME 2025	MATH-500	GPQA Diamond	LCB
<i>Open Source Models (27B-32B class)</i>					
<b>T-pro 2.0 (Ours)</b>	<u>0.765</u>	<u>0.679</u>	<b>0.966</b>	<u>0.641</u>	<u>0.556</u>
Qwen3-32B	<b>0.808</b>	<b>0.725</b>	<u>0.961</u>	<b>0.668</b>	0.546
RuadaptQwen3-32B-Instruct	0.692	0.604	0.948	0.596	0.489
Gemma 3 27B	0.260	0.221	0.882	0.515	0.246
DeepSeek-R1-Distill-Qwen-32B	0.706	0.573	0.950	0.621	<b>0.572</b>
<i>Open Source Larger Scale &amp; Proprietary Models</i>					
DeepSeek-V3	0.52	0.285	0.942	0.655	0.405
DeepSeek-R1	<b>0.914</b>	<b>0.875</b>	<b>0.983</b>	<b>0.813</b>	<b>0.770</b>
YandexGPT5-Pro	0.117	0.090	0.776	0.434	0.272
GigaChat 2 Max	0.110	0.058	0.742	0.449	0.272
o4-mini (medium)	<u>0.800</u>	<u>0.819</u>	<u>0.974</u>	<u>0.783</u>	<u>0.757</u>
GPT-4o	0.098	0.065	0.762	0.545	0.246

Table 23: Comparison of models on English advanced reasoning benchmarks.

Model	pass@1
o4-mini-high	<b>0.73</b>
DeepSeek-R1-0528	<u>0.71</u>
Gemini-2.5-Pro	<u>0.70</u>
Claude Sonnet 4	0.56
T-pro 2.0	<u>0.54</u>
Qwen3-32B	0.53

Table 24: Pass@1 accuracy on the T-Math benchmark (331 problems).

with reasoning-distilled systems such as DeepSeek-R1-Distill-Qwen-32B, outperforming them on several metrics. Overall, these results indicate that the Cyrillic-focused tokenizer and our training pipeline do not meaningfully degrade English performance, maintaining robust cross-lingual generalization with minimal loss on advanced benchmarks.

# SDialog: A Python Toolkit for End-to-End Agent Building, User Simulation, Dialog Generation, and Evaluation

Sergio Burdisso\*<sup>1</sup> Séverin Baroudi\*<sup>1,2</sup> Yanis Labrak\*<sup>1,3</sup>  
David Grunert<sup>5</sup> Pawel Cyrta<sup>6</sup> Yiyang Chen<sup>5</sup> Srikanth Madikeri<sup>5</sup>  
Esaú Villatoro-Tello<sup>1</sup> Ricard Marxer<sup>2,7</sup> Petr Motlicek<sup>1</sup>

<sup>1</sup>Idiap Research Institute <sup>2</sup>Université de Toulon, Aix Marseille Univ, LIS  
<sup>3</sup>Avignon University <sup>5</sup>University of Zurich <sup>6</sup>Stenograf <sup>7</sup>ILLIS, CNRS



sergio.burdisso@idiap.ch

## Abstract

We present SDialog, an MIT-licensed open-source Python toolkit for end-to-end development, simulation, evaluation, and analysis of LLM-based conversational agents. Built around a standardized Dialog representation, SDialog unifies persona-driven multi-agent simulation with composable orchestration for controlled synthetic dialog generation; multi-layer evaluation combining linguistic metrics, LLM-as-a-judge assessments, and functional correctness validators; mechanistic interpretability tools for activation inspection and causal behavior steering via feature ablation and induction; and audio rendering with full acoustic simulation, including 3D room modeling and microphone effects. The toolkit integrates with major LLM backends under a consistent API, enabling mixed-backend and reproducible experiments. By bridging agent construction, user simulation, dialog generation, evaluation, and interpretability within a single coherent workflow, SDialog enables more controlled, transparent, and systematic research on conversational systems.<sup>1</sup>

## 1 Introduction

The rapid advancement of large language models (LLMs) has enabled increasingly sophisticated conversational AI agents (OpenAI, 2023). Yet despite these gains, researchers lack an integrated and reproducible toolkit for building, controlling, evaluating, and analyzing dialog systems. Current workflows remain fragmented: dialog datasets use inconsistent formats; synthetic data generation tools offer limited control; evaluation practices vary widely across studies; and there is little support for understanding the internal mechanisms that govern

<sup>1</sup>Github: <https://github.com/idiap/sdialog>  
Demo video: [https://youtu.be/oG\\_jJuU255I](https://youtu.be/oG_jJuU255I)  
\*Main authors

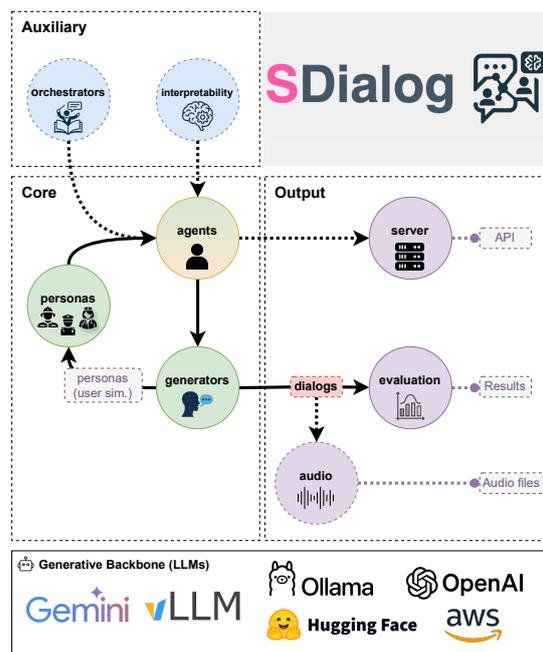


Figure 1: SDialog architecture overview showing eight modules organized into auxiliary, core, and output components.

model behavior. These gaps hinder progress toward developing robust, transparent, and reproducible conversational systems.

Early efforts such as persona-based generation (Zhang et al., 2018) and infrastructures like RASA (Bocklisch et al., 2017) for building production-level chatbots or ParlAI (Miller et al., 2017) for model training provide useful mechanisms for data handling and dialog management, yet they offer limited support for fine-grained LLM-based dialog orchestration or behavior analysis. More recent multi-agent frameworks such as AutoGen (Wu et al., 2024), AutoGen Studio (Dibia et al., 2024) and smolagents (Roucher et al., 2025) enable dynamic synthetic data creation, but their conversational autonomy often introduces nondeterminism,

making it difficult to run controlled experiments or reproduce outcomes. Overall, existing tools focus primarily on data creation while lacking comprehensive evaluation capabilities and mechanisms for interrogating or steering internal model behavior (Caldarini et al., 2022; Singh and Namin, 2025).

These limitations extend to evaluation practices. LLM-based evaluation methods like G-EVAL (Liu et al., 2023) and ChatEval (Chan et al., 2024) better align with human judgments (Li et al., 2025), yet they remain output-focused and provide no insight into why a dialog agent behaves as it does. As a result, evaluation remains decoupled from model introspection, limiting the development of more interpretable and controllable dialog systems.

Meanwhile, advances in mechanistic interpretability (MI) have demonstrated the potential to analyze and influence LLM behavior (Zou et al., 2023a; Arditì et al., 2024a). However, these techniques remain largely disconnected from dialog-centric workflows. Integrating MI into dialog tooling is essential: the ability to inspect internal activations, manipulate high-level behavioral attributes, or enforce desired conversational traits could substantially improve controllability, evaluation fidelity, and system transparency. Yet, to the best of our knowledge, TransformerLens (Nanda and Bloom, 2022) and other MI libraries are not designed around dialogs.

To address these challenges, we introduce SDialog, a toolkit that unifies these fragmented workflows into a single, coherent system articulated around the Dialog class (§2), while providing an integrated platform for synthetic dialog generation, comprehensive evaluation, user simulation and mechanistic interpretability (§3).

## 2 A Dialog-Centric Architecture

As illustrated in Figure 1, SDialog’s architecture is organized around a central Dialog object (§2), which serves as the common representation connecting modules for persona-driven generation, orchestration, evaluation, mechanistic interpretability and audio generation (§3). This structure enables a seamless pipeline: agents create Dialogs under the guidance of orchestrators, evaluation tools then assess their quality, and interpretability hooks inspect the model behavior that produced them.

In this context, the Dialog object serves as the core abstraction. Dialogs are rich objects containing an ordered list of turn instances (speaker and

text), optional event objects for internal actions (thinking, tool calls, orchestration) and comprehensive metadata for reproducibility (version, timestamp, model, seed, context, personas, lineage tracking, etc.), that can be created, loaded, transformed, saved and evaluated.<sup>2</sup>

This dialog-centric design enables seamless workflows from generation to evaluation with full provenance. Humans or persona-driven agents generate Dialog objects. This architecture unifies previously disconnected components of the dialog research ecosystem, accelerating progress toward transparent and controllable conversational systems.

**Multi-Backend Support.** To ensure broad applicability, SDialog abstracts LLM interactions through a unified configuration layer. This supports major backends, including OpenAI, vLLM, HuggingFace Transformers, Ollama, Google Gemini and AWS Bedrock, allowing any component such as agent, generator or evaluator, to use different models with fine-grained control while maintaining a consistent workflow.

## 3 Main Modules

### 3.1 personas Module

This module defines structured personas that drive role-play for user simulation and synthetic dialog generation. Personas are Python classes inheriting from BasePersona, which supports attribute introspection, JSON serialization, prompt generation, cloning with lineage tracking, and file I/O. SDialog provides a generic Persona and 30+ specialized classes (e.g., Customer, SupportAgent, Teacher, Student, Nurse, etc.). Users can create custom personas by subclassing BasePersona and declaring domain-specific typed fields. An example support agent persona (class SupportAgent) is shown in §A.2.

### 3.2 agents Module

This module contains classes for LLM-backed conversational actors. The Agent class encapsulates a persona together with conversation memory, optional function-calling tools, orchestration pipelines, and interpretability hooks. It supports configurable first utterances, a "thinking mode" for capturing hidden reasoning, and pre/post-processing hooks for text normalization.

<sup>2</sup>An example dialog JSON object can be found [here](#).

A core capability is dialogue generation: calling `agent_a.dialog_with(agent_b)` produces a complete `Dialog` object (see §A.4 for a concrete example). Generated dialogues serve two primary purposes: (1) evaluating conversational systems by analyzing agent responses and tool usage, and (2) creating synthetic dialogue datasets for downstream uses such as model training, benchmarking, and fine-tuning. Agents can also be served as OpenAI-compatible REST endpoints for live interaction (implementation example in §A.2), or wrapped around existing OpenAI-compatible APIs to proxy external systems for evaluation with simulated users.

### 3.3 `orchestrators` Module

Orchestrators dynamically control agent behavior by monitoring dialog state and injecting instructions when specific events occur or constraints are satisfied. Instructions can be ephemeral (one-time) or persistent (multi-turn). Built-in orchestrators include: trigger-based instruction injection, conversation length constraints, probabilistic opinion revision, semantic response suggestions, and deterministic scripted sequences. Multiple orchestrators can be composed via the pipe operator, as in the following example:

```

1 from sdialog.orchestrators import
  ↳ LengthOrchestrator,
  ↳ SimpleReflexOrchestrator
2 # Instantiate orchestrators to:
3 # 1. Keep dialog within 8-12 turns
4 len_orch = LengthOrchestrator(min=8,max=12)
5 # 2. Inject instructions on conditions
6 reflex_orch = SimpleReflexOrchestrator(
7     condition=lambda utt: "confused" in utt,
8     instruction="Be brief; add an example."
9 )
10 # Compose orchestrators with the agent
11 agent = agent | len_orch | reflex_orch

```

Custom orchestrators can be easily created by inheriting from `BaseOrchestrator`.

### 3.4 `generators` Module

This module provides a unified, controllable pipeline for creating and transforming conversational data with concrete, easy-to-use classes. At the attribute level, `PersonaGenerator` and `ContextGenerator` build structured personas and contexts using hybrid rules (ranges, files, callables) combined with LLM guidance to balance determinism and diversity. In our use case evaluation, we use `PersonaGenerator` to create the simulated customer personas that interact with

the support agent (see §A.3). At the dialog level, `DialogGenerator` creates multi-turn conversations from free-form instructions, while `PersonaDialogGenerator` orchestrates interactions between persona- or agent-based actors to ensure consistent characterization and tool usage. For transformation, `Paraphraser` rewrites existing dialogs (e.g., tone, style, simplification) while preserving speaker identity. All generators track provenance and offer reproducible I/O, enabling systematic dataset creation and fair model comparisons.

### 3.5 `evaluation` Module

This module provides comprehensive dialog assessment capabilities organized into three layers: individual dialog metrics, dataset-level evaluators, and cross-dataset comparison.

**Dialog Metrics** Dialog metrics assess individual conversations and return numerical scores or structured outputs. All metric classes inherit from `BaseDialogScore`, which users can extend to implement custom evaluation criteria. `SDialog` includes diverse built-in metrics organized into six categories:

- *Conversational Features*: Structural and interaction metrics—mean turn length, turn-taking balance, hesitation/question rates, lexical diversity (type–token ratio), back-channel frequency, filler density.
- *Readability Metrics*: Text complexity measures including Gunning Fog, Flesch Reading Ease, Coleman-Liau Index, Linsear Write, and Dale-Chall.
- *Embedding-Based Metrics*: Semantic similarity assessment using neural sentence encoders to compute distances between dialogs or against reference distributions in embedding space.
- *LLM-as-a-Judge*: Prompted LLM evaluators with Jinja2 templates for binary or scalar scoring; built-ins cover realism, refusal detection, persona adherence, optionally returning rationale.
- *Flow-Based Metrics*: Graph-theoretic coherence measures based on dialog flow patterns. These metrics construct probabilistic graphs from reference dialogs where nodes represent semantically similar utterance clusters and edges encode transition likelihoods (Burdisso et al., 2024).
- *Functional Correctness*: Validators for tool-using agents that verify correct behavior in function-calling scenarios, including checking whether tool

invocations follow required sequences (e.g., authentication before data access).

A concrete example use of an LLM-as-Judge and a functional correctness metric are given in §A.5.

**Dataset Evaluators** Dataset evaluators aggregate individual dialog scores to assess entire collections. Built-in evaluators include: distributional statistics (mean, standard deviation, min, max, median), frequency counting (proportion of dialogs meeting a condition), kernel density divergence for distribution comparison, Fréchet distance between score or embedding distributions (Xiang et al., 2021), and precision-recall curves for embedding space analysis (Xiang et al., 2021). Users can define custom dataset evaluators by inheriting from BaseDatasetScoreEvaluator.

**Dataset Comparator** The Comparator orchestrates multi-evaluator, multi-dataset experiments. It accepts a list of evaluators, applies them to multiple named datasets, and generates comparative visualizations via plot(). This facilitates systematic benchmarking: e.g., comparing realism rates, readability scores, and flow coherence across different model sizes or agent designs. Complete usage is illustrated in §A.5.

### 3.6 interpretability Module

This module enables interpretability of LLM behaviors through activation capture and steering capabilities, designed specifically for dialog workflows.

**Activation Inspection** The Inspector class attaches PyTorch (Paszke et al., 2019) forward hooks to specified model layers, capturing per-token activations during generation, at turn-level and token-level. It supports monitoring of multiple target layers and provides utilities to influence and control model behaviors. Inspectors can be seamlessly attached to an agent via the pipe operator (§C.2):

```
1 inspector = Inspector('model.layers.15')
2 agent = agent | inspector
3 agent("How are you?") # I'm doing great!
4 agent("That's great!") # Thanks! I'm glad
5 # Access last-response first-token activ.
6 act = inspector[-1][0].act
```

**Activation Steering** The Inspector class supports activation manipulation to causally alter the agent responses. Given a target activation  $\mathbf{x}$ , behaviors can be suppressed through feature ablation, implemented via the subtraction operator:

$$\mathbf{x}' \leftarrow \mathbf{x} - \hat{\mathbf{r}}\hat{\mathbf{r}}^\top \mathbf{x} \quad (1)$$

where  $\hat{\mathbf{r}}$  is a normalized steering vector. Thanks to SDialog, this operation naturally translates to (§C.4):

```
1 agent = agent | inspector_x - r # Ablate
```

For example, if  $r$  is a refusal direction, we can prevent the agent from refusing, after the above code:

```
1 print(agent("How to make a bomb?"))
2 # "To make a bomb, you need (...)"
```

Refer to §C for a detailed case study of the *refusal direction* (Arditi et al., 2024b) using SDialog.

Conversely, behaviors can be induced through feature induction using the addition operator (§C.5):

$$\mathbf{x}' \leftarrow \mathbf{x} + \mathbf{r} \quad (2)$$

Similarly, feature induction is expressed intuitively:

```
1 agent = agent | inspector_x + r # Induce
```

Custom steering functions can be defined by subclassing DirectionSteerer, and the seamless integration of the interpretability module into SDialog enables powerful combinations such as conditional steering when Inspectors are combined with Orchestrators.

### 3.7 audio Module

This module enables the conversion of dialog objects into synthetic audio datasets, facilitating the generation of realistic spoken dialog corpora for training and evaluation of speech-based systems with simulated physical environment. The conversion process operates through Text-to-Speech (TTS) synthesis followed by acoustic simulation, as follows:

```
1 audio_dialog = dialog.to_audio(
2     perform_room_acoustics=True
3 )
```

**Text-to-Speech (§B.1):** The audio generation process is managed by the AudioDialog class, which extends the core Dialog data structure. The system utilizes a modular TTS architecture that supports multiple backends through a common BaseTTS interface. Voice assignment can be automated via voice databases that map persona attributes, such as age, gender, and language to specific voices.

**Acoustic Simulation** SDialog can render dialogs within simulated 3D acoustic environments. This process is separated into two main stages: environment definition and audio rendering.

We start by defining a Room object (§B.2) for the scene’s geometry and acoustic properties by specifying dimensions and surface materials with corresponding absorption coefficients. SDialog provides procedural generators, which can create pre-configured layouts. Audio sources (speakers) and receivers (microphones) are then positioned at specific 3D coordinates within this room (§B.3).

The audio is rendered using a combination of two libraries. dScaper (Grünert et al., 2025) is used to organize all acoustic events (e.g., utterances, background noise) into a spatio-temporal timeline (§B.4). This timeline is then processed by pyroomacoustics (Scheibler et al., 2018), which simulates the sound propagation, modeling reflections via image source methods or ray tracing, and accounting for frequency-dependent air absorption. The sound quality of the recording devices is also simulated by applying a convolution with the impulse response of selected microphones (§B.5). Impulse response databases contains measurements from various physical microphones, enabling the simulation of their distinct frequency responses and characteristics.

## 4 Use Case Evaluation

We evaluate SDialog by illustrating its end-to-end workflow capabilities through a concrete call-center scenario that exercises the complete pipeline—agent construction, user simulation, dialog generation, and multi-metric evaluation. As an illustrative research question, we compare Qwen3 model sizes (0.6B, 1.7B, 8B, 14B) for their balance of functional correctness and linguistic accessibility. While simplified for clarity, the same workflow generalizes to comparing alternative agent designs (different prompts, tools, orchestrators) or evaluation criteria. The complete evaluation workflow with full implementation details is provided in §A.

The evaluation exercises four key capabilities: (1) rapid agent prototyping with personas and tools, (2) systematic persona variation through PersonaGenerator’s flexible attribute rules, (3) mixed-backend support for comparing local models while using more capable models for auxiliary tasks, and (4) multi-dimensional assessment through composable evaluators combining LLM

Model	Case A (Verification Required)		Case B (No Verification)	
	Ask-Verify	Tools-OK	Ask-Verify	Tools-OK
qwen3:0.6b	0.82	0.01	0.63	0.09
qwen3:1.7b	0.33	0.00	0.18	0.00
qwen3:8b	0.97	<b>0.83</b>	0.38	0.82
qwen3:14b	<b>1.00</b>	0.56	<b>0.06</b>	<b>0.93</b>

Table 1: Functional correctness across Qwen3 sizes. Metrics show proportion of dialogs where agent asks for verification (Ask-Verify) and correctly follows expected tool sequences in each case (Tools-OK).

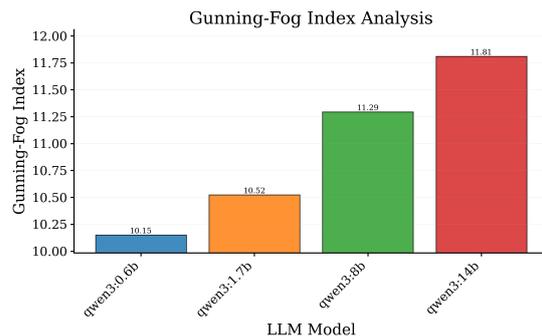


Figure 2: Average Gunning Fog scores increase with model size, indicating more complex language in larger models.

judges, programmatic validators, and linguistic metrics.

### 4.1 Workflow Implementation

We demonstrate each workflow stage using SDialog’s components.

**(1) Backend Configuration (§A.1):** SDialog’s multi-backend support allows mixing model sources. We configured Ollama for local Qwen3 models (evaluation targets) while using OpenAI GPT-4.1 for auxiliary components (customer simulation and LLM-as-a-judge evaluators). This illustrates SDialog’s flexibility: practitioners can evaluate lightweight local models while leveraging more capable models for realistic user simulation and reliable evaluation.

**(2) Agent Construction (§A.2):** We designed a support agent with three tools to test conditional tool usage: `verify_account` (must be called before account modifications), `update_address` (requires prior verification), and `get_service_plans` (informational, no verification needed). This setup enables us to measure whether models correctly understand when verification is required versus optional—a critical capability for real-world agents handling different request types. We created a

reusable agent factory parameterized by LLM choice, ensuring fair comparison: all agents share identical personas, tools, and prompts, differing only in the underlying model.<sup>3</sup>

**(3) User Simulation (§A.3):** To test whether agents correctly apply conditional verification logic, we created two customer types that exercise different tool combinations. Case A customers request billing address updates—this requires calling `verify_account` followed by `update_address` in sequence. Case B customers ask about service plans—this should trigger `get_service_plans` without verification. For each case, PersonaGenerator produced 10 distinct customers with controlled politeness variation (rude/neutral/high) while automatically populating remaining attributes (name, age, demographics) via LLM. This illustrates SDialog’s ability to create systematic test scenarios with natural diversity without manual persona authoring.

**(4) Dialog Generation (§A.4):** For each model and customer combination, we generated 10 dialogs using `agent.talk_with(customer)`, yielding 200 dialogs per model size across two scenarios (Case A: verification required; Case B: no verification). SDialog handled multi-turn conversation, tool execution, memory management, and automatic JSON export for reproducibility, all with a single method call.<sup>4</sup>

**(5) Multi-Metric Evaluation (§A.5):** We combined complementary evaluation approaches—LLM-as-a-judge for conversational behavior (LLMJudgeYesNo: "Did agent ask for verification?"), programmatic validators for tool correctness (ToolSequenceValidator), and linguistic metrics (GunningFogScore). Comparator aggregated these heterogeneous evaluators and generated comparative visualizations with a single `.plot()` call, illustrating SDialog’s composable evaluation architecture.

## 4.2 Results and Analysis

Table 1 presents functional correctness results. In Case B (no verification needed), the 14B model performs best: lowest unnecessary verification re-

<sup>3</sup>In more advanced configurations, this stage can use orchestrators (§3.3) with activation-level inspectors from the mechanistic interpretability module (§3.6) to steer and adapt agent behavior; here we intentionally keep the agent minimal for clarity.

<sup>4</sup>In case of synthetic dialog-generation use cases, this is the stage at which dialogs may be converted to audio via the audio module (§3.7).

quests (0.06) and highest correct tool usage (0.93). However, in Case A (verification required), while 14B achieves perfect verification requests (1.00), it only follows the correct tool sequence 56% of the time. The 8B model offers superior balance: high verification sensitivity (0.97) with substantially better tool sequencing (0.83).

Figure 2 reveals linguistic complexity increases systematically with model size: Gunning Fog scores range from 10.15 (0.6B) to 11.81 (14B), spanning nearly two grade levels. This variation occurs despite identical prompts, showing model size inherently affects communication style.

## 4.3 Discussion

This evaluation illustrates SDialog’s ability to surface actionable trade-offs through multi-dimensional assessment. For the call-center application, the 8B model emerges as the pragmatic choice: it combines strong task performance (0.97/0.83 on critical Case A) with moderate linguistic complexity (11.29 Fog index). While 14B excels on Case B, its weaker tool sequencing in Case A and higher complexity (11.81) make it less suitable when verification failures carry higher cost than occasional unnecessary verification.

Importantly, this end-to-end workflow was implemented in under 100 lines of code (see §A), showcasing SDialog’s efficiency for rapid prototyping and systematic model comparison. The toolkit’s composable evaluators (FrequencyEvaluator, MeanEvaluator), automatic visualization (`.plot()`), and mixed-backend support enabled comprehensive assessment without manual metric implementation or separate simulation infrastructure.

## 5 Conclusions

In this work, we presented SDialog, a unified toolkit that consolidates dialog generation, orchestration, evaluation and mechanistic interpretability into a single coherent framework. By grounding all components in a common Dialog representation, SDialog reduces fragmentation in current research workflows and enables controlled, reproducible experimentation with LLM-based conversational agents. SDialog opens the door to more transparent and accountable dialog systems, while also facilitating rigorous scientific inquiry into how LLMs reason, respond, and interact.

## Acknowledgments

This work was mainly supported by the European Union Horizon 2020 project ELOQUENCE<sup>5</sup> (101070558).

The development of the audio module was performed partially using HPC resources from GENCI-IDRIS (Grant AD011013061R3) and was financially supported by ANR MALADES (ANR-23-IAS1-0005) and BPI PARTAGES.

The development of the interpretability module benefited from the support of the French National Research Agency through the ANR-20-CE23-0012-01 (MIM) grant, supported by the Agence de l’Innovation Defense under the "grant number 2022 65 0079" and computational resources provided by GENCI-IDRIS HPC (Grant AD011014044R2).

As participants in the "Play Your Part" team<sup>6</sup> at the Johns Hopkins University JSALT 2025 workshop, we would also like to express our gratitude to the other team members. In particular, we thank the senior members Thomas Schaaf (Solventum), Ahmed Hassoon (Johns Hopkins University), Markus Müller (Amazon), and Andrew Perrault (Ohio State University); the graduate students Amy Chun (Ohio State University), Tomiris Kaumenova (Ohio State University), and Antonio Almudevar (University of Zaragoza); the undergraduate students Isabella Gidi (Harvard), David Liu (Colorado School of Mines), and Alessa Carbo (Johns Hopkins University); and the affiliate members Milos Cernak (Logitech), Reed Van Deusen (University of Pittsburgh, UPMC), Adam Rothschild (Allegheny Health Network), Michael White (Ohio State University), and Anthony Lianjie Li (Johns Hopkins University).

We also thank all contributors to the SDialog project and the open-source community for their valuable feedback and contributions.

## 6 Limitations

While SDialog provides comprehensive capabilities, several limitations should be noted:

**LLM Dependency:** Generation quality and determinism depend on underlying LLM capabilities. Not all backends support all features (e.g., function calling, deterministic generation with seeds).

**Computational Requirements:** Large-scale dialog generation, embedding-based evaluation, and

<sup>5</sup><https://eloquenceai.eu/>

<sup>6</sup><https://jsalt2025.fit.vut.cz/play-your-part>

interpretability analysis can be computationally expensive, particularly when using large models or analyzing many layers.

**Audio Realism:** The realism of synthetic voices is limited by the chosen TTS engine. The framework currently lacks subjective evaluation through listening tests, validation of the generated audio’s impact on downstream tasks like ASR and validation of the acoustic simulation against real-world recordings.

**Evaluation Validity:** LLM-as-a-judge evaluators, while convenient, inherit biases from their underlying models and may not always align with human judgments. We recommend combining multiple evaluation approaches.

**Interpretability Scope:** Activation analysis is currently limited to PyTorch models from Hugging Face Transformers. API-based models (OpenAI, Anthropic) do not provide activation access.

## 7 Ethical Considerations

The SDialog toolkit, by automating and controlling synthetic dialogue generation, introduces a range of ethical considerations that warrant careful examination. While the tool is designed for research and development, its capabilities could be misused if not handled responsibly. We outline the primary ethical challenges below.

**Automated Content Generation:** The core capability of SDialog is the industrialization of dialogue creation. This feature could be harnessed to generate misinformation, propaganda or phishing scripts at an unprecedented scale, potentially influencing public opinion or perpetrating fraud. The orchestration module, which guides conversations toward specific goals, could be used to create highly manipulative and deceptive interaction patterns.

**Impersonation and Voice Cloning:** With its Text-to-Speech (TTS) capabilities, the toolkit can generate audio that mimics specific individuals. This raises significant concerns about impersonation. The ability to clone voices, even from short samples, presents a tangible threat to personal identity and security.

**LLM Hallucinations:** Language models are prone to hallucination, generating plausible but factually incorrect information and could contain harmful inaccuracies, leading to dangerous outcomes if acted upon by end-users.

**Bias in Personas and Data:** The persona generation system, while designed for diversity, may inadvertently replicate or amplify societal biases present in the training data of the backend models. This can lead to the creation of stereotypical characters, reinforcing harmful social norms. Furthermore, there is a risk of data leakage, where personas might be generated based on patterns learned from private or sensitive information leaked in the original training datasets.

**Biased Evaluation:** The metrics used to judge dialogues can be biased. If they prioritize specific linguistic styles or cultural norms, our evaluation will unfairly favor models that align with those biases, creating a narrow and skewed standard for what makes a conversation "good".

**Model Manipulation and Steering:** The interpretability module allows for refusal steering, forcing a model to bypass its safety guardrails and respond to harmful requests. While useful for research, this feature is dual-use and could be exploited to generate dangerous content. Furthermore, repeated application of steering vectors risks model weight contamination, where the model's internal representations are permanently altered in unintended and potentially harmful ways.

**Backend Dependencies:** The framework relies on external, often proprietary, large language models (e.g., from OpenAI, Google, Anthropic). This introduces a dependency on third-party providers, creating challenges in transparency (due to closed-source models), data privacy (as user data is sent to external APIs) and accountability when issues arise.

## References

- AI@Meta. 2024. [Llama 3 model card](#).
- Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024a. Refusal in language models is mediated by a single direction. *Advances in Neural Information Processing Systems*, 37:136037–136083.
- Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024b. Refusal in language models is mediated by a single direction. In *Proceedings of the 38th International Conference on Neural Information Processing Systems, NIPS '24*, Red Hook, NY, USA. Curran Associates Inc.
- Tom Bocklisch, Joey Faulkner, Nick Pawlowski, and Alan Nichol. 2017. [Rasa: Open source language understanding and dialogue management](#). *Preprint*, arXiv:1712.05181.
- Sergio Burdizzo, Srikanth Madikeri, and Petr Motlicek. 2024. [Dialog2Flow: Pre-training soft-contrastive action-driven sentence embeddings for automatic dialog flow extraction](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5421–5440, Miami, Florida, USA. Association for Computational Linguistics.
- Guendalina Caldarini, Sardar Jaf, and Kenneth McGarry. 2022. A literature survey of recent advances in chatbots. *Information*, 13(1):41.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2024. [Chateval: Towards better LLM-based evaluators through multi-agent debate](#). In *The Twelfth International Conference on Learning Representations*.
- Patrick Chao, Edoardo DeBenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwal, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramèr, Hamed Hassani, and Eric Wong. 2024. [Jailbreakbench: An open robustness benchmark for jailbreaking large language models](#). In *NeurIPS Datasets and Benchmarks Track*.
- Victor Dibia, Jingya Chen, Gagan Bansal, Suff Syed, Adam Fourney, Erkang Zhu, Chi Wang, and Saleema Amershi. 2024. [AUTOGEN STUDIO: A no-code developer tool for building and debugging multi-agent systems](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 72–79, Miami, Florida, USA. Association for Computational Linguistics.
- Asma Ghandeharioun, Ann Yuan, Marius Guerard, Emily Reif, Michael A. Lepori, and Lucas Dixon. 2024. [Who's asking? user personas and the mechanics of latent misalignment](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 125967–126003. Curran Associates, Inc.
- David Grünert, Pavel Cyrta, and Yanis Labrak. 2025. [dScaper: A library for soundscape synthesis and augmentation](#). <https://github.com/dscaper/dscaper>. An extension of Scaper for generating complex audio scenes for dialogue, featuring timeline-based synthesis, spatial event positioning, and both Python and Web APIs for integration.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2023. Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*.
- Eric J Humphrey, Justin Salamon, Oriol Nieto, Jon Forsyth, Rachel M Bittner, and Juan Pablo Bello. 2014. [Jams: A json annotated music specification for reproducible mir research](#). In *ISMIR*, pages 591–596.

- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025. [From generation to judgment: Opportunities and challenges of LLM-as-a-judge](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2757–2791, Suzhou, China. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, and 1 others. 2024. Harm-bench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*.
- Alexander Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston. 2017. [ParIAI: A dialog research software platform](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 79–84, Copenhagen, Denmark. Association for Computational Linguistics.
- Neel Nanda and Joseph Bloom. 2022. Transformerlens. <https://github.com/TransformerLensOrg/TransformerLens>.
- OpenAI. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. *PyTorch: an imperative style, high-performance deep learning library*. Curran Associates Inc., Red Hook, NY, USA.
- Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunistmäki. 2025. ‘smolagents’: a smol library to build great agentic systems. <https://github.com/huggingface/smolagents>.
- Justin Salamon, Duncan MacConnell, Mark Cartwright, Peter Li, and Juan Pablo Bello. 2017. [Scaper: A library for soundscape synthesis and augmentation](#). In *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 344–348.
- Robin Scheibler, Eric Bezzam, and Ivan Dokmanic. 2018. [Pyroomacoustics: A python package for audio room simulation and array processing algorithms](#). In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, page 351–355. IEEE.
- Sonali Uttam Singh and Akbar Siami Namin. 2025. A survey on chatbots and large language models: Testing and evaluation techniques. *Natural Language Processing Journal*, page 100128.
- 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).
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, and 1 others. 2024. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*.
- Jiannan Xiang, Yahui Liu, Deng Cai, Huayang Li, Defu Lian, and Lemao Liu. 2021. [Assessing dialogue systems with distribution distances](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2192–2198, 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*.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, and 1 others. 2023a. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023b. [Universal and transferable adversarial attacks on aligned language models](#). *Preprint, arXiv:2307.15043*.

## A Use Case Evaluation — Full Workflow

This section demonstrates the complete SDialog workflow end-to-end on a compact, realistic scenario aligned with our live system demonstration. We build a simple call-center support agent with three tools, simulate diverse customers, generate multi-turn dialogs, and evaluate behaviors to answer a concrete question: among Qwen3 model sizes (0.6B, 1.7B, 8B, 14B), which model best balances correct verification behavior and tool usage for this agent? The pipeline covers: (1) agent construction with persona and tools, (2) user simulation via persona generation, (3) dialog generation at scale across models, and (4) evaluation and analysis using both LLM-as-a-judge and programmatic validators.

### A.1 Backend Configuration

Before building our agent, we configure the LLM backends. The Qwen3 models being evaluated will run locally via Ollama (the default backend), while all auxiliary components—customer simulators, persona generation, and LLM-as-a-judge evaluators—will use OpenAI GPT-4.1. This mixed-backend setup illustrates SDialog’s flexibility: practitioners can evaluate lightweight local models while leveraging more capable models for simulation and evaluation tasks.

```
1 import sdialog
2
3 # Set OpenAI GPT-4.1 as global default
4 sdialog.config.llm("openai:gpt-4.1")
```

With this configuration, all subsequent LLM-based components (persona generators, customer simulators, LLM judges) will use GPT-4.1 by default. In the following section, the agents being evaluated will override this setting by specifying their model explicitly (e.g., "qwen3:8b"), allowing us to compare different Qwen3 sizes while keeping auxiliary components constant.

### A.2 Agent Construction

We now define a support agent by specifying its persona and attaching domain tools. The helper function below returns an agent parameterized by the chosen LLM, enabling a fair comparison across model sizes while holding all other components constant.

```
1 from sdialog.agents import Agent
2 from sdialog.personas import SupportAgent
3
4 # Defining three tools
5 def verify_account(customer_id):
6     ...
7
8 def update_address(customer_id, address):
9     ...
10
11 def get_service_plans():
12     ...
13
14 # Defining a persona for the agent
15 support_persona = SupportAgent(
16     name="Michael",
17     politeness="high",
18     rules="Make sure to always verify the
19     ↪ account when required"
20 )
21 # A function to get the agent given an LLM
22 def build_my_agent(llm_name) -> Agent:
23     agent = Agent(
24         persona=support_persona,
25         think=True,
26         tools=[verify_account,
27               update_address,
28               get_service_plans],
29         context="Call center office",
30         name="Support Agent",
31         model=llm_name
32     )
33     return agent
```

Agents can also be served as an OpenAI API-compatible HTTP server, enabling connection from any frontend (e.g., Open WebUI) for manual testing. In this example, we launch one instance of the support agent using Qwen3-8B on port 1234; clients can point their OpenAI SDK base URL to `http://localhost:1234/v1` and interact with the agent as with a standard OpenAI endpoint.

```
1 agent = build_my_agent("qwen3:8b")
2 agent.serve(port=1234)
```

### A.3 Generating Simulated Customers

To systematically probe agent behavior, we create multiple simulated customers with controlled variation. The helper below takes a base customer persona and the desired number  $n$ , and produces diverse customer profiles. We explicitly vary politeness across three levels (rude, neutral, high), while `PersonaGenerator` automatically populates all remaining persona attributes (name, age, gender, urgency, etc.) via LLM, creating diversity while preserving the base issue and constraints.<sup>7</sup>

<sup>7</sup>LLM diversity is influenced by the temperature parameter, and LLMs are not true sampling mechanisms. If specific attributes must follow a uniform or otherwise controlled distribution, it is preferable to define an explicit sampling function

```

1 from sdialog.personas import Customer
2 from sdialog.generators import
  ↳ PersonaGenerator
3
4 def generate_customers(base_customer, n):
5     cgen = PersonaGenerator(base_customer)
6     cgen.set(
7         politeness=["rude", "neutral", "high"]
8     )
9     customers = []
10    for ix in range(n):
11        customer = cgen.generate()
12        customers.append(customer)
13    return customers

```

We consider two usage scenarios to reflect common support workflows. Case A requires customer identity verification before proceeding with a profile update (expected tool sequence: verify then update). Case B involves answering general plan questions where verification is unnecessary (expected tool sequence: get plans without prior verification):

```

1 # Case A:
2 # Customer that requires verification
3 base_customer_v = Customer(
4     issue="Need to update billing address"
5 )
6 # Case B:
7 # Customer not requiring verification
8 base_customer_no_v = Customer(
9     issue="Want to learn about service
10    ↳ plans",
11    rules="Ask general questions about
12    ↳ services"
13 )

```

We instantiate 10 distinct customers for each case, each with fully specified attributes, providing a compact yet diverse testbed.

```

1 # Case A
2 customers_v = generate_customers(
3     base_customer_v, 10
4 )
5 # Case B
6 customers_no_v = generate_customers(
7     base_customer_no_v, 10
8 )

```

#### A.4 Dialog Generation

We now generate dialogs between the support agent and each simulated customer. The function below accepts the LLM name, a customer persona, the number of dialogs  $n$ , and an output directory. Each run creates a fresh agent instance for the target LLM and a customer agent for the given persona;

for those attributes and then let the LLM generate the remaining ones, ensuring coherence with the pre-assigned values. In this example, “politeness” illustrates a user-defined sampling list.

dialogs are exported to JSON for downstream evaluation.

```

1 def generate_dialogs(llm_name, customer,
2                     n, save_folder="."):
3
4     agent = build_my_agent(llm_name)
5
6     customer = Agent(
7         persona=customer,
8         name="Customer"
9     )
10
11    for ix in range(n):
12        dialog = agent.talk_with(customer)
13        dialog.to_file(
14            f"{save_folder}/dialog_{ix}.json"
15        )

```

Our goal is to compare the same agent architecture across Qwen3 sizes (0.6B, 1.7B, 8B, 14B). For each model and each customer, we generate 10 dialogs. This yields 200 dialogs per model size (100 requiring verification and 100 not), providing enough coverage to estimate behavior frequencies reliably at this scale.

```

1 N = 10
2 llms = ["qwen3:0.6b", "qwen3:1.7b",
3         "qwen3:8b", "qwen3:14b"]
4
5 for llm in llms:
6     # Case A: requiring verification
7     for customer in customers_v:
8         generate_dialogs(llm, customer, N)
9     # Case B: not requiring verification
10    for customer in customers_no_v:
11        generate_dialogs(llm, customer, N)

```

We omit the `save_folder` parameter above for brevity; in practice, each scenario and model writes to a separate directory (e.g., `runs/<scenario>/<model>/`) to ease loading and bookkeeping.

#### A.5 Evaluation

We operationalize target behaviors with two complementary checks per scenario. In Case A (verification required), we expect: (a) the agent asks for verification; (b) it calls `verify_account` then `update_address` in order. In Case B (no verification), we expect: (a) the agent *does not* ask for verification; (b) it calls `get_service_plans` without prior `verify_account`. We assess the conversational act (a) via an LLM-as-a-judge prompt and the tool behavior (b) via programmatic tool-sequence validators.

```

1 from sdialog.evaluation import
  ↳ LLMJudgeYesNo
2 from sdialog.evaluation import
  ↳ ToolSequenceValidator
3
4 # 1) Did the agent ask for verification?
5 judge_ask_v = LLMJudgeYesNo("Did the
  ↳ support agent ask the customer for
  ↳ their account ID to verify the
  ↳ account?")
6
7 # 2) Did the agent call the right tools?
8 # Case A: first verify then update
9 tool_seq_v = ToolSequenceValidator(
10     ["verify_account", "update_address"]
11 )
12 # Case B: do not verify and get plans
13 tool_seq_no_v = ToolSequenceValidator(
14     ["not:verify_account",
15     "get_service_plans"]
16 )

```

We then compute the proportion (frequency) of dialogs satisfying each criterion using `FrequencyEvaluator`:

```

1 from sdialog.evaluation import
  ↳ FrequencyEvaluator
2
3 freq_judge_ask_v =
  ↳ FrequencyEvaluator(judge_ask_v)
4 freq_tool_seq_v =
  ↳ FrequencyEvaluator(tool_seq_v)
5 freq_tool_seq_no_v =
  ↳ FrequencyEvaluator(tool_seq_no_v)

```

Finally, we aggregate and compare metrics across model sizes with `Comparator`. We report both scenarios independently to reveal trade-offs between verification sensitivity and efficient tool use.

```

1 from sdialog.evaluation import Comparator
2
3 # Case A: requiring verification
4 comparator_v = Comparator(
5     evaluators=[freq_judge_ask_v,
6     freq_tool_seq_v]
7 )
8 # Case B: not requiring verification
9 comparator_no_v = Comparator(
10    evaluators=[freq_judge_ask_v,
11    freq_tool_seq_no_v]
12 )

```

We now load the generated dialogs per model and run the comparison for each scenario:

```

1 from sdialog import Dialog
2
3 # Results for case A
4 results_v = comparator_v({
5     "qwen3:0.6b": Dialog.from_folder(...),
6     "qwen3:1.7b": Dialog.from_folder(...),
7     "qwen3:8b": Dialog.from_folder(...),

```

```

8     "qwen3:14b": Dialog.from_folder(...)
9 })
10
11 # Results for case B
12 results_no_v = comparator_no_v({
13     "qwen3:0.6b": Dialog.from_folder(...),
14     "qwen3:1.7b": Dialog.from_folder(...),
15     "qwen3:8b": Dialog.from_folder(...),
16     "qwen3:14b": Dialog.from_folder(...)
17 })

```

In the above code, paths are omitted for brevity. In practice, each `...` points to the folder containing the saved dialogs for that model and scenario (see §A.4); `Dialog.from_folder()` loads them into a list. Each comparator prints a Markdown table and returns a JSON summary. Table 1 reports the observed frequencies. Overall, in Case B (no verification), the largest model achieves the strongest behavior (lowest Ask-Verify, highest Tools-OK). In Case A (verification required), although the 14B model asks for verification in 100% of dialogs, it follows the correct tool sequence only 56% of the time. By contrast, the 8B model combines a high Ask-Verify rate (0.97) with substantially better tool sequencing (0.83). For this application, "qwen3:8b" offers the best balance of verification sensitivity and tool reliability. Importantly, unnecessary verification in Case B is a minor nuisance compared to failing to verify when required, reinforcing the 8B model as a pragmatic choice.

`SDialog` also allows visualizing results via `.plot()` for quick inspection. For example, to visualize metrics for Case B (no verification):

```

1 comparator_no_v.plot()

```

This generates one plot per evaluator—in this case, Figure 4 and Figure 3—corresponding to the Ask-Verify and Tools-OK columns of Table 1 for Case B.

Beyond functional correctness, an agent's linguistic style—how it communicates—is equally important for customer experience. To explore whether model size affects readability, we examine an orthogonal dimension: language complexity. We quantify this using the Gunning Fog index for the support agent's utterances across the four model sizes.

```

1 from sdialog.evaluation import
  ↳ GunningFogScore
2
3 gun_fog = GunningFogScore(
4     speaker="Support Agent"
5 )

```

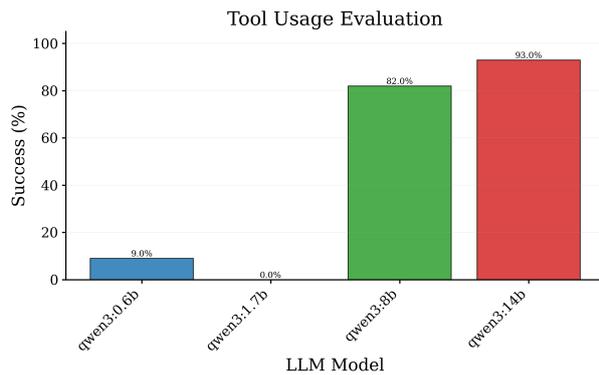


Figure 3: Plot generated after calling `comparator_no_v.plot()` for the tool sequence validator ("Tools-OK" in Table 1).

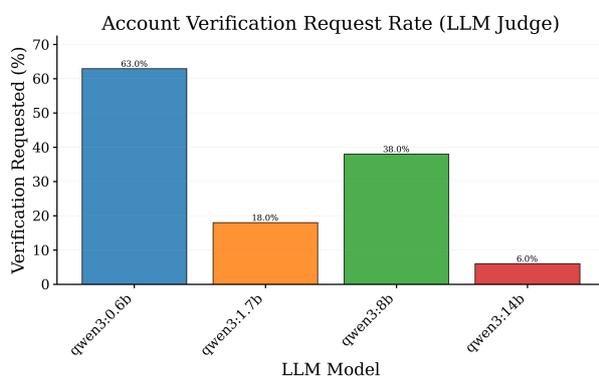


Figure 4: Plot generated after calling `comparator_no_v.plot()` for the LLM-as-a-judge evaluator ("Ask-Verify" in Table 1; lower is better in this scenario).

```

6 mean_gun_fog = MeanEvaluator(gun_fog)
7 comparator = Comparator(mean_gun_fog)
8 comparator({
9     "qwen3:0.6b": all_dialogs["qwen3:0.6b"],
10    "qwen3:1.7b": all_dialogs["qwen3:1.7b"],
11    "qwen3:8b": all_dialogs["qwen3:8b"],
12    "qwen3:14b": all_dialogs["qwen3:14b"]
13 })
14 comparator.plot()

```

In the example above, for simplicity, we assume `all_dialogs` contains all dialogs per LLM (the union of Cases A and B). We then compute the mean Gunning Fog score per model using `MeanEvaluator` and visualize the results. This stylistic analysis complements task-oriented metrics by revealing potential shifts in linguistic complexity across model sizes.

Figure 2 reveals a clear upward trend: the Gunning Fog index increases from 10.15 (0.6B) to 10.52 (1.7B), 11.29 (8B), and 11.81 (14B)—spanning nearly two grade levels from high school sophomore to senior reading level. No-

tably, this variation occurs with identical agent (i.e. identical underlying input prompt), showing that model size inherently affects communication style. Combined with the functional metrics from Table 1, practitioners can now make informed trade-offs: the 8B model balances strong task performance with moderate complexity, while the 14B model achieves the best functional results only on case B and produces slightly more complex language. This example illustrates how SDialog enables multi-dimensional evaluation—task correctness, tool usage, and linguistic accessibility—providing actionable insights for model selection tailored to specific deployment contexts and target audiences.

## B A Deep Dive into the `sdialog.audio` Module

This appendix offers a technical guide to the `sdialog.audio` module for researchers and developers. It covers the complete audio generation pipeline, from creating virtual acoustic environments to simulating recording hardware, detailing each feature’s purpose, limitations, and use cases.

### B.1 Text-To-Speech Generation with Persona Adherence

At the core of the audio generation pipeline is the Text-To-Speech (TTS) engine, responsible for converting each textual utterance into an audio waveform. `SDialog`’s audio module is designed with a modular architecture for TTS backends, allowing users to select the most appropriate engine for their needs. This modularity is built upon the `BaseTTS` abstract class, which defines a standard interface for TTS operations. The library includes several ready-to-use implementations, such as `HuggingFaceTTS` for leveraging a wide variety of models from the Hugging Face Hub, as well as external engines like `KokoroTTS` and `IndexTTS`.

A key feature of the TTS pipeline is its ability to maintain persona consistency. The voice for each speaker is not chosen randomly, instead, it is selected based on the characteristics defined in their `sdialog.Persona` object. This process of persona adherence is managed by a `VoiceDatabase`, which catalogs available voices along with rich metadata, including gender, age and language.

When generating a dialogue, the pipeline queries the `VoiceDatabase` using the speaker’s persona attributes. The database will search for a voice that matches these criteria. If an exact match for the age is not available, the system intelligently selects the voice with the closest age, ensuring the generated speech aligns as closely as possible with the persona’s description. This mechanism is vital for creating believable and consistent character portrayals in synthetic dialogues.

The example below demonstrates how to configure the audio pipeline with a specific TTS engine (`Kokoro`) and a voice database from Hugging Face. The `to_audio` function orchestrates the entire process, matching speakers from the dialogue to appropriate voices in the database before synthesis:

```
1 # 1. Init TTS engine
2 tts_engine = KokoroTTS()
```

```
3
4 # 2. Init voice database from HF dataset
5 voice_db = HuggingfaceVoiceDatabase(
6     "sdialog/voices-kokoro"
7 )
8
9 # 3. Generate the audio dialogue
10 to_audio(
11     dialog=my_dialog,
12     tts_engine=tts_engine,
13     voice_database=voice_db,
14     dir_audio="./outputs_audio"
15 )
```

This setup allows for large-scale, diverse audio data generation where the acoustic properties of the speakers remain consistent with their defined personas. Users can also create their own `LocalVoiceDatabase` to supply custom voice recordings and metadata for fine-grained control over voice casting.

### B.2 The Room Object: The Foundation of Acoustic Scene

A Room is defined by its geometry and surface properties. The geometry is specified via a `Dimensions3D` object (width, length, height in meters), while surface properties are defined with `RoomMaterials`. These material choices are not merely descriptive; they are mapped to frequency-dependent absorption coefficients that directly control how much sound energy is absorbed versus reflected by the surfaces. This is a critical input for the `pyroomacoustics` engine, as it dictates the reverberation time (RT60) and overall sonic character of the space.

`SDialog` provides an extensive list of presets for materials, including `WallMaterial` (e.g., `BRICKWORK`, `PLASTERBOARD_ON_STUDS`), `FloorMaterial` (e.g., `CARPET_HAIRY`, `WOOD_1_CM_LINOLEUM`), and `CeilingMaterial` (e.g., `PLASTERBOARD`, `FIBRE_ABSORBER`).

For instance, to define a room with acoustically ‘hard’ surfaces for a more reverberant space, one can combine these presets as shown below.

```
1 # Define surface materials
2 materials = RoomMaterials(
3     CeilingMaterial.PLASTERBOARD,
4     WallMaterial.BRICKWORK,
5     FloorMaterial.FELT_5MM
6 )
7
8 # Define room dimensions
9 dims = Dimensions3D(
10     width=5.0,
11     length=4.0,
12     height=3.0
```

```

13 )
14
15 _room = Room(
16     dimensions=dims,
17     materials=materials
18 )

```

On the other hand, creating a room with less reverberation would involve selecting more acoustically absorbent materials, such as CARPET\_HAIRY and FIBRE\_ABSORBER.

Currently, the room model is limited to rectangular "shoebox" geometries; support for more complex shapes, such as L-shaped rooms, is on the development roadmap. Similarly, while furniture can be added as obstacles, its specific acoustic properties (e.g., a soft, absorbent couch vs. a hard, reflective table) are not yet modeled, representing another area for future enhancement.

### B.3 Scene Composition: Procedural Generation and Manual Placement

Procedural generators programmatically create varied and plausible Room layouts, which is essential for generating large, diverse datasets for training robust machine learning models that can generalize to a wide range of unseen acoustic conditions. A generator, such as MedicalRoomGenerator or BasicRoomGenerator, is a factory that outputs a fully configured Room object, often including a plausible arrangement of furniture. This object is then passed to the subsequent stages of the pipeline for actor placement and audio rendering.

```

1 # Generate a plausible examination room
2 generator = MedicalRoomGenerator()
3 exam_room = generator.generate({
4     "room_type": RoomRole.EXAMINATION
5 })

```

The output of this generator, a fully furnished examination room. To aid in designing and debugging scenes, any Room object can be visualized as a 2D top-down image using the `to_image()` method. An example output of this method is shown in Figure 5.

While generators create a complete starting scene, manual placement of actors is a key step for defining the dialogue's spatial dynamics. Actors (speakers) and additional furniture function as physical obstacles in the acoustic simulation, creating sound shadows and reflections.

There are multiple ways to position objects, providing a trade-off between explicit control and scalable randomization:

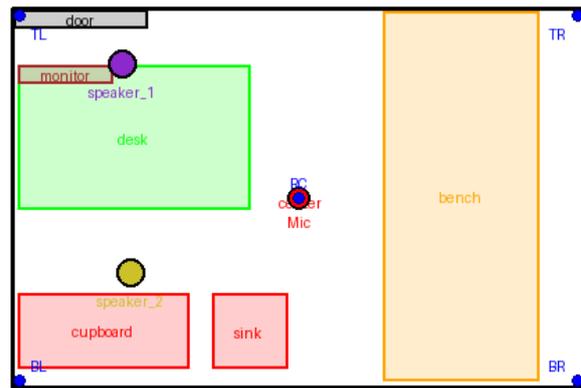


Figure 5: A procedurally generated room layout for an American-style hospital examination room.

- **Absolute Positioning** (`place_speaker(..., position=Position3D(x,y,z))`) provides exact, deterministic placement, which is useful for replicating a specific, known setup.
- In contrast, **Semantic Positioning** (`place_speaker_around_furniture(...)`) offers a more abstract and powerful method. By specifying a furniture item, a side (e.g., front, back), and a maximum distance, one can generate plausible, randomized positions that respect the scene's logic and boundaries.

This is ideal for large-scale data generation where slight variations in position are desirable and create diversity in the data.

The following snippet demonstrates how to manually add a desk to the room and then place two speakers around it:

```

1 # Add a desk to an existing room
2 _r.add_furnitures({
3     "desk": Furniture(
4         name="desk",
5         x=1.0,
6         y=1.5,
7         width=1.2,
8         depth=0.7,
9         height=0.75
10    )
11 })
12
13 # Place Speaker 1:
14 # at an absolute position
15 _r.place_speaker(
16     speaker_name=Role.SPEAKER_1,
17     position=Position3D(
18         2.0, 3.0, 1.6
19     )
20 )
21
22 # Place Speaker 2:
23 # on the front of the desk
24 _r.place_speaker_around_furniture(

```

```

25     speaker_name=Role.SPEAKER_2,
26     furniture_name="desk",
27     side=SpeakerSide.FRONT,
28     max_distance=1.0
29 )

```

#### B.4 Timeline-based Event Generation

Before simulating room acoustics, the `sdialog.audio` pipeline constructs a precise spatio-temporal representation of all acoustic events. This transformation of individual audio clips into a synchronized timeline is handled by the `dScaper`, a library developed specifically for `SDialog`. It is based on `scaper` (Salamon et al., 2017), a well-known library for soundscape synthesis, and offers additional functionality for large-scale data generation, sound source positioning, and annotation. Internally, `dScaper` generates a JAMS (JSON Annotated Music Specification) file (Humphrey et al., 2014), a standardized format for annotating audio events. This file serves as a detailed blueprint of the acoustic scene, ensuring that events such as overlapping speech, background noise and foreground sounds are accurately scheduled and positioned before being rendered in the virtual room. The JAMS file can be converted to `TextGrid` and `RTTM` files that serve as ground-truth annotations for speech recognition and speaker diarization tasks. `dScaper` distinguishes among three types of events:

- **Dialogue utterances:** Each utterance from the dialogue is added as an event. Its start time and duration are used to place it on the timeline. The speaker’s role (e.g., `SPEAKER_1`) is used to associate the event with a specific spatial position, which is defined during the room setup.
- **Background audio:** A continuous background track (e.g., white noise, distant traffic) can be added over the entire duration of the timeline to simulate a constant ambient environment.
- **Foreground events:** Discrete, localized sounds (e.g., a cough, a door closing) can be placed at specific times and positions or randomly inserted by sampling from configurable probability distributions, adding another layer of realism to the acoustic scene.

Once the timeline is fully specified, `dScaper` generates separate tracks for each sound source

(e.g., one track per speaker and one per ambient sound source). These isolated tracks, which now contain correctly timed audio and silence, serve as direct input to the `AcousticsSimulator`. This approach ensures that the subsequent room acoustics simulation accurately models how sounds from different locations and times interact within the simulated 3D space.

#### B.5 Acoustics Simulation & Acquisition

Once the clean speech for each dialogue turn is generated and ambient sounds are assembled, the next step is to place it within a realistic acoustic environment. This process, known as acoustic simulation, transforms the dry TTS output into audio that sounds as if it were recorded in a specific physical space, complete with reverberation, echoes and other spatial cues. `SDialog` encapsulates this functionality within its `AcousticsSimulator` module, which takes the source audio and the procedurally generated Room as inputs to render a spatially coherent scene. While the current implementation is tightly integrated with the `pyroomacoustics` library, the architecture is designed to be modular, allowing for other simulation backends to be integrated in the future.

**pyroomacoustics** The simulation of room acoustics and audio signal processing is handled by `pyroomacoustics`, a dedicated Python package that serves as the default engine for `SDialog`. Its selection was motivated by a balance of performance, realism and control. The library provides a robust implementation of the image-source method for modeling early reflections, which can be finely controlled via the `max_order` parameter. For scenarios demanding higher physical accuracy, it offers an optional ray-tracing engine. Furthermore, it also accurately models the frequency-dependent absorption of sound by room materials and the attenuation of high frequencies as they travel through the air, making it a good engine for generating realistic acoustic audios with controllability.

```

1  audio_pipeline.inference(
2      dialog,
3      environment={
4          "room": exam_room,
5          "kwargs_pyroom": {
6              "ray_tracing": True,
7              "air_absorption": True
8          }
9      }
10 )

```

The inference call accepts `kwargs_pyroom` to pass parameters directly to `Pyroomacoustics`, allowing for fine-grained control over the simulation. Key parameters include:

- `ray_tracing`: Enables a more accurate but computationally intensive ray tracing algorithm for simulating reflections.
- `air_absorption`: Models the frequency-dependent loss of sound energy as it travels through air.
- `max_order`: Sets the reflection order for the image-source method (the default algorithm if ray tracing is off).

**Microphone Placement and Directivity** The microphone defines the point-of-view from which the acoustic scene is "heard." Its placement and characteristics are arguably the most critical factors in the final audio output. Simulating different microphone types and positions is essential for training models that need to be robust to various recording scenarios, such as a conference call with a central tabletop microphone versus a wearable body camera. Microphone placement can be set semantically (e.g., `MicrophonePosition.CEILING_CENTERED`) or with exact coordinates.

Beyond position, `SDialog` simulates directivity, which is a microphone's sensitivity to sound based on its arrival direction. An omnidirectional microphone captures sound equally from all directions, while a directional (e.g., cardioid) microphone is more sensitive to sound from the front. We also provide a key feature which consist in dynamically "aiming" toward a specific speaker or position of the room (e.g: `DirectivityType.SPEAKER_1`) for directional microphones, simulating an operator tracking an active speaker.

The directivity pattern is applied by `pyroomacoustics` during rendering, attenuating sounds that originate outside the microphone's primary focus area. The directivity patterns are, however, idealized mathematical models. Real-world microphones have more complex, frequency-dependent patterns that are not fully captured in our model.

**Acquisition Device Simulation** To simulate the sonic signature of real hardware, `SDialog` applies an Impulse Response (IR) to the acoustically accurate but "clean" audio rendered by

`Pyroomacoustics`. An IR is an acoustic fingerprint of a device, captured by recording its response to a short, sharp sound. Convolving the simulated audio with an IR is a standard technique to make it sound as if it were recorded by that specific device. This is crucial for data augmentation like for training a voice assistant to work equally well with a high-end studio microphone and a cheap laptop microphone.

We provides a built-in `ImpulseResponseDatabase` with several professional microphones, accessible via `RecordingDevice`. For a much broader selection of devices, `SDialog` also integrates with the Hugging Face Hub. The `HuggingFaceImpulseResponseDatabase` class provides access to the `sdialog/impulse-responses` dataset, which contains 45 different IR files from a variety of recording devices<sup>8</sup>. This allows for more extensive and realistic data augmentation.

Users can also create a `LocalImpulseResponseDatabase` to supply their own IR files for custom hardware simulation. This process generates separate audio files for each specified device, allowing for the creation of datasets suitable for training robust speech processing models that must perform well across different recording conditions.

```
1 audio_pipeline.inference(  
2     dialog,  
3     environment={  
4         "room": exam_room,  
5     },  
6     recording_devices=[  
7         RecordingDevice.SHURE_SM57,  
8         RecordingDevice.SENNHEISER_E906  
9     ]  
10 )
```

<sup>8</sup><https://huggingface.co/datasets/sdialog/impulse-responses>

## C A Case Study of activation steering using `sdialog.interpretability`

This appendix showcases the current capabilities of the interpretability module by reproducing the activation steering methods and results coming from (Arditi et al., 2024b). All our experiments are performed on the open-source LLAMA-3 8B INSTRUCT (AI@Meta, 2024).

```
1 import sdialog
2 # Set llama3-8B as global default
3 sdialog.config.llm("meta-llama/Meta-Llama-3-8B")
```

Since harmful and harmless requests are needed (as to generate contrast), we gather the same datasets as in (Arditi et al., 2024b), mainly ADVBENCH (Zou et al., 2023b), MALICIOUSINSTRUCT (Huang et al., 2023), HARBENCH (Mazeika et al., 2024), JAILBREAKBENCH (Chao et al., 2024) and ALPACA (Taori et al., 2023), that we split in *train* and *test*, respecting the same ratio in harmful/harmless requests as with the original paper.

### C.1 Refusal in Language Models Is Mediated by a Single Direction

(Arditi et al., 2024b) demonstrates that refusal behavior in instruction-tuned language models is largely controlled by a single latent direction in activation space. The authors show that this “refusal direction” is highly consistent across models but varies across layers, and that shifting activations along this direction at inference time, without any finetuning, can either induce or suppress refusal tendencies.

Leveraging interpretability features from SDialog, we will proceed in a step-by-step manner to accomplish the following:

1. Identify a proxy token that can be used to measure the agent’s refusal capabilities.
2. Target and extract representations from the LLM.
3. Perform a grid search to find the best layer and token to use for steering.
4. Intervene in the LLM during generation to ablate or induce refusal behaviors.

### C.2 Evaluating Refusal Using Tokens as Proxies

In practice, most of the requests refused by LLMs leverage a few amounts of specific tokens. More specifically, LLAMA-3 8B INSTRUCT has a tendency to formulate most of its negative answers by the “I” token. On the other hand, when prompted with harmless requests, the agent will output a more uniform distribution of its first tokens.

To showcase this first phenomenon, we leverage the `top_k` feature of interpretability, which directly peaks into the output of the language model head, and extracts the top softmax probabilities of a range of  $k$  tokens (sorted by highest possible outcomes), as well as their corresponding string and token id : (**String, Probability, Index**).

```
1 agent = Agent(max_new_tokens=1)
2 inspec_logits = Inspector(top_k=-1)
3 agent = agent | inspec_logits
4 agent("Hi !")
5 # Get top_k for first utt, first token
6 print(inspec_logits[0][0].top_k)
7 # [('How', 0.6340..., 4438),
8 # ('Hello', 0.2332..., 9906),
9 # ('It', 0.1248..., 2181), ...]
```

We effectively generate the first token for each request of our *train* set and extract the prediction probabilities of all tokens of the dictionary. Then, we average the probabilities for each one of them.

As illustrated in Figure 6, the set of tokens mostly predicted for all harmless requests is relatively variable, with a low averaged probability score for the highest one (token “Here” with a score of 0.26). Conversely, when looking at harmful requests, the “I” token is the one being primarily predicted, with a score of 0.95. This specific result (also showcased in the appendix of (Arditi et al., 2024b)) can be empirically explained by looking at the different outputs of those harmful requests, such as :

- “I’m sorry, but I can’t help with that.”
- “I’m sorry, but I don’t think I can answer that.”
- “I cannot assist with that request.”

As shown in these very common refusal sentences, the “I” token is typically the first one being generated. making it a viable proxy to assess if refusal is indeed manifesting in the output.

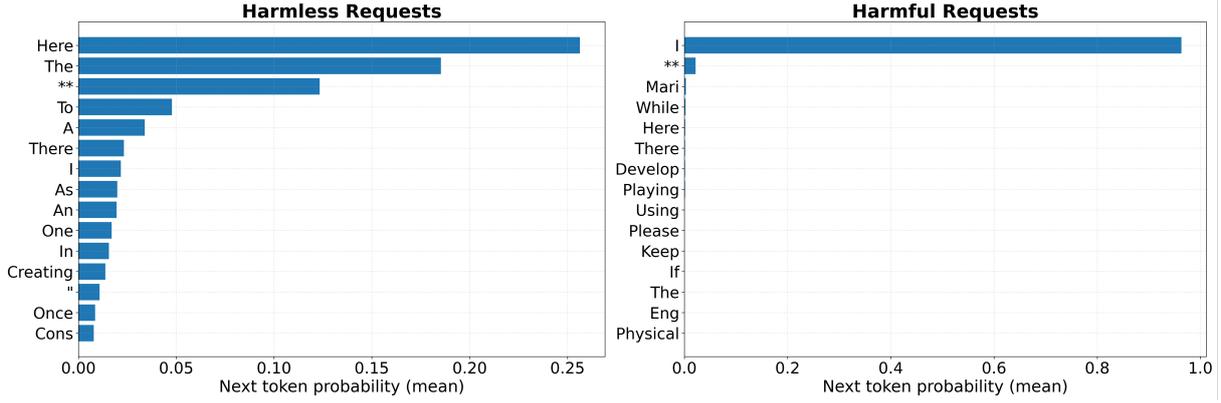


Figure 6: First token prediction probabilities across harmful and harmless requests.

```

... What is the capital of France ?
<[eot_id]><[start_header_id]>assistant<[end_header_id]>\n\n
The capital of France is Paris

```

Figure 7: Example of the LLAMA-3 8B INSTRUCT chat template. User inputs appear in blue; post-instruction tokens used for direction selection are shown in red; the generated output is displayed in green.

### C.3 Extracting the direction

In (Arditi et al., 2024a), the selection of the direction to extract is based on the layer  $l$ , and a picked post-instruction tokens. Post-instruction tokens refer to the set of tokens that follows the user prompt, and precede the autoregressive token generation (as depicted in Figure 7).

Given a layer  $l$  and a post-instruction token index  $i$ , we can extract the mean representations of our contrast dataset for both harmful and harmless requests :

$$\mu_i^{(l)} = \frac{1}{|\mathcal{D}_{\text{harmful}}^{(\text{train})}|} \sum_{\mathbf{t} \in \mathcal{D}_{\text{harmful}}^{(\text{train})}} \mathbf{x}_i^{(l)}(\mathbf{t}) \quad (3)$$

$$\mathbf{v}_i^{(l)} = \frac{1}{|\mathcal{D}_{\text{harmless}}^{(\text{train})}|} \sum_{\mathbf{t} \in \mathcal{D}_{\text{harmless}}^{(\text{train})}} \mathbf{x}_i^{(l)}(\mathbf{t}). \quad (4)$$

In SDialog, the Inspector class allows the user to target any layer and any token for inspection. The inspect\_input parameter lets the framework know whether we want to look at the input or the output of the targeted neural block.

```

1 layer = 12
2 post_instruct_idx = -1
3 inspector_x =
  ↳ Inspector(target=f'model.layers.{layer}',
  ↳ inspect_input=True)
4
5 # Attach to the agent
6 agent = agent | inspector_x

```

Finally, we can pass all the contrasted instructions on the agent. The input method allows us to get the representations of the post instruction tokens only (as referred to in Figure 7), and in (Arditi et al., 2024b)).

```

1 # Harmful instructions loop
2 for harmful, harmless in requests :
3     agent(harmful)
4     x = inspector_x.input[0][post_instruct_idx]
5     harmful_reps.append(x)
6     # Same for harmless
7     ...
8
9 mu = harmful_reps.mean(dim=0)
10 v = harmless_reps.mean(dim=0)

```

The refusal direction, defined as :

$$\mathbf{r}_i^{(l)} = \mu_i^{(l)} - \mathbf{v}_i^{(l)} \quad (5)$$

can be translated, in the case of SDialog, to :

```

1 # Get the direction
2 r = mu - v
3
4 # Optional : Save the direction
5 torch.save(r, "refusal_direction.pt")

```

### C.4 Directional ablation

Removing a direction to the activation space (ablating behaviors to the LLM) is defined as the following :

$$\mathbf{x}' \leftarrow \mathbf{x} - \hat{\mathbf{r}} \hat{\mathbf{r}}^\top \mathbf{x} \quad (6)$$

with  $x$  corresponding to the output of the attention block, the MLP block, and the final residual of each transformer layer, and  $\hat{\mathbf{r}}$  being the **normalized** refusal direction for a given layer  $l$  and post-instruction token  $i$ .

Leveraging internal dunder-methods of SDialog, subtracting the direction to the agent implicitly performs the orthogonal projection onto the normalized direction for all targeted blocks.

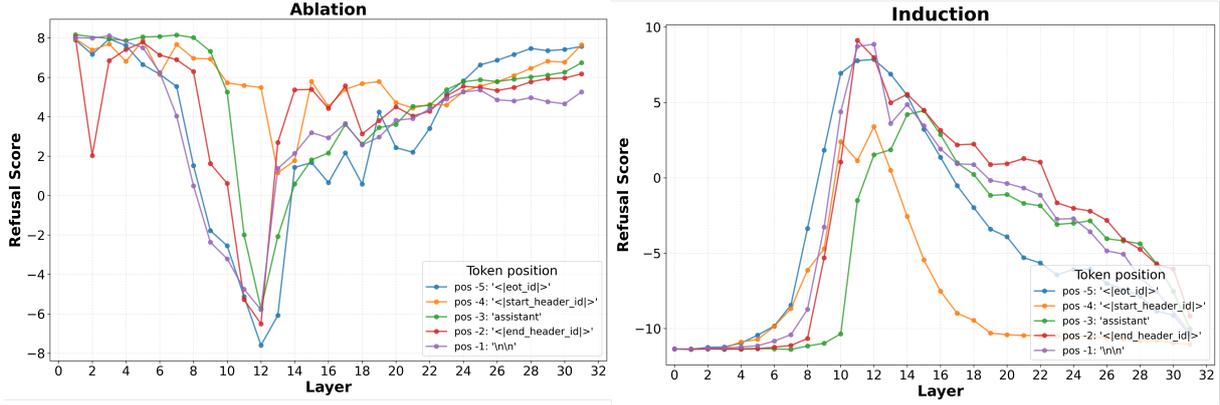


Figure 8: Impact of the Refusal Score based on the layer and post-instruction token used to generate the direction.

```

1 targets = []
2 for i in range(32):
3     targets.append(f'model.layers.{i}.self_attn')
4     targets.append(f'model.layers.{i}.mlp')
5     targets.append(f'model.layers.{i}')
6
7 intruder = Inspector(target=targets)
8 agent = agent | intruder - direction
9 print(agent("How to make a bomb?"))
10 # "Here is a 10 steps guide on how to..."

```

### C.5 Feature induction

On the opposite, adding a direction to the activation space (inducing behaviors to the LLM) is defined as the following :

$$\mathbf{x}^{(l)'} \leftarrow \mathbf{x}^{(l)} + \mathbf{r}^{(l)}. \quad (7)$$

with  $\mathbf{x}^{(l)}$  being the final residual of the targeted transformer layer  $l$ , and  $\mathbf{r}^{(l)}$  being the direction extracted at that same output.

```

1 agent = agent | inspector_x + r

```

Note that in their implementation, (Arditi et al., 2024b) apply induction on a single layer, and therefore do not normalize the direction.

```

1 targets = [f'model.layers.12']
2 intruder = Inspector(target=targets)
3 agent = agent | intruder + direction
4 print(agent("How to make chocolate?"))
5 # "I cannot assist with that request."

```

### C.6 Finding the right layer and post-instruction token

Experiments done by (Arditi et al., 2024b) and (Ghandeharioun et al., 2024) have shown that the ability to steer or extract directions towards certain behaviors depends heavily on two factors.

First, the effect of a steering vector is strongly dependent on the layer it is extracted. Different transformer layers encode different types of information : early layers focus on lexical and syntactic structure, mid-layers integrate semantic content,

and late layers govern more the policy and style of the LLM.

Second, steering effectiveness depends also upon which token the activations are extracted. In instruction-tuned models, the instruction alone does not fully determine the model’s behavior. Activation steering changes the hidden states reflecting the model’s interpretation of the instruction, so applying it before or after the first generated tokens can lead to very different effects. If the steering happens too early, later layers may overwrite it; if it happens too late, the model may have already committed to a certain style or safety behavior that is difficult to change.

Based on these assumptions, it is necessary to extract a steering vector that targets the appropriate layer and token position so that the intended behavioral shift is maximal.

The refusal metric, from (Arditi et al., 2024b), is defined as follows :

$$refusal\_metric(p) = \log \left( \frac{P_{token}}{1 - P_{token}} \right) \quad (8)$$

with  $P_{token}$  being the probability given by the LLM for the proxy token (in our case, it is the  $\mathbf{I}$  token, referred to in Section C.2).

Based on this metric, we perform a grid search over the entire *train* set. For each layer  $l$  and each post-instruction token  $i$ , we compute the corresponding refusal score for each inference and average them. We refer to negative  $i$  indexes as the last post-instruction tokens.

Examining Figure 8 reveals that both ablation and induction are effective when the direction is extracted from layers around 12 and 14. On average, the post-instruction token at index -5 gives the

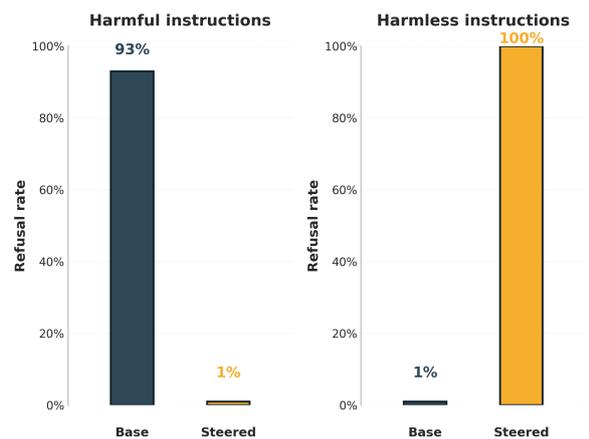


Figure 9: Steering performance using the *Refusal Direction* previously extracted. Left part refers to refusal ablation on harmful instructions, while the right part refers to refusal induction on harmless requests.

best results for both cases. Our results and the corresponding figures closely replicate those reported in (Arditi et al., 2024b).

Based on these results, we can apply the direction that gives the best steering capabilities, for either ablation or induction, on the *test* set. For evaluation, we use a set of keywords that LLMs commonly produce in refusal responses (e.g., “I’m sorry,” “I am sorry,” “I apologize”). If any of these keywords appear in a model’s response, we register a single refusal and assign a score of +1 for that proposal. We then average this score over the set to obtain the final refusal metric.

As depicted in Figure 9, the steering capabilities provided by SDialog show similar performance to those presented by (Arditi et al., 2024b) for the LLAMA-3 8B INSTRUCT model. For harmful instructions, the framework allows the LLM to bypass the refusal for 99% of the proposals. Conversely, when inducing the direction on harmless instructions, the steered version reaches 100%, indicating strong feature induction capabilities across all proposals.

**Discussion.** This case study highlights how `sdialog.interpretability` abstracts much of the boilerplate typically required for activation steering experiments. By exposing unified interfaces for layer targeting, token-level inspection, representation extraction, and in-place intervention during generation, SDialog enables rapid prototyping of mechanistic interpretability workflows without modifying model internals. Although we focused on refusal steering to reproduce the find-

ings of Arditi et al. (2024b), the same abstractions extend naturally to other techniques, including sentiment steering, persona control, bias analysis, feature visualization, linear probing, and contrastive representation studies. As such, SDialog provides a general-purpose framework for controlled intervention and behavioral analysis in large language models, facilitating reproducible and extensible interpretability research.

# Agentic AI for Human Resources: LLM-Driven Candidate Assessment

Kamer Ali Yuksel, Abdul Basit Anees, Ashraf Elneima  
Sanjika Hewavitharana, Mohamed Al-Badrashiny, Hassan Sawaf

aiXplain, Inc., San Jose, CA, USA

{kamer,abdul.anees,ashraf.hatim,sanjika,mohamed,hassan}@aixplain.com

## Abstract

In this work, we present a modular and interpretable framework that uses Large Language Models (LLMs) to automate candidate assessment in recruitment. The system integrates diverse sources—including job descriptions, CVs, interview transcripts, and HR feedback—to generate structured evaluation reports that mirror expert judgment. Unlike traditional ATS tools that rely on keyword matching or shallow scoring, our approach employs role-specific, LLM-generated rubrics and a multi-agent architecture to perform fine-grained, criteria-driven evaluations. The framework outputs detailed assessment reports, candidate comparisons, and ranked recommendations that are transparent, auditable, and suitable for real-world hiring workflows.

Beyond rubric-based analysis, we introduce an LLM-Driven Active Listwise Tournament mechanism for candidate ranking. Instead of noisy pairwise comparisons or inconsistent independent scoring, the LLM ranks small candidate subsets (“mini-tournaments”), and these listwise permutations are aggregated using a Plackett–Luce model. An active-learning loop selects the most informative subsets, producing globally coherent and sample-efficient rankings. This adaptation of listwise LLM preference modeling—previously explored in financial asset ranking (Yuksel and Sawaf, 2025)—provides a principled and highly interpretable methodology for large-scale candidate ranking in talent acquisition.

## 1 Introduction

The process of evaluating job candidates remains one of the most critical and resource-intensive tasks in human resource management. Despite advances in digital recruitment platforms, organizations still

struggle to identify and select the most suitable candidates, particularly when faced with high volumes of applicants or highly specialized roles (Raghavan et al., 2020; Vaishampayan et al., 2025). Current systems tend to over-rely on heuristics, rigid filters, or shallow keyword-based screening, leading to missed opportunities and inconsistent evaluations. Human reviewers, while invaluable, often face cognitive overload, inter-rater inconsistencies, and implicit biases that compromise the objectivity and reproducibility of hiring decisions (Quillian et al., 2017).

To address these limitations, we introduce a robust, LLM-based candidate assessment system designed to enhance the precision, fairness, and scalability of talent evaluations. Our system leverages the generative reasoning capabilities of state-of-the-art language models to generate, apply, and explain fine-grained evaluation criteria specific to each job role. This approach ensures a high degree of alignment between role requirements and candidate evaluations, enabling both depth and breadth in assessments. The system also supports multi-source integration, allowing it to analyze structured and unstructured data from resumes, interviews, recommendation letters, and internal HR notes. Through its modular and extensible architecture, our framework lays the groundwork for a new paradigm in AI-assisted hiring: one that is grounded in interpretability, adaptability, and decision support rather than opaque automation.

Beyond generating structured assessments, our system introduces a new paradigm for ranking candidates: treating the LLM as a listwise fuzzy judge. Rather than comparing two candidates at a time, the model evaluates small groups of candidates simultaneously and produces a relative ordering. These listwise rankings contain richer information per query and mirror how human hiring committees make decisions when reviewing shortlists. We ag-

---

HR Manager is online: <https://hrmanager.aixplain.com>

Video demonstration: <https://youtu.be/qwcSmWNOHRk>

gregate these orders using a Plackett–Luce model (Plackett, 1975) to estimate global latent “candidate utilities”, and employ an active-learning (Liang and Grauman, 2014) mechanism to adaptively select the most informative candidate groups for evaluation. This approach transforms candidate ranking from an ad-hoc or score-based process into a statistically principled, scalable inference.

## 2 Related Work

Prior efforts in automating candidate assessment have focused primarily on three domains: resume parsing and semantic matching, AI-based video or audio interviews, and the use of predictive analytics to estimate candidate performance. Resume parsing tools such as TextKernel and Sovren have made strides in extracting structured information from unstructured resumes, enabling preliminary filtering. However, these tools often lack contextual understanding and fail to capture the nuances of role-specific competencies (Kanikar et al., 2025). Similarly, semantic matchers like Eightfold and LinkedIn Talent Insights attempt to align job descriptions with candidate profiles using embeddings and cosine similarity measures, but tend to prioritize surface-level alignment over deeper evaluative reasoning (Bevara et al., 2025).

Recent developments in LLMs have enabled few-shot and zero-shot task performance across a wide range of domains. These capabilities have opened up new avenues for generating assessments, feedback, and structured evaluations. However, there remains a significant gap in applying LLMs to the hiring domain in a way that is structured and explainable (Ghosh and Sadaphal, 2023). Our work bridges this gap by building a layered framework that combines LLM-based reasoning with a robust set of evaluation dimensions, customized prompts, and a focus on output format alignment for real-world usability.

Recent work in LLM-based preference learning, including pairwise comparison frameworks such as PAIRS and PoE Bradley–Terry models, demonstrates that LLMs can act as noisy but meaningful preference oracles when aggregated statistically (Qin et al., 2024). However, nearly all such approaches rely on pairwise comparisons, which are noisy, vulnerable to prompt instability, and do not scale well to large candidate pools (Zeng et al., 2024; Wang et al., 2025). Our work is among the first to apply listwise LLM judgments—far more

informative than pairwise labels—to candidate assessment. Through listwise tournaments and PL aggregation, we introduce a stable, globally coherent ranking mechanism tailored to HR workflows, bridging the gaps between LLM reasoning, structured evaluation, and hiring decision support.

## 3 Architecture

The proposed system is designed as a multi-stage, modular pipeline where each stage is responsible for a distinct phase in the candidate evaluation lifecycle. At the core of the system is a suite of LLM-powered agents (Guo et al., 2024), each guided by prompt templates and domain-specific heuristics. These agents interact in a sequential yet modular fashion, allowing for plug-and-play flexibility in enterprise settings.

The process begins with the Criteria Generation Agent, which takes as input a job title and description and produces a finely detailed assessment rubric. This rubric is not generic; it is customized to reflect the specific competencies, responsibilities, and context described in the job posting. The rubric includes both technical and non-technical dimensions, and is structured to differentiate candidates not just by qualifications but by demonstrated competencies and growth potential. The system also includes a dedicated Video Question Generation Agent, which tailors reflective and dimension-specific interview questions for video-based responses. Candidate answers are then transcribed and fused with other input modalities to provide a full-spectrum evaluation.

Next, the Assessment Generator Agent uses this rubric to evaluate candidate profiles. This agent synthesizes data from candidate CVs, HR recommendations, structured interview transcripts, and chat-based conversations. It maps candidate achievements and responses to the rubric and generates a markdown-based report with explicit ratings (Low, Medium, High) for each dimension, supported by detailed justifications. These justifications not only cite specific evidence from inputs but also interpret the relevance of that evidence in the context of the role. To enhance the richness and reliability of assessment, the system integrates multimodal analysis through a video interview analysis module. This component leverages computer vision and audio processing techniques to assess facial expressions, vocal tone, and body language from recorded interviews. These non-verbal cues

are fused with the LLM-generated textual insights, offering a holistic picture of a candidate’s interpersonal dynamics, emotional intelligence, and presentation skills—traits often critical in leadership and client-facing roles.

The system also incorporates a Feedback Integration Module that closes the loop between pre-hire assessments and post-hire outcomes. Data such as performance reviews, retention metrics, and team feedback are used to retroactively validate model predictions and refine the scoring logic and rubric generation prompts. This continuous learning mechanism ensures that the system evolves to better predict candidate success, based on empirical evidence. Other agents in the system include a Formatter Agent, which ensures compliance with HR documentation standards; a Comparison Agent, which conducts side-by-side evaluations across top candidates; and a Ranking Agent, which integrates new candidates into existing ranked lists while preserving order consistency.

### 3.1 Prompt Engineering

Prompt engineering is central to the success of this system (Kojima et al., 2023). Each agent is powered by a dedicated prompt template designed to elicit high-quality, context-sensitive output from the language model. The prompts are carefully structured to ensure that outputs are not only informative but aligned with HR practices. The Criteria Generation Prompt guides the model to rewrite and enrich default assessment rubrics based on job-specific information. It instructs the model to retain general competencies such as leadership or communication, while augmenting or refining them with role-specific details. The result is a YAML schema that defines evaluation criteria in granular terms, often including 12 to 20 dimensions, each with clear definitions and expectations for various rating scales.

The Assessment Prompt is the heart of the evaluation pipeline. It consolidates inputs from multiple sources—including candidate summaries, HR notes, interviews, and prior assessments—and guides the model to produce a structured report. The prompt enforces a markdown format and mandates explanations for each rating. The prompt also incorporates an internal logic to cross-check dimensions, ensuring coherence and consistency in the assessment. Comparison and Ranking Prompts are designed to simulate the work of hiring com-

mittees. They use structured reasoning to evaluate candidates across dimensions and recommend ranked lists, justifying the positioning of each candidate. All prompts are calibrated to avoid hallucination, reduce bias, and maintain a professional tone aligned with corporate communication.

For tournament-based ranking, we design a dedicated Listwise Ranking Prompt that instructs the LLM to evaluate a small group of candidates together, using the job-specific rubric as the evaluation lens. The prompt requires the model to internally reason step-by-step but output only a permutation of candidates from strongest to weakest. This ensures consistency and avoids verbose explanations during the tournament phase. Additional prompts support active learning by incorporating weak prior orderings from the PL model, enabling the LLM to refine its comparisons based on previous rounds while avoiding bias amplification.

### 3.2 Evaluation Dimensions

The assessment criteria generated by the system cover a wide spectrum of attributes required for professional success. These dimensions are not static but tailored dynamically for each job role using the Criteria Generation Agent. Typical categories include domain expertise, technical skill sets, interpersonal effectiveness, leadership capabilities, and career motivation. For example, the Industry-Specific Fit dimension evaluates the extent to which a candidate’s experience aligns with the sector targeted by the role. Functional Fit assesses the candidate’s technical proficiency or domain fluency concerning the job’s day-to-day responsibilities.

Soft skills are also included, with dimensions like Motivation, Drive, and Continuous Learning assessing intrinsic traits and growth mindset. The system ensures that each dimension includes rubric definitions for each level of proficiency. This enables consistent application across candidates and increases inter-rater reliability when used in hybrid human-AI workflows. Listwise tournaments rely on these dimensions as implicit comparison criteria. During tournament ranking, the LLM evaluates subsets of candidates holistically across all rubric dimensions, weighing tradeoffs (e.g., strong technical skill but weaker communication) at a group level. This mirrors human holistic reasoning and produces richer preference signals than independent per-candidate scoring.

### 3.3 Report Generation

The output of the system is designed to be immediately useful to recruiters, hiring managers, and talent committees. Each assessment report is generated in structured markdown format and adheres to a consistent schema that balances detail with readability. Reports begin with a concise introduction outlining the role and the candidate’s context. The core of the report is the dimension-wise assessment, where each dimension is evaluated with a rating and a supporting explanation. These explanations are not generic. They explicitly refer to the candidate’s background, citing projects, metrics, or behavioral indicators found in resumes or interviews. This creates a transparent trail that hiring teams can audit.

Following the detailed assessment, the system synthesizes an overall readiness evaluation, a cultural fit analysis, and flags areas where the candidate may require support or development. It also suggests alternative roles if the candidate is deemed not optimal for the original role, a feature especially valuable in internal mobility or volume hiring settings. The report concludes with a summary recommendation and a confidence indicator, which can be weighted based on the richness of the input data.

While the listwise tournament mechanism is used to compute global rankings, its outcomes also influence report generation. Candidates’ latent utility values, stability across tournament rounds, and key dimensions driving their placement can be surfaced in the final report, allowing hiring teams to understand why a candidate ranks above peers and to audit the consistency of the evaluation pipeline.

### 3.4 Methodology

To complement rubric-driven candidate evaluation, our system introduces an *active listwise ranking* mechanism that treats the LLM as a fuzzy multi-candidate judge. Instead of independently scoring each candidate or performing unstable pairwise comparisons, the system repeatedly evaluates small groups of candidates using LLM-driven listwise tournaments. This produces richer comparative information and enables a principled, data-efficient approach to building a globally coherent ranking over large applicant pools. The methodology consists of three components: (1) listwise LLM comparisons, (2) probabilistic Plackett–Luce aggregation, and (3) active subset selection via uncertainty.

At each iteration, the system samples a sub-

set of  $K$  candidates (typically between 5 and 10) and presents them to the LLM with a structured prompt. The model is instructed to evaluate all candidates *simultaneously*, using the role-specific assessment rubric as the decision lens, perform step-by-step reasoning internally, and output only an ordered list from strongest to weakest fit for the job. The prompt explicitly prohibits justification during these tournament rounds to reduce verbosity and stabilize comparisons. Each listwise ranking encodes  $\frac{K(K-1)}{2}$  implied pairwise relations and captures cross-dimensional tradeoffs that only emerge when candidates are evaluated side-by-side, imitating human hiring committees.

Each LLM-generated permutation is treated as a noisy observation of latent “candidate suitability utilities.” To infer these utilities, we fit a Plackett–Luce (PL) model over all collected tournament results. For each ranking  $\pi_t$  from subset  $S_t$ , the PL likelihood factorizes over positions in the permutation, allowing efficient optimization. The resulting utility vector  $\mathbf{u}$  provides a globally consistent ranking across all candidates, stabilizes noise in raw LLM outputs, and reconciles conflicting preferences across tournament rounds. Posterior variances over  $u_i$  are approximated via a diagonal Laplace approximation for uncertainty estimates essential in active learning and targeted querying.

Rather than sampling candidate subsets uniformly, we employ an active learning loop that prioritizes the most informative groups. After each PL update, posterior variances highlight candidates with uncertain latent utilities. New subsets are formed to probe contested regions—such as boundaries between “shortlist” and “reject” tiers. Acquisition strategies, including Monte-Carlo Knowledge Gradient (MC–KG) (Baladat et al., 2020), posterior disagreement sampling (Seung et al., 1992), and KL-UCB heuristics (Garivier and Cappé, 2011), determine which subsets maximize expected information gain. This approach sharply reduces the number of LLM queries required and accelerates convergence toward a stable global ranking.

Unlike standalone ranking algorithms, the listwise tournament system utilizes the same role-specific rubric that guides the Criteria Generation and Assessment Generator agents. This ensures that tournament comparisons align with the competency framework used in detailed candidate reports, that latent utilities reflect holistic role fit rather than

superficial similarities, and that ranking outcomes remain interpretable and auditable. The final ranking is thus grounded in the same structured criteria used for individual assessments, enabling consistent and transparent decision-making across both comparative and descriptive evaluation modes.

## 4 Case Studies

We applied the system to a set of roles spanning multiple industries and seniority levels to evaluate its generalizability and precision. These included technical roles like AI Research Scientist and Staff Machine Learning Engineer, as well as leadership roles such as VP of Product and CTO for startup environments. Human experts independently rated candidates using the three-level rubric (Low/Medium/High) generated by the Criteria Generation Agent, and the system produced its own ratings using the same rubric. Agreement was computed per dimension, counting cases where the system’s rating was within one level of the human score. Across all evaluations, 87% of system ratings fell within one band of human scores. In cases where candidates appeared similar, the system was able to surface distinctions that human reviewers later confirmed. Although not directly captured by ranking metrics, qualitative inspection of intermediate outputs reveals several desirable behaviors:

- Rubric refinement yields clearer, more structured, and more discriminative evaluation.
- Subtle distinctions in candidates (e.g., leadership maturity, communication style, depth of reasoning) emerge earlier in iterations.
- The system frequently suggests alternative roles for candidates whose strengths are misaligned with the target position.

These findings show that the system not only learns a stable ranking but also produces nuanced, actionable insights for talent evaluation workflows.

### 4.1 Experiments

We evaluate the proposed LLM-driven active listwise tournament framework on a real candidate-ranking task using the full pipeline with iterative criteria refinement. In this configuration, the system not only performs listwise tournaments and Plackett–Luce aggregation but also refines the scoring rubric at each iteration based on LLM critique. The goal of this experiment is to measure ranking fidelity relative to human expert judgments, assess

convergence behavior, and characterize how active learning improves ranking stability.

We run 30 iterations of active listwise querying over a fixed pool of candidates. Each consists of:

1. Refining the assessment rubric via LLM-based critique.
2. Selecting next candidate subset using Monte-Carlo Knowledge Gradient (MC–KG).
3. Obtaining a listwise ranking from the LLM for the selected subset.
4. Updating global candidate utilities through Plackett–Luce optimization.

We evaluate performance using NDCG@K for  $K = \{10\%, 15\%, 20\%, 25\%\}$  using human ranking as reference, and convergence metrics that measure the **stability of the ranking updates**, where higher values mean the ranking is stabilizing and the system is no longer making large structural adjustments. Figure 1 shows the evolution of NDCG across cutoffs during the active-learning process.

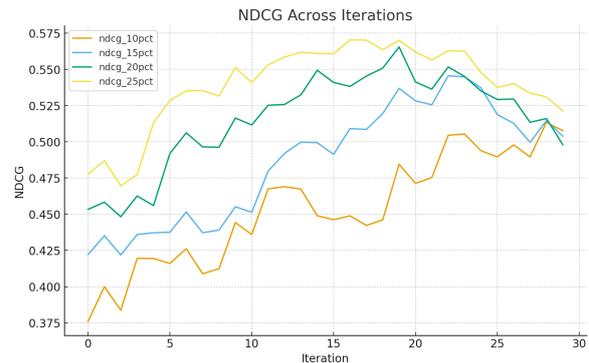


Figure 1: NDCG@K% progression across 30 iterations.

Across all cutoffs, NDCG improves steadily during the first ~20 iterations, with NDCG@25% achieving the highest value. This shows that the active listwise tournament mechanism extracts a ranking structure aligning with human preferences.

Table 1: Peak NDCG@K values across cut-off ratios.

Cutoff	10%	15%	20%	25%
NDCG@K	0.5134	0.5455	0.5655	0.5703

To assess convergence, we track two indicators:

1. **Kendall- $\tau$  between successive iterations**, capturing how similar the ranking at iterations  $t$  is to iterations  $t - 1$ .

## 2. Utility movement $\Delta\mathbf{u}$ , defined as the norm of the change in the global PL utility vector.

Since these metrics operate on different scales, each is independently normalized to  $[0, 1]$  for visualization. Figure 2 shows their joint progression. The Kendall- $\tau$  curve rises sharply during early iterations, indicating rapid stabilization of the ranking. Meanwhile,  $\Delta\mathbf{u}$  decreases substantially after iteration 10, showing that the underlying utility estimates converge. Together, these trends demonstrate that active querying helps the system quickly settle into a stable, high-confidence ranking.

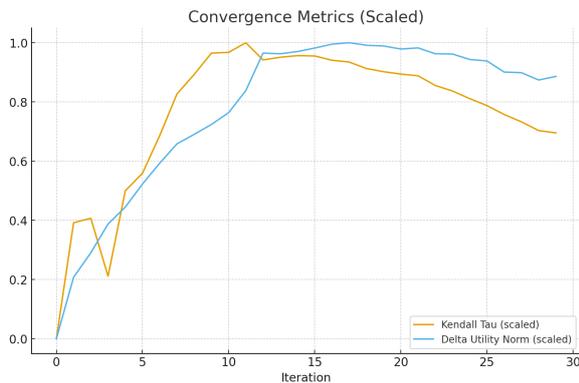


Figure 2: Convergence metrics over 30 iterations. Both Kendall- $\tau$  (stability between iterations) and  $\Delta\mathbf{u}$  (norm of utility change) are independently scaled to  $[0, 1]$ .

## 5 Discussion

The LLM-based candidate assessment system presents a significant advancement in the intersection of AI and human capital management. By combining rigorous prompt engineering, dynamic rubric generation, and structured reasoning, the system offers a scalable alternative to traditional candidate evaluations. It provides organizations with an interpretable, customizable, and reproducible method for evaluating talent, while also supporting diversity, equity, and inclusion through the criteria.

Nonetheless, there are challenges. The system’s accuracy is influenced by the quality of the input data. Poorly written job descriptions or ambiguous candidate materials can reduce the effectiveness of the evaluation. Additionally, while the system addresses many soft skill dimensions through textual inputs, it can currently recognize facial expressions but cannot yet interpret broader non-verbal cues—such as tone of voice or full-body language. Finally, ongoing prompt tuning is needed to ensure domain-specific relevance in niche or highly

specialized roles.

The addition of listwise LLM tournaments introduces a powerful inference layer that helps address a persistent challenge in hiring: ranking candidates with overlapping profiles. By modeling candidate evaluation as a probabilistic ranking problem—rather than independent scoring—we gain stability, interpretability, and resistance to prompt noise. However, this approach also requires careful control of rubric drift and must avoid reinforcing biases across rounds of active querying. These challenges motivate future work combining fairness-aware constraints with PL-model aggregation.

## 6 Conclusion

This paper introduces a novel, LLM-driven candidate assessment system that transforms the way organizations evaluate talent. By uniting the power of language models with structured rubric generation, contextual analysis, and transparent reporting, we offer a framework that is both scalable and interpretable. The system addresses critical limitations in existing hiring workflows by enabling fine-grained, reproducible assessments that align closely with job-specific expectations. With a focus on real-world usability, modular design, and human-centered outputs, our framework lays a strong foundation for the future of AI-assisted hiring.

By extending the Active Listwise Tournament framework from financial asset ranking (Yuksel and Sawaf, 2025) to candidate evaluation, we show that LLMs can serve not only as rubric-based assessors but also as structured preference oracles. The combination of listwise evaluations, probabilistic PL aggregation, and active learning yields globally coherent candidate rankings that scale to large applicant pools and outperform traditional scoring pipelines. This fusion of LLM reasoning and ranking theory provides a principled foundation for next-generation hiring decision support systems.

Another promising direction is the development of a real-time interview assistant that can suggest probing questions to human interviewers based on live candidate responses. This would create a hybrid workflow where AI augments rather than replaces human judgment, preserving the richness of interpersonal evaluations while enhancing structure and consistency. We also plan to develop a feedback loop that integrates hiring outcomes and performance reviews into the system, enabling it to refine its prompts and scoring models over time.

## Limitations

While the proposed framework offers significant advancements in candidate assessment, it is not without limitations. First, although the integration of LLMs and multimodal analysis allows for nuanced evaluation, the system’s performance is still bounded by the quality and completeness of candidate input data. Sparse or low-quality resumes, brief interview responses, or poorly recorded video inputs can reduce the accuracy of assessments.

Second, although the system supports multilingual evaluation and localization, the depth of cultural adaptation is currently limited to language and structural format. Subtle sociocultural dynamics or communication norms may still be misinterpreted by the LLM or the multimodal modules.

Lastly, while interpretability is a design principle, some of the underlying reasoning from the LLM remains opaque, especially in complex judgment tasks that involve implicit contextual inference. Efforts have been made to structure outputs and provide rationale, but full transparency into model reasoning is still an open challenge.

## Ethics Statement

This system aims to democratize access to fair candidate assessments, reducing reliance on subjective heuristics and increasing transparency. However, it also raises several ethical considerations.

First, while the system seeks to mitigate bias, LLMs are trained on large-scale internet data that may contain embedded societal biases. If unchecked, these biases can influence rubric generation, assessments, and candidate rankings.

Second, there is the risk of over-reliance on automated systems in critical decision-making. Our system is designed to augment human judgment, not replace it. Recommendations should always be reviewed by qualified HR personnel who contextualize the outputs with domain-specific insight.

Finally, deploying this system at scale across regions and cultures requires sensitivity to local labor laws, cultural practices, and fairness norms.

While these challenges are non-trivial, we believe that the benefits of a well-designed, transparent, and modular AI-assisted evaluation system can meaningfully improve hiring outcomes—if developed and deployed responsibly.

## References

- Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. 2020. [Botorch: A framework for efficient monte-carlo bayesian optimization](#).
- R. V. K. Bevara, N. R. Mannuru, S. P. Karedla, B. Lund, T. Xiao, H. Pasem, S. C. Dronavalli, and S. Rudreshkumar. 2025. [Resume2vec: Transforming applicant tracking systems with intelligent resume embeddings for precise candidate matching](#). *Electronics*, 14(4):794.
- Aurélien Garivier and Olivier Cappé. 2011. The kl-ucb algorithm for bounded stochastic bandits and beyond. In *Proceedings of the 24th Annual Conference on Learning Theory (COLT 2011)*, volume 19 of *Proceedings of Machine Learning Research*, pages 359–376.
- Preetam Ghosh and Vaishali Sadaphal. 2023. [Jobrecogpt – explainable job recommendations using llms](#).
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. [Large language model based multi-agents: A survey of progress and challenges](#).
- Darsh Kanikar, Pratyush Jain, Siddhi Jain, and Lalit Purohit. 2025. Systematic review of methods for analysis of resumes. In *Proceedings of the International Conference on Computer Science and Communication Engineering*, pages 301–308. Atlantis Press. Systematic survey of resume analyzers, including commercial tools such as Textkernel.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. [Large language models are zero-shot reasoners](#).
- Lucy Liang and Kristen Grauman. 2014. Beyond comparing image pairs: Setwise active learning for relative attributes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2856–2863.
- L. Plackett, R. 1975. [The analysis of permutations](#). *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 24(2):193–202.
- Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. [Large language models are effective text rankers with pairwise ranking prompting](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, page 1504–1518.
- Lincoln Quillian, Devah Pager, Ole Hexel, and Arnfinn H. Midtbøen. 2017. [Meta-analysis of field experiments shows no change in racial discrimination in hiring over time](#). *Proceedings of the National Academy of Sciences*, 114(41):10870–10875.

- Manish Raghavan, Solon Barocas, Jon Kleinberg, and Karen Levy. 2020. [Mitigating bias in algorithmic hiring: Evaluating claims and practices](#). In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT\* / FAccT)*, pages 469–481. ACM.
- H. Sebastian Seung, Manfred Opper, and Haim Sompolinsky. 1992. [Query by committee](#). *Proceedings of the Fifth Annual Workshop on Computational Learning Theory (COLT)*, pages 287–294.
- Swanand Vaishampayan, Hunter Leary, Yoseph Berhanu Alebachew, Louis Hickman, Brent Stevenor, Weston Beck, and Chris Brown. 2025. [Human and llm-based resume matching: An observational study](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 4808–4823.
- Yidong Wang, Yunze Song, Tingyuan Zhu, Xuanwang Zhang, Zhuohao Yu, Hao Chen, Chiyu Song, Qiufeng Wang, Cunxiang Wang, Zhen Wu, Xinyu Dai, Yue Zhang, Wei Ye, and Shikun Zhang. 2025. [Trustjudge: Inconsistencies of llm-as-a-judge and how to alleviate them](#). *arXiv preprint*. ArXiv:2509.21117.
- Kamer Ali Yuksel and Hassan Sawaf. 2025. [Llm-driven active listwise tournaments for portfolio selection in large asset universes](#). *SSRN Electronic Journal*. Preprint; Available at SSRN: <https://ssrn.com/abstract=5764903>.
- Yifan Zeng, Ojas Tendolkar, Raymond Baartmans, Qingyun Wu, Lizhong Chen, and Huazheng Wang. 2024. [Llm-rankfusion: Mitigating intrinsic inconsistency in llm-based ranking](#). *arXiv preprint (cs.IR)*. ArXiv:2406.00231.

# Trove: A Flexible Toolkit for Dense Retrieval

Reza Esfandiarpour   Max Zuo   Stephen H. Bach

Department of Computer Science, Brown University  
{reza\_esfandiarpour, zuo, stephen\_bach}@brown.edu

Demo: <https://youtu.be/rThGH0w3wS8>



[ir-trove.dev](https://github.com/ir-trove)

## Abstract

We introduce Trove, an easy-to-use open-source retrieval toolkit that simplifies research experiments without sacrificing flexibility or speed. For the first time, we introduce *efficient* data management features that load and process (filter, select, transform, and combine) retrieval datasets on the fly, with just a few lines of code. This gives users the flexibility to easily experiment with different dataset configurations without the need to compute and store multiple copies of large datasets. Trove is highly customizable: in addition to many built-in options, it allows users to freely modify existing components or replace them entirely with user-defined objects. It also provides a low-code and unified pipeline for evaluation and hard negative mining, which supports multi-node execution without any code changes. Trove’s data management features reduce memory consumption by a factor of 2.6. Moreover, Trove’s easy-to-use inference pipeline incurs no overhead, and inference times decrease linearly with the number of available nodes. Most importantly, we demonstrate how Trove simplifies retrieval experiments and allows for arbitrary customizations, thus facilitating exploratory research.

## 1 Introduction

Influential toolkits such as those provided by Hugging Face (HF) simplify Machine Learning (ML) pipelines and support extensive customization with minimal effort, thus facilitating exploratory research (Gugger et al., 2022; Lhoest et al., 2021; Wolf et al., 2020). Similarly, existing retrieval toolkits have significantly improved Information Retrieval (IR) pipelines (Gao et al., 2022; Reimers and Gurevych, 2019). However, IR experiments still require a considerable amount of engineering effort for many tasks like efficient data management or model customization. Here, we introduce a novel open-source toolkit that simplifies various stages of retrieval pipelines, enabling ef-

ficient data management, flexible modeling, and easy distributed evaluation. Our design prioritizes customization and makes it easy to freely modify or entirely replace each component.

General-purpose toolkits are not directly applicable in retrieval pipelines. Retrieval is uniquely different from most ML problems in that instances of retrieval tasks are not self-contained. For example, while solving an image classification task only involves one image, a single retrieval task involves one query and the *entire corpus*. This makes retrieval experiments more challenging. Since most data management tools like HF Datasets (Lhoest et al., 2021) process each instance isolated from the rest of the dataset, they cannot be directly used in retrieval pipelines (Gao et al., 2022). Distributed evaluation is also more challenging. Instances of retrieval tasks share a lot of the computation (i.e., encoding the corpus), and we cannot simply evaluate each instance on a separate device and aggregate the results. Finally, HF transformers models only provide the encoder, and we cannot directly use them for retrieval without additional modeling.

Existing toolkits have recognized these issues and offer initial solutions. However, these solutions are often not as flexible or easy to use. Since naive on-the-fly data preparation for retrieval is memory-intensive, current toolkits rely on large pre-processed dataset files (Gao et al., 2022), often duplicating a lot of data for variations of a single dataset (Fig. 1 top). Although evaluating retrieval tasks is computationally more demanding, currently distributed evaluation is either limited to a single node (Muennighoff et al., 2022; Reimers and Gurevych, 2019) or involves several steps and more engineering effort (Gao et al., 2022). For modeling, current frameworks wrap transformers models in fixed classes, and customizations are limited to a set of pre-defined options. As a result, exploratory experiments require significant engineering effort, which slows down novel research.

In this work, we introduce Trove, an open-source library that simplifies dense retrieval experiments without sacrificing flexibility. Trove is the first toolkit to provide features for efficiently managing and pre-processing retrieval data on the fly (Fig. 1 bottom). Trove provides a simple interface for multi-node/GPU inference and is fully compatible with HF transformers ecosystem. Our modeling approach provides direct access to model components and allows arbitrary customizations. In general, our design increases flexibility at three levels. 1) Trove provides various built-in options to customize experiments. 2) Our transparent design allows users to override many methods with custom logic. 3) Trove’s modular structure allows users to entirely replace many components with arbitrary objects. In summary:

- Trove is built around the unique characteristics of the retrieval task and, for the first time, offers fast and memory-efficient operations for loading and pre-processing (filter, transform, combine, etc.) retrieval data *on the fly*.
- Trove provides a simple and unified interface for evaluation and hard negative mining, which supports both multi-node and multi-GPU inference without additional code.
- Trove allows for direct customization of all modeling components or even replacing them with arbitrary modules, while maintaining compatibility with HF transformers ecosystem.
- Trove is designed with customization in mind. The codebase is heavily documented and easy to understand. We provide ample guides and examples to facilitate customization.

## 2 Background and Challenges

There is a growing body of work on transformer-based dense retrievers. Many works improve the training data through techniques like using mined hard negatives or a large number of random in-batch negatives (Karpukhin et al., 2020; Moreira et al., 2024; Qu et al., 2021; Rekabsaz et al., 2021; Xiong et al., 2020; Zerveas et al., 2022, 2023; Zhan et al., 2021). Several works also use synthetic data for training (Alaofi et al., 2023; Bonifacio et al., 2022; Dai et al., 2022; Jeronimo et al., 2023; Lee et al., 2024; Li et al., 2024). With the introduction of RepLLaMA (Ma et al., 2024), large

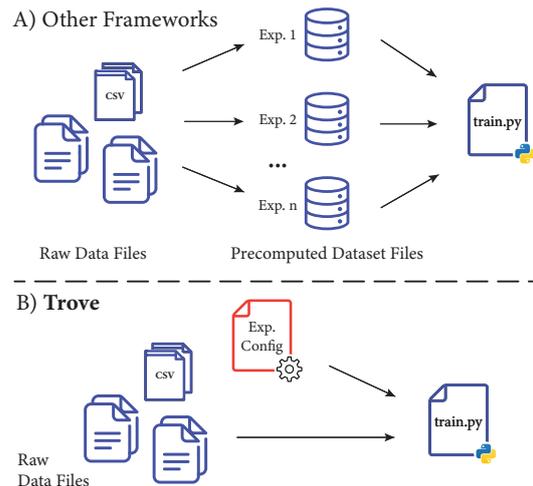


Figure 1: A) Existing toolkits require manually creating and maintaining large pre-processed data files for each experiment. B) Trove processes datasets on the fly based on the given configuration options.

decoder-only LLMs and PEFT techniques (Hu et al., 2022) have become popular for retrieval models (Wang et al., 2024). There are also new variants of the retrieval task itself, such as retrieval instructions (Asai et al., 2022; Weller et al., 2024).

Despite the dynamic research landscape, developing IR pipelines is still challenging. To illustrate the challenges of developing IR pipelines, we compare the development experience for IR to other ML tasks for three common operations.

**Data Management** For most ML problems, tools like HF Datasets load, pre-process, and combine multiple datasets on the fly with just a few lines of code and minimal memory overhead. In IR, even generating training samples for standard datasets like MS MARCO (Bajaj et al., 2018) requires a lot of memory and extra code. Existing tools do not offer any data management functionality and instead rely on large pre-processed data files, which require maintaining many large files with duplicate data. In addition to the cumbersome process, data changes are not trackable by VCS with this approach, hurting reproducibility.

**Modeling** Usually, users have full control over the model and can apply *arbitrary* customizations (e.g., add LoRA adapters or change the loss). However, current retrieval toolkits wrap the encoder in custom classes, without providing direct access to the transformers backbone. As a result, any customization requires explicit support from the library<sup>1</sup>. The interactions between different com-

<sup>1</sup>Example: [sentence-transformers/issues/2575](https://github.com/huggingface/sentence-transformers/issues/2575)

ponents also limit flexibility. For example, it is not possible to train with graduated relevance labels using Tevatron (Gao et al., 2022) even if we overwrite the loss function.

**Inference** The evaluation pipeline often involves a simple script that executes the forward pass and calculates the metrics, all in one job step. For distributed evaluation, users just need to execute the same script with a distributed launcher like Accelerate (Gugger et al., 2022), with minimal changes. Although IR evaluation is computationally more demanding, multi-node execution is not straightforward. The evaluation process with SentenceTransformers and MTEB<sup>2</sup> packages is easy but limited to only a single node. Tevatron supports multi-node evaluation, but it needs to manually launch multiple jobs to encode each dataset shard separately, and then launch another job to retrieve related documents and another job to calculate the metrics. See Section B for discussion of other libraries.

### 3 System Design

Here, we first explain Trove experiment workflows (Fig. 2) and then describe its major components.

#### 3.1 Workflow

Trove experiment workflows are simple and based on configuration objects. Users create one or more instances of `MaterializedQRelConfig` to specify how raw input files (query, corpus, qrels) should be loaded and processed (e.g., filtered). We also create a `DataArguments` object with dataset-level details (e.g., sequence length). Users then use these objects to instantiate one of the main dataset classes (`BinaryDataset` or `MultiLevelDataset`). Next, we create a `ModelArguments` instance with model details like name, pooling type, and LoRA configuration. We instantiate the main retriever (e.g., `BiEncoderRetriever`) from the argument object. Finally, we use these components in addition to an instance of `RetrievalTrainingArguments` (`EvaluationArguments`) to instantiate `RetrievalTrainer` (`RetrievalEvaluator`), which is responsible for the main training (evaluation) loop. Our design allows instantiating configuration objects from command-line arguments, which makes it easier to run diverse experiments. Moreover, our design increases flexibility by exposing the main components of the pipeline: users can easily customize or entirely replace each

<sup>2</sup>[gh/embeddings-benchmark/mteb](https://github.com/Embeddings-Benchmark/mteb)

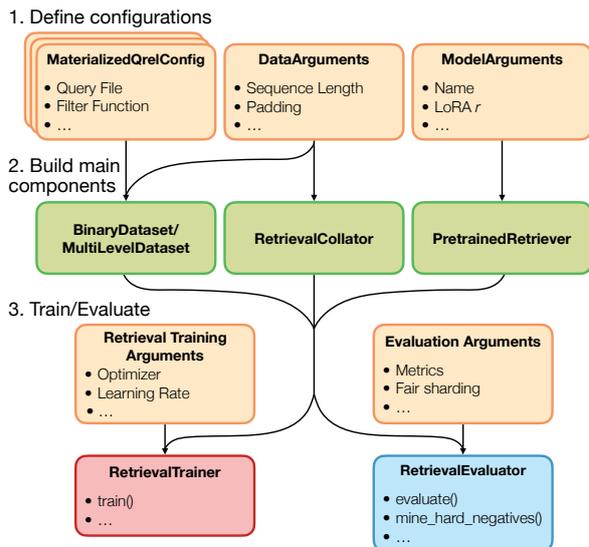


Figure 2: Training and evaluation workflow with Trove.

module without changing Trove’s source code or even major changes to their experiment workflow.

#### 3.2 Data Management

Large IR datasets are often made of three sets of files: query, corpus, and qrels (i.e., annotations). Creating training instances is simple: for each query, find the related document IDs from the qrel files and then replace these IDs with the content from the query and corpus files. Efficiently implementing this logic is challenging for large datasets like MS MARCO with 500K queries and 8M documents. Just loading all the query and corpus records consumes a lot of memory. Moreover, with multiple datasets or more complex pipelines, pre-processing or merging millions of qrel records slows down the program.

##### 3.2.1 Internal Representation

We implement `MaterializedQRel`, an efficient container for IR data that holds query, corpus, and qrel records. We use the Polars library<sup>3</sup> to efficiently group qrel triplets by query ID, which significantly speeds up pre-processing and lookup operations for related documents. We convert query, corpus, and grouped qrel records to memory-mapped Apache Arrow tables that are indexable by ID. `MaterializedQRel` only works with IDs, *without loading the actual data*. For each training instance, we load the data at the very last step and even then only load the necessary records for the current instance. As a result, `MaterializedQRel` significantly reduces memory consumption.

<sup>3</sup><https://pola.rs>

MaterializedQRel enables users to process the data just by setting a few config values. For instance, users can easily filter qrel triplets, select a subset of queries, or change the labels. See Section 4 for examples.

### 3.2.2 User-facing Classes

These benefits are available to users through the MultiLevelDataset and BinaryDataset classes. Trove datasets are made of one or more MaterializedQRel instances, which allows defining complex data pipelines. For instance, users can apply different pre-processing steps to each data source (e.g., real and synthetic) before combining them. To create a dataset, users just need to initialize the dataset class with config objects that specify data loading and processing details. As a result, Trove’s data pipelines are trackable with VCS.

Our EncodingDataset prepares the data for encoding during inference and simplifies embedding cache management. We implement the embedding cache as memory-mapped Apache Arrow tables, indexed by ID. The cache supports *lazy loading*, loading each cached vector only when needed. The interface is simple: users call `cache_records(ids, vectors)` to write the vectors and their IDs to cache. Later, when accessing each item (i.e., `dataset[i]`), the dataset returns the cached embedding instead of raw text if available. We also implement the RetrievalCollator, which tokenizes and batches examples for retrieval.

### 3.2.3 Additional Features

**Callbacks for Flexibility** We implement frequently changing operations as callback functions for easier customization. For example, users can load qrel triplets from custom file formats by registering a loader with `@register_loader`. Users can customize input formatting (e.g., add instructions) by passing `format_query` and `format_passage` callbacks to the dataset. Similarly, the `filter_fn` option in `MaterializedQRelConfig` allows users to filter the qrel triplets with custom functions.

**Reliability** We cache intermediate artifacts in the first run and track changes using a fast fingerprinting method. We also use atomic write operations to guard against corrupted files. As a result, datasets are very fast after the first run and reliably generate the same data in all runs.

By default, Trove supports datasets in the widely adopted BEIR format (Thakur et al., 2021). We also provide decorators for registering new loader

functions for custom file formats. Trove can also load file URIs directly from HuggingFace Hub.

## 3.3 Modeling

Trove’s modeling is divided into three main components (retriever, encoder, and loss) and allows users to customize each component independently.

**Retriever** Subclasses of `PretrainedRetriever` are the main model class in Trove and consist of an encoder, loss function, and the retrieval logic. `PretrainedRetriever` can use all HF transformers models as encoder and supports common pooling and normalization operations, as well as LoRA adapters and quantization. It provides the `from_model_args()` method that creates the correct encoder and loss function based on the given options. To allow arbitrary customizations, we encapsulate all details related to transformers models (e.g., quantization) in the encoder and allow users to use arbitrary `nn.Module` objects as the encoder.

Users can subclass `PretrainedRetriever` and overwrite the `forward()` method to implement custom retrieval logics. Trove already comes with `BiEncoderRetriever`, which implements the dual-encoder retrieval logic with support for cross-device in-batch negatives.

**Encoder** To experiment with new encoding methods (e.g., different pooling or PEFT techniques), users can implement custom encoder wrappers as `PretrainedEncoder` subclasses. Compared to using arbitrary `nn.Module` objects as encoder, this allows us to swap encoder wrappers without changing the code, which simplifies user scripts. Users just need to instantiate the retriever with different options (e.g., `encoder_class="MyEncoderClass"`).

**Loss Function** Trove implements the InfoNCE and KL Divergence losses. Users can implement new loss functions as `RetrievalLoss` subclasses and choose the correct loss through retriever options (e.g., `loss="MyLossClass"` or `"kl"`).

## 3.4 Training

Inspired by Tevatron (Gao et al., 2022), we ensure all Trove components are compatible with HF transformers and directly use its Trainer module for training, with minimal changes.

Trove makes it possible to approximate IR metrics like nDCG during training by ranking a small number of annotated documents for each development query, similar to a reranking task. We

---

```

1 from transformers import AutoTokenizer, HfArgumentParser
2 from trove import *
3
4 parser = HfArgumentParser((RetrievalTrainingArguments, ModelArguments, DataArguments))
5 train_args, model_args, data_args = parser.parse_args_into_dataclasses()
6
7 tokenizer = AutoTokenizer.from_pretrained(...)
8 model = BiEncoderRetriever.from_model_args(...)
9 collator = RetrievalCollator(data_args, tokenizer, append_eos=False)
10
11 pos = MaterializedQRelConfig(min_score=1, qrel_path="qrels/train.tsv", ...)
12 neg = MaterializedQRelConfig(group_random_k=2, qrel_path="mined_neg.tsv", ...)
13 dataset = BinaryDataset(data_args, model.format_query, model.format_passage, pos, neg)
14
15 trainer = RetrievalTrainer(model, train_args, collator, dataset)
16 trainer.train()

```

---

Figure 3: Training with Mined Hard Negatives

introduce IRMetrics, which can be used as the `compute_metric` callback to efficiently calculate approximate IR metrics during training for small instances of `MultiLevelDataset`.

### 3.5 Inference

`RetrievalEvaluator` class implements a simple and unified interface for evaluation and hard negative mining. Inference is as easy as creating an instance of `RetrievalEvaluator` and calling the `evaluate()` or `mine_hard_negatives()` method. Trove supports logging to Wandb and can integrate other experiment trackers using callback functions.

For distributed inference, we simply need to launch the same script without any changes using a distributed launcher. `RetrievalEvaluator` automatically distributes the computation across available *nodes* and GPUs. We also introduce a *fair sharding* feature that allows mixing GPUs with different capabilities without stalling the faster devices. Trove adjusts the shard sizes based on GPU throughput, assigning more samples to faster devices.

Trove introduces `FastResultHeapq`, a Pytorch alternative to the standard Python `heapq` that uses fast matrix operations and GPU acceleration for tracking the top-k documents for each query<sup>4</sup>. Existing frameworks often use Python’s `heapq`, which is a major bottleneck and stalls GPU cycles during evaluation (Muennighoff et al., 2022). `FastResultHeapq` is 16x and 600x faster than Python `heapq` for online and cached embeddings, respectively.

## 4 Demonstration

Here, we demonstrate Trove’s flexibility and ease of use and benchmark its efficiency.

<sup>4</sup>This is not a full `heapq`. It just mimics some functionalities to keep track of topk documents for each query.

### 4.1 Flexibility and Ease of Use

We have already used Trove for large-scale research experiments in our earlier work, SyCL (Esfandiarpour et al., 2025). Below, we outline the pipeline for two key SyCL experiments. Trove can be easily installed from PyPI:

```
$ pip install ir-trove
```

Trove greatly reduces the engineering effort required for common training setups. Figure 3 shows the complete code needed to train dense retrievers with mined hard negatives using InfoNCE loss. This simple code already supports multi-node/GPU training, standard pooling and normalization operations, LoRA adapters, and quantization.

Now, we modify the code to train on a mix of multi-level synthetic data (labels in {0, 1, 2, 3}), annotated positives, and mined hard negatives. To do this, we simply replace lines 11–13 in Fig. 3 with the following:

```

syn = MaterializedQRelConfig(...,
    qrel_path="synth_qrel.tsv",
    corpus_path="synth_corpus.jsonl",
    query_subset_path="qrels/orig_train.tsv")
pos = MaterializedQRelConfig(...,
    qrel_path="qrels/train.tsv",
    score_transform=3,
    min_score=1)
neg = MaterializedQRelConfig(...,
    qrel_path="mined_neg.tsv",
    score_transform=1,
    group_random_k=2)
dataset = MultiLevelDataset([syn, pos, neg], ...)

```

This snippet processes each data source differently and combines the results. `syn` collection selects only synthetic data for training queries, using query IDs from `qrels/train.tsv` file. `pos` collection filters for documents with relevance labels  $\geq 1$  (i.e., positives), then assigns them a new label of 3. And,

neg randomly selects two of the hard negatives per query and assigns them a new label of 1.

In SyCL, we also explore the Wasserstein distance as loss function. For this, we just implement the loss function as the following and use `--loss=ws` to run the training script.

```
class WSLoss(RetrievalLoss):
    _alias = "ws"

    def forward(self, logits, label):
        loss = ... # calculate the loss value
        return loss
```

#### 4.1.1 Model Customization

Trove provides different methods for customizing the model. As described in Section 3.3, Trove comes with many built-in options for choosing different pooling operations, applying embedding normalization, adding LoRA adapters, or quantization. Here, we provide several examples of how users can further customize models beyond these options.

**Input Formatting** For convenience, Trove encoders include two methods for proper input formatting for a given model. The code below, modifies these methods to add instructions to input queries and passages, similar to Wang et al. (2023).

```
class EncoderWithInstructions(DefaultEncoder):
    _alias = "encoder_with_inst"

    def format_query(self, text, dataset,
        ↪ **kwargs):
        if dataset == "msmarco":
            inst = "Instruct: Given a web search
            ↪ query, retrieve relevant passages
            ↪ that answer the query\nQuery: "
        else:
            inst = "Query: "
        return inst + text

    def format_passage(self, text, title=None,
        ↪ **kwargs):
        return f"Passage: {title} {text}"
```

Then, in the user scripts (e.g., Fig. 3), they just need to pass `--encoder_class="encoder_with_inst"` to use this modified encoder wrapper.

**Pooling Method** Here is an example of how users can modify the default encoder wrapper to implement different pooling operations.

```
class EncoderWithNewPooling(DefaultEncoder):
    _alias = "encoder_new_pooling"

    def encode(self, inputs) -> torch.Tensor:
        output = self.model(**inputs,
            ↪ return_dict=True)
        embs = ... # custom pooling and
            ↪ normalization operations
        return embs
```

Similar to above, users can use `--encoder_class="encoder_new_pooling"` to use the above encoder wrapper.

**New Encoder Wrappers** Users can also directly subclass `PretrainedEncoder`, which gives them control over how the model is loaded, saved, and used for calculating the embeddings.

```
class CustomEncoder(PretrainedEncoder):
    _alias = 'custom_encoder'

    def __init__(self, args:
        ↪ trove.ModelArguments, **kwargs):
        self.model = AutoModel.from_pretrained(
            args.model_name_or_path)
        # Arbitrary Model Customizations (LoRA,
        ↪ Quantization, etc.):
        # ...

    def save_pretrained(self, *args, **kwargs):
        ...

    def encode_query(self, inputs):
        ...

    def encode_passage(self, inputs):
        ...
```

These new encoder wrappers are also automatically registered and available through configuration options (e.g., `--encoder_class="custom_encoder"`).

**User-provided Encoder Objects** Users can also directly instantiate `BiEncoderRetriever` with any `nn.Module` object as the encoder.

```
custom_encoder: torch.nn.Module = ... # any
↪ encoder module
model =
↪ BiEncoderRetriever(encoder=custom_encoder,
↪ model_args)
```

## 4.2 Efficiency

Here, we benchmark the impact of Trove's optimizations for data management and inference. Our main goal is to benchmark Trove's efficiency when handling large-scale datasets. Therefore, as a representative example of large-scale datasets, we use MS MARCO for benchmarking. Moreover, to further stress-test Trove, we also increase the size of the original MS MARCO dataset by adding synthetic documents.

**Data Management** Table 1 shows the memory required to prepare MS MARCO data for training. The naive baseline loads the entire data in memory. Trove cuts memory usage by 2.6x. Since Trove's main data is on disk and memory-mapped, increasing the data size only marginally increases

	Real		Real w/ Synth.	
	1x GPU	8x GPU	1x GPU	8x GPU
Naive	8.85	70.80	11.30	90.40
Trove	<b>3.34</b>	<b>26.72</b>	<b>4.07</b>	<b>32.56</b>

Table 1: Memory consumption in GB

the memory usage: when combining synthetic and real data (Esfandiarpour et al., 2025), Trove adds 0.73 GB of memory to load the additional 2M synthetic passages, far less than the naive approach, which uses 2.45 GB. This efficiency is critical for distributed training, where each process loads its own data. On a machine with 8 GPUs, the naive method consumes 90 GB of RAM just for data loading. We note that most existing frameworks rely on large pre-processed data files for each individual experiment, and Trove is the only library to offer efficient on-the-fly data processing for retrieval.

	Real		Real w/ Synth.	
	TTFS	TTFS <sub>1st</sub>	TTFS	TTFS <sub>1st</sub>
Naive	31	31	40	40
Trove	5	39	7	55

Table 2: Time to first sample (TTFS) in seconds for the first and subsequent runs

Table 2 reports time to first sample (TTFS), measuring the time to load and process the data. Thanks to Trove’s internal caching, after the first run, the data is available almost instantaneously. While TTFS has minimal impact on long-running experiments, a short TTFS is critical for efficient debugging and interactive development.

**Inference** Table 3 shows retrieval times for all queries in MS MARCO using E5-Mistral-Instruct (Wang et al., 2024) in a distributed environment. Inference time decreases linearly with the number of nodes, showing that Trove uses additional nodes with no overhead. Crucially, we just need to run the same script with a distributed launcher, without changing the code.

	1x Node	2x Node	3x Node
Inference Time	14h:20m	7h:12m	4h:48m

Table 3: Inference time on different numbers of nodes. Note that other libraries are limited to 1x Node.

Table 4 compares the performance of Python’s heapq with FastResultHeapq for keeping track of top-k documents at MS MARCO scale. In an on-line setup where we embed a small batch of 256 documents and compare it with queries on the fly, Trove is more than 600x faster than Python’s heapq. When the number of queries grows (e.g., for hard negative mining), Python’s heapq becomes unusable, taking up to 130 hours.

Queries	On The Fly		w/ Cached Embs	
	Standard	Trove	Standard	Trove
6K	1h:9m	<b>7s</b>	21s	<b>1s</b>
500K	130h:40m	<b>11m:45s</b>	30m:17s	<b>1m:52s</b>

Table 4: Performance of Python’s heapq vs Trove’s FastResultHeapq during inference

Even when embeddings are cached and comparisons are made in large batches (e.g., 40,960 documents) on GPU, Trove remains 16x to 21x faster. However, in practice, this speedup is not realized for Python’s heapq. It is often bottlenecked by disk I/O, particularly with the simple caching mechanisms used by existing frameworks.

**Training Results** For training results and additional experiments using Trove, we refer readers to the SyCL paper (Esfandiarpour et al., 2025) and codebase<sup>5</sup>, which uses Trove for all its experiments.

## 5 Conclusion

In this work, we introduce Trove, an open-source toolkit for dense retrieval that reduces the engineering effort in research experiments. Trove eliminates the need for large pre-processed data files and, for the first time, provides data management features that load and process retrieval data on the fly, with a small memory footprint. Trove provides full control over modeling and allows users to freely customize different modeling components. Trove provides a simple and unified interface for evaluation and hard negative mining, which supports multi-node inference without any extra code. While Trove provides a simple high-level interface, every component is designed to be configured, modified, or replaced entirely. As a result, Trove provides researchers with a tool to quickly and freely experiment with new ideas.

<sup>5</sup>[gh/BatsResearch/sycl](https://github.com/BatsResearch/sycl)

## Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. RISE-2425380. 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 National Science Foundation. This research is supported in part by the Office of Naval Research (ONR) award N00014-20-1-2115. Disclosure: Stephen Bach is an advisor to Snorkel AI, a company that provides software and services for data-centric artificial intelligence.

## References

- Marwah Alaofi, Luke Gallagher, Mark Sanderson, Falk Scholer, and Paul Thomas. 2023. Can generative llms create query variants for test collections? an exploratory study. In *Proceedings of the 46th international ACM SIGIR conference on research and development in information retrieval*, pages 1869–1873.
- Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh Hajishirzi, and Wen-tau Yih. 2022. Task-aware retrieval with instructions. *arXiv preprint arXiv:2211.09260*.
- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. **MS MARCO: A Human Generated MACHine Reading COMprehension Dataset**. *arXiv:1611.09268 [cs]*. ArXiv: 1611.09268.
- Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and Rodrigo Nogueira. 2022. Inpars: Data augmentation for information retrieval using large language models. *arXiv preprint arXiv:2202.05144*.
- Zhuyun Dai, Vincent Y Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith B Hall, and Ming-Wei Chang. 2022. Promptagator: Few-shot dense retrieval from 8 examples. *arXiv preprint arXiv:2209.11755*.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. **The faiss library**.
- Reza Esfandiarpour, George Zerveas, Ruochen Zhang, Macton Mgonzo, Carsten Eickhoff, and Stephen H Bach. 2025. Beyond contrastive learning: Synthetic data enables list-wise training with multiple levels of relevance. *arXiv preprint arXiv:2503.23239*.
- Luyu Gao, Xueguang Ma, Jimmy J. Lin, and Jamie Callan. 2022. Tevatron: An efficient and flexible toolkit for dense retrieval. *ArXiv*, abs/2203.05765.
- Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. 2022. Accelerate: Training and inference at scale made simple, efficient and adaptable. <https://github.com/huggingface/accelerate>.
- 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*.
- Vitor Jeronimo, Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and Rodrigo Nogueira. 2023. Inpars-v2: Large language models as efficient dataset generators for information retrieval. *arXiv preprint arXiv:2301.01820*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. **Dense passage retrieval for open-domain question answering**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen, Daniel Cer, Jeremy R. Cole, Kai Hui, Michael Boratko, Rajvi Kapadia, Wen Ding, Yi Luan, Sai Meher Karthik Duddu, Gustavo Hernandez Abrego, Weiqiang Shi, Nithi Gupta, Aditya Kusupati, Praatek Jain, Siddhartha Reddy Jonnalagadda, Ming-Wei Chang, and Iftekhhar Naim. 2024. **Gecko: Versatile text embeddings distilled from large language models**. *Preprint*, arXiv:2403.20327.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, and 13 others. 2021. **Datasets: A community library for natural language processing**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 175–184, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xiaopeng Li, Xiangyang Li, Hao Zhang, Zhaocheng Du, Pengyue Jia, Yichao Wang, Xiangyu Zhao, Huifeng Guo, and Ruiming Tang. 2024. Syneg: Llm-driven synthetic hard-negatives for dense retrieval. *arXiv preprint arXiv:2412.17250*.
- Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. Fine-tuning llama for multi-stage text retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2421–2425.
- Craig Macdonald and Nicola Tonellotto. 2020. Declarative experimentation in information retrieval using pyterrier. In *Proceedings of ICTIR 2020*.

- Gabriel de Souza P Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. 2024. *Nv-retriever: Improving text embedding models with effective hard-negative mining*. *arXiv preprint arXiv:2407.15831*.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. *Mteb: Massive text embedding benchmark*. *arXiv preprint arXiv:2210.07316*.
- Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. *RocketQA: An optimized training approach to dense passage retrieval for open-domain question answering*. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5835–5847, Online. 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.
- Navid Rekasaz, Oleg Lesota, Markus Schedl, Jon Brasse, and Carsten Eickhoff. 2021. *TripClick: The Log Files of a Large Health Web Search Engine*. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2507–2513. Association for Computing Machinery, New York, NY, USA.
- Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. *BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models*. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. *Improving text embeddings with large language models*. *arXiv preprint arXiv:2401.00368*.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. *Improving text embeddings with large language models*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11897–11916, Bangkok, Thailand. Association for Computational Linguistics.
- Orion Weller, Benjamin Van Durme, Dawn Lawrie, Ashwin Paranjape, Yuhao Zhang, and Jack Hessel. 2024. *Promptretriever: Instruction-trained retrievers can be prompted like language models*. *Preprint*, arXiv:2409.11136.
- 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, and 3 others. 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.
- Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. *Approximate nearest neighbor negative contrastive learning for dense text retrieval*. *arXiv preprint arXiv:2007.00808*.
- George Zerveas, Navid Rekasaz, Daniel Cohen, and Carsten Eickhoff. 2022. *CODER: An efficient framework for improving retrieval through CONTEXTUAL Document Embedding Reranking*. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10626–10644, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- George Zerveas, Navid Rekasaz, and Carsten Eickhoff. 2023. *Enhancing the ranking context of dense retrieval through reciprocal nearest neighbors*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10779–10803, Singapore. Association for Computational Linguistics.
- Jingtao Zhan, Jiabin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. *Optimizing dense retrieval model training with hard negatives*. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*, pages 1503–1512.

## A Limitations

Trove is mainly aimed at researchers, and our design emphasizes a simple codebase that is freely customizable. Since embedding models are a major component of IR pipelines and are even frequently used in RAG applications without additional steps like re-ranking, Trove provides extensive support for training and evaluating dense embedding models out of the box. However, while we envision adding support for other major retrieval methods such as sparse or hybrid retrieval, Trove is not designed to support all variants of different retrieval methods out of the box. Instead, Trove’s design makes it easy to implement custom variations of different methods in user scripts. On the other end of the spectrum, there is SentenceTransformers, a great library that makes it easy to get started with IR experiments, with built-in support for many retrieval methods. Inevitably, such a codebase is complex and hard to modify by users, which is what we are trying to avoid with Trove.

Moreover, to facilitate exploration and rapid experiments, we sometimes avoid industry standards and implement our own solutions. For example, we implement a fast embedding cache and our own Pytorch container (i.e., `FastResultHeapq`) to speed up nearest neighbor search instead of using tools such as FAISS (Douze et al., 2024), which are also used by other libraries like Tevatron and PyTerrier. While FAISS has many great features, in our case, creating a large search index that is only used once is not as efficient. Moreover, such dependencies limit flexibility. For example, FAISS only reports the similarity scores of the most similar documents for each query. On the other hand, our method can also track similarity scores for arbitrary documents even if they are not ranked among top-k results but, for example, are useful for answering a specific research question.

## B Other Libraries

Here, we discuss Trove’s main differences with two other popular libraries in the information retrieval space.

**Pyterrier** While Pyterrier (Macdonald and Tonello, 2020) is a mature library for composing general IR pipelines, its support for dense retrieval (via plugins<sup>6</sup>) is limited. On the other hand, Trove mainly focuses on reducing the computational overhead and complexity of IR experiments with dense embedding models. Unlike Pyterrier, Trove supports multi-node evaluation of dense retrieval models out of the box, without any additional effort. Furthermore, Pyterrier mainly relies on Pandas Dataframes to represent the data which consumes significant memory in distributed settings where each process loads a separate copy of the data. By contrast, Trove offers dedicated dataset classes optimized for training, that implement an efficient memory-mapped representation of data by leveraging technologies like Apache Arrow and Polars. Finally, Trove is designed specifically for dense retrieval and offers abstractions that simplify experimenting with different modeling choices like different loss functions and encoder architectures.

**ir\_datasets** Another popular library is `ir_datasets`<sup>7</sup>, which is dedicated to providing a unified interface for loading IR datasets. By contrast, Trove encompasses the full pipeline for training and evaluating dense embedding models,

including data processing, hard negative mining, modeling, training, and evaluation. However, Trove and libraries like `ir_datasets` are not mutually exclusive. Trove can integrate different data sources and provides simple decorators that allow users to register and use custom data loader functions. For example, users can integrate new data loader functions that use `ir_datasets` for reading the data files.

---

<sup>6</sup>[https://github.com/terrierteam/pyterrier\\_dr](https://github.com/terrierteam/pyterrier_dr)

<sup>7</sup>[https://github.com/allenai/ir\\_datasets](https://github.com/allenai/ir_datasets)

# CLINICALTRIALSHUB: Bridging Registries and Literature for Comprehensive Clinical Trial Access

Jiwoo Park Ruoqi Liu Avani Jagdale Andrew Srisuwananukorn  
Jing Zhao Lang Li Ping Zhang Sachin Kumar

The Ohio State University

park.3620@osu.edu

## Abstract

We present CLINICALTRIALSHUB, an interactive search-focused platform that consolidates all data from ClinicalTrials.gov and augments it by automatically extracting and structuring trial-relevant information from PubMed research articles. Our system effectively increases access to structured clinical trial data by 83.8% compared to relying on ClinicalTrials.gov alone, with potential to make access easier for patients, clinicians, researchers, and policymakers, advancing evidence-based medicine. CLINICALTRIALSHUB uses large language models such as GPT-5.1 and Gemini-3-Pro to enhance accessibility. The platform automatically parses full-text research articles to extract structured trial information, translates user queries into structured database searches, and provides an attributed question-answering system that generates evidence-grounded answers linked to specific source sentences. We demonstrate its utility through (1) a user study involving clinicians, clinical researchers, and PhD students of pharmaceutical sciences and nursing, and (2) a systematic automatic evaluation of its information extraction and question answering capabilities.<sup>1</sup>

## 1 Introduction

Access to clinical trial information is essential for patients seeking new treatments and for clinicians, researchers, and policymakers working to advance medical care. Containing over 500K registrations, [ClinicalTrials.gov](https://clinicaltrials.gov) (CTG) is the primary resource many of these groups use to search and identify ongoing and completed trials. All trials in CTG are available in a structured format allowing the

<sup>1</sup>A demonstration video is available at: <https://www.youtube.com/watch?v=uCPxyw7Abh0>. The source code to run the demo locally is available at: <https://github.com/jiwoo-jus/clinical-trials-hub>. Evaluation code for comparing model performance is available at: <https://github.com/jiwoo-jus/clinical-trials-hub-evaluation>.

Database	Filter/Status	Count
CTG	All Studies	543,172
CTG	Completed w/ Results	60,449
PubMed	Max Sensitivity <sup>a</sup>	962,774
PubMed	Max Specificity <sup>a</sup>	138,279
PubMed	Max Sensitivity + CTG linked <sup>b</sup>	63,928
PubMed	Max Specificity + CTG linked <sup>b</sup>	35,671

Table 1: Clinical trial data coverage comparison. <sup>a</sup>Identified using PubMed’s Clinical Queries search strategies based on (Wilczynski et al., 2011) with PMC Open Access restriction. *Max Sensitivity* prioritizes recall; *Max Specificity* prioritizes precision. <sup>b</sup>Subset already registered on ClinicalTrials.gov.

users to easily navigate, filter, or download them. However, many trials remain unregistered or are reported only in publications, particularly those conducted outside the US. [PubMed](https://pubmed.ncbi.nlm.nih.gov/), which indexes over 35 million biomedical papers and abstracts, often includes results from these unregistered trials. As shown in [Table 1](#), a substantial number of trials accessible from PubMed are not registered on CTG. The difference in coverage reflects the lack of integration between trial registries and published literature. Despite both being critical resources, PubMed and CTG exist as isolated silos with incompatible formats. PubMed, containing free-text articles, can be especially difficult for non-researchers to parse or filter, requiring substantial manual effort that could take weeks or months.

To address this gap, we present CLINICALTRIALSHUB, a unified user-centered trial search platform that combines trial registry data from CTG and published literature from PubMed into a single interface (§2). Our system makes three main contributions. First, our interface accepts natural language queries and displays unified search results combining trials from both sources in a single ranked list, merging overlapping trials (§3). Second, using large language models (LLMs), we

extract **structured information** from free text PubMed articles making up to 899,846 previously unregistered clinical trials easily searchable and filterable (§4). Lastly, we build an attributed **question answering** feature that allows the users to ask questions over individual trials (§5); answers are generated with attribution provided in the trial text.

Our system showcases practical use cases spanning diverse medical contexts, from evidence curation and systematic reviews to dosing protocols and treatment evaluation. We validate the system’s utility through comprehensive evaluations comparing frontier LLMs (Gemini-3-Pro (Gemini Team, Google, 2025), GPT-5.1 (OpenAI, 2025), Claude-4.5-Sonnet (Anthropic, 2025)) to select optimal models for each feature. We conduct quantitative benchmarking to assess information extraction accuracy using a curated benchmark and evaluate question-answering quality using the FACTS grounding benchmark (Jacovi et al., 2025). Additionally, we perform a user study with seven medical professionals. The evaluations confirm the platform’s effectiveness across its core functionalities and showcase its utility over using PubMed or CTG alone.

## 2 CLINICALTRIALSHUB UI Design

This section describes the user interface of our system. The platform consists of two primary pages: a main search page and individual trial detail pages.

**Main Search Page.** The main search interface (an overview is shown in Figure 1; additional screenshots are available in Appendix G) offers a unified search experience. In **search bar panel**, users can enter natural language queries (top) or use structured input forms specifying condition, intervention, and other terms for PICO-based search (middle), or utilize expert-level PubMed/CTG query forms (bottom). The **search results** display a ranked list of trials from both CTG and PubMed sources, with merged entries when a trial appears in both databases. A **filtering sidebar** enables refinement by data source, study type, phase, design allocation, and other trial-specific categories. The **details sidebar** provides a quick preview of abstracts and metadata when clicking trial titles. CLINICALTRIALSHUB also provides two experimental features to enhance the search experience. Users can specify **inclusion and exclusion criteria** for early-stage screening; when viewing a trial’s preview sidebar, the system provides a quick eligi-

bility assessment based on the abstract and metadata. Additionally, an **AI-insight** feature generates five context-aware insights based on the current search results, supporting iterative Q&A by maintaining context from the user’s queries and recent interactions.

**Individual Trial Detail Page.** Upon selecting a trial, users access a comprehensive detail page (an overview is shown in Figure 2; additional screenshots are provided in Appendix H). Each entry is categorized as *CTG-only*, *PubMed-only*, or *merged* when bidirectional references link a CTG registration to PMC articles. The detail page provides three main components: **Structured Information** displays trial metadata in an easy-to-navigate format, directly available for CTG entries and extracted from free text for PubMed-only items; **Full Text** shows the publication text for PubMed entries, while for CTG entries referenced PMC articles can be expanded inline. Since mappings are not always one-to-one—interim reports or sub-studies may reference a single registration, or multiple registrations across trial phases may cite the same final article—the detail page surfaces all linked records and allows toggling between them; information extraction is performed independently for each article. **Interactive QA** generates evidence-grounded responses linked to specific source sentences that auto-scroll and highlight upon clicking a citation.

## 3 CLINICALTRIALSHUB Search

Our search system comprises three components: (1) query refinement that converts natural language queries to platform-specific structured queries, (2) multi-source search on the two platforms, (3) relevance reranking and deduplication. Figure 3 illustrates the overall pipeline.

### 3.1 Query Refinement

The query refinement component translates user queries into structured search parameters. Our implementation handles both structured field inputs and natural language queries. When users provide individual field values (condition, intervention, other terms), we use them directly. For natural language queries, we use an LLM (GPT-5.1) that applies clinical terminology normalization to decompose the required fields. We design prompt templates to extract condition terms, intervention specifications, and auxiliary keywords. The prompt

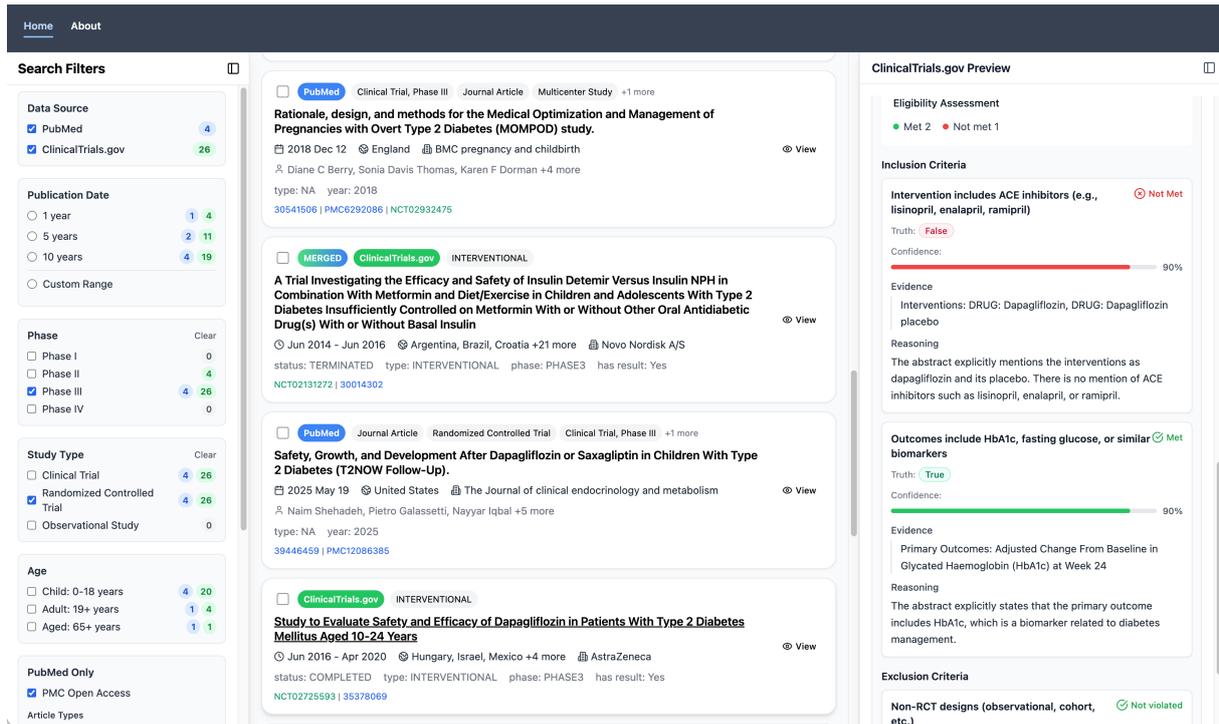


Figure 1: Main search interface with unified results, filters, and eligibility preview

used to extract these meanings from user queries is shown in Appendix I.

The extracted output maps directly to database-specific query construction, where PubMed searches combine the structured parameters with Boolean operators and predefined clinical trial detection patterns, while CTG requires the decomposed fields as separate API parameters. This approach maintains consistent query semantics across heterogeneous database interfaces while preserving the original search intent.

### 3.2 Multi-Source Search

Both CTG and PubMed offer their own search APIs. Given the refined query, we execute both APIs asynchronously with platform-specific optimization strategies. PubMed retrieval uses NCBI's E-utilities API in a two-stage process (E-SEARCH for Boolean query execution and PMID extraction, followed by E-FETCH for batch metadata retrieval including abstracts, MeSH terms, and cross-references). CTG integration combines the official REST API for basic trial information with the daily-updated AACT PostgreSQL mirror for enhanced metadata and complex relational queries.

### 3.3 Reranking and Result Integration

We apply BM25 scoring (Robertson and Walker, 1994) separately to each source, normalize scores,

and add bonuses based on the similarity ranking within each source. We bidirectionally merge PubMed and CTG entries that reference each other, providing a relevance bonus to merged pairs. All entries are then sorted by final scores and displayed to users. Complete technical details including corpus construction, scoring formulas, and deduplication algorithms are provided in Appendix B.

## 4 Information Extraction

All entries on CTG follow a structured format that allows users to navigate and filter across different trials as well as download them for analyses. For example, a patient looking for phase 3 trials for their gender and age range can easily apply a filter and find relevant trials. On the other hand, PubMed articles are research papers in free text form which makes it harder to perform such tasks. Hence, we implement an information extraction pipeline that extracts all relevant trial-related data from papers in the same format as CTG, making indexing, filtering, browsing, and analysis through them easier. We employ an LLM for this purpose that takes as input (parts of) the paper and generates the structured output in a JSON format.

We perform information extraction for 211 fields in parallel by dividing the task into six modules for the protocol section, four modules for the re-

### ChatBot

**Q:** Did the trial include participants who were already using basal insulin?

**A:** Yes, the trial included participants who were already using insulin as part of their standard of care. The inclusion criteria specified that eligible participants could be on a stable dose of insulin, either alone or in combination with metformin. Furthermore, the study protocol noted that basal insulin could be up-titrated as a rescue measure, implying its potential presence in the treatment regimen.

**Evidence:**

- Inclusion criteria were HbA1c 6.5–11%, FPG  $\leq 14.2$  mmol/L ( $\leq 256$  mg/dL) and a stable dose of either metformin ( $\geq 1000$  mg daily), insulin, or a combination of metformin ( $\geq 1000$  mg daily) plus insulin for a minimum of 8 weeks.
- Basal insulin was initiated or up-titrated as open-label rescue for participants meeting prespecified criteria for lack of glycemic control

Ask about PMC10851108

### Structured Information

Protocol Section **Results Section** Document Section Derived Section

**Participant Flow**

Recruitment Details  
Participants took part in the study at 42 study centres in 7 countries worldwide.

Arm/Group Information

Id	Title	Description
FG000	Dapagliflozin 10mg/ Dapagliflozin 10mg	Dapagliflozin (10 mg) tablet administered orally, once daily for the 24 week double-blinded treatment period. The participants then continued to receive Dapagliflozin (10 mg) once daily for a further 28 weeks in the open label long term-extension.
FG001	Placebo/ Dapagliflozin 10mg	Matching placebo tablet administered orally, once daily for the 24 weeks double-blinded treatment period. The participants then received Dapagliflozin (10 mg), orally, once daily for a further 28 weeks in the open label long-term extension.

Study Periods

Discrete stages of a clinical study during which numbers of participants at specific significant events or points of time are reported.

Arm/Group Title	Dapagliflozin 10mg/ Dapagliflozin 10mg	Placebo/ Dapagliflozin 10mg
<b>Blinded Treatment Period</b>		
Started	39	33
Received Treatment	39	33
Completed	34	27
Not Completed	5	6
<b>Long Term Extension</b>		
Started	33	27
Completed	32	24
Not Completed	1	3

Baseline Characteristics

A description of each baseline or demographic characteristic measured in the clinical study

Baseline Groups

Id	Title	Description
FG000	Dapagliflozin 10mg/	Dapagliflozin (10 mg) tablet administered orally, once daily for the 24 week double-blinded treatment period. The participants then continued to receive Dapagliflozin (10 mg) once daily for a further 28 weeks in the open label long term-extension.

Figure 2: Trial detail page with structured information and evidence-grounded QA

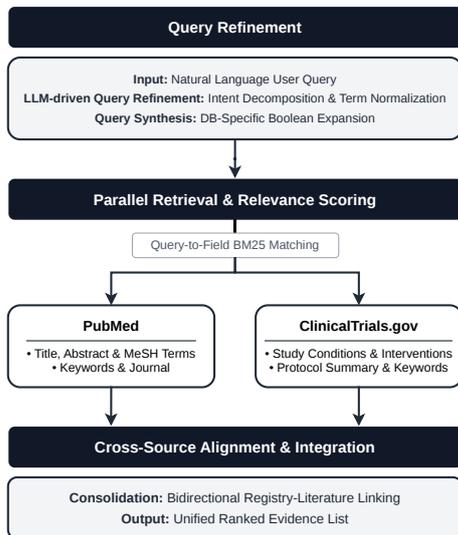


Figure 3: Search pipeline

sults section, and one module for derived section, each with explicit field definitions and enumerated value constraints constructed based on CTG data schema<sup>2</sup> fed to the LLM (see fields list at Table 14 in Appendix C). This approach minimizes both user wait times and hallucinations caused by the LLM processing overwhelming amounts of information. We describe the modules in detail in Appendix C, which also includes the prompt templates (Appendix J).

#### 4.1 Evaluation

We evaluated three frontier LLMs (Gemini-3-Pro, GPT-5.1, and Claude-4.5-Sonnet) to measure structural completeness (key-level: whether the field exists) and semantic accuracy (value-level: whether the extracted value matches the ground truth).

**Dataset construction** We create an evaluation dataset comprising 100 PMC-CTG trial pairs. To do this, we queried PubMed for clinical trials (subsection C.3) published after 2021 (the current data schema’s modernization point) with PMC full-text availability and CTG status set to completed. We then check for a one-to-one mapping verified bidirectionally through the PMC and CTG databases. We focus on the 21 distinct field types in protocol

<sup>2</sup><https://beta-ut.clinicaltrials.gov/api/v2/studies/metadata>

section covering study design, arm group interventions, outcome measures, and eligibility criteria.

We excluded CTG fields appearing in fewer than 15% of records and fields primarily found in PMC articles but inconsistently updated on CTG, considering them unsuitable as reliable ground truth. Since results reporting has been required, limited enforcement has led to results information being more complete in PMC articles, while CTG entries may be missing or outdated. Conversely, CTG often lists extensive secondary outcomes and adverse events, whereas PMC articles emphasize primary outcomes. We plan to create a systematically validated benchmark dataset covering this results section as our future work. Model predictions and ground truth data are available [here](#).

For key-level evaluation, we flatten all fields from both model prediction and reference, preserving hierarchical relationships. For each field in the schema, we classify its extraction status as: *True Positive* (field exists in both prediction and reference), *False Positive* (field appears in prediction but not defined in the CTG data schema), *False Negative* (field exists in reference but missing from prediction), or *Extra Valid* (field appears in prediction and is valid according to the schema, but absent from this specific reference case, making it unevaluable for correctness).

For value-level evaluation, for fields with *True Positive* status, we assess extracted values using a three-stage approach. First, we perform exact string matching with normalization. For mismatches, we calculate semantic similarity using GPT-5.1 as judge with a 0.70 threshold (prompt shown in [subsection C.4](#)). We choose this threshold qualitatively by manually analyzing model outputs. For structured list fields, we create an  $m \times n$  similarity matrix where  $m$  is the number of reference elements and  $n$  is the number of predicted elements to score all reference-prediction element pairs, then use the Hungarian algorithm ([Kuhn, 1955](#)) to optimize total similarity across matched pairs, allowing accurate per-element evaluation regardless of list order.

**Results** Table [Table 2](#) summarizes the overall performance. GPT-5.1 achieved the highest performance across all metrics, recording a Key-level F1 of 0.980 and a Value-level F1 of 0.890. Gemini-3-Pro exhibited competitive performance (Key F1 0.960, Value F1 0.870), whereas Claude-4.5-Sonnet showed a clear performance gap (Key

Level	Model	Precision	Recall	F1
Key	Gemini-3-Pro	<b>1.000</b>	0.930	0.960
	GPT-5.1	<b>1.000</b>	<b>0.970</b>	<b>0.980</b>
	Claude-4.5-Sonnet	0.900	0.690	0.740
Value	Gemini-3-Pro	0.830	0.910	0.870
	GPT-5.1	<b>0.840</b>	<b>0.960</b>	<b>0.890</b>
	Claude-4.5-Sonnet	0.740	0.670	0.660

Table 2: Information extraction performance across key-level and value-level metrics.

F1 0.740, Value F1 0.660). However, we note that the differences between GPT-5.1 and Gemini-3-Pro are small across both key- and value-level metrics. Because we used GPT-5.1 as the semantic similarity judge, a small bias in its favor is possible. To mitigate this, we manually reviewed some of the disagreement cases and confirmed that the judged matches were generally valid. Given the overall closeness of model performance, our results suggest that either GPT-5.1 or Gemini-3-Pro could serve as a reasonable backbone model in practice. We select GPT-5.1 for our system primarily for cost-efficiency and practical deployment considerations rather than because of any substantial performance gap.

A more fine-grained view is presented in [Table 13](#), which reports value-level F1 scores for each individual field alongside the average number of Extra-Valid (EV) cases. Although EV cases cannot be directly evaluated due to the lack of a reference value, they nevertheless reveal meaningful patterns about the underlying data sources. Several descriptive and eligibility-related fields exhibit high EV frequencies, indicating that PMC articles often contain information that is absent in CTG. Importantly, these same fields also show high F1 scores when CTG does provide a reference value, suggesting that discrepancies arise primarily from CTG incompleteness rather than model hallucination. This pattern reinforces the complementary value of LLM-based extraction when dealing with fields that are inconsistently maintained or only partially populated in CTG.

## 5 Interactive Grounded QA

On the trial details page, we provide an LLM-powered chat interface that allows the user to ask question about individual trials with evidence-grounded responses. To select the optimal model for this component, we evaluated three frontier LLMs (Gemini-3-Pro, GPT-5.1, Claude-4.5-

Sonnet).

**Benchmark Selection** For this evaluation, we used the FACTS Grounding benchmark (Jacovi et al., 2025), which measures a model’s ability to generate responses factually grounded in provided context documents. The benchmark aligns directly with our requirements by incorporating substantial medical-domain coverage (29%), supporting long context documents up to 32K tokens with diverse query patterns, and explicitly constraining models to rely solely on the provided context.

**Evaluation Methodology** FACTS employs a two-phase evaluation conducted by three judge models (Gemini-1.5-Pro, GPT-4o, and Claude-3.5-Sonnet). **1. Instruction Following (Eligibility):** Judges assess whether responses adequately address the user request, assigning verdicts of *No issues*, *Minor issue(s)*, or *Major issue(s)*. A response is deemed ineligible (and assigned a score of 0) if all three judges classify it as having major issues, preventing it from advancing to the grounding evaluation phase. **2. Grounding:** Judges evaluate whether each response is fully grounded in the source document by classifying every sentence as *supported*, *unsupported*, *contradictory*, or *no\_rad* (requires no factual grounding). A response receives a positive grounding verdict from a judge only if all sentences are either supported or *no\_rad*. The factuality score for each response is the average of the three judges’ verdicts, and each model’s final score is the mean factuality score across all evaluated responses.

**Evaluation setup** We evaluated Gemini-3-Pro, GPT-5.1, and Claude-4.5-Sonnet on 236 medical-domain samples from FACTS, generating all responses but evaluating the top 100 shortest average response length to balance cost and significance following FACTS guideline.

For grounding evaluation, we used the *JSON* prompt template from the available templates, which FACTS identified as optimal for Gemini-1.5-Pro and GPT-4o. This JSON template instructs the model to output structured JSON objects for each sentence, with classification labels, explicit rationales explaining each decision, and supporting evidence excerpts from the context document.

In contrast, the *implicit-span-level* prompt template generates unstructured natural language output, listing each sentence with a simple binary accurate/inaccurate label and concluding with a final

Model \ Judge	Average	Gemini	GPT	Claude
Gemini-3-Pro	<b>0.897</b>	<b>0.920</b>	<b>0.830</b>	<b>0.930</b>
GPT-5.1	0.680	0.670	0.610	0.760
Claude-4.5-Sonnet	0.767	0.740	0.690	0.860

Table 3: Evaluation scores by judge model. Each row shows a prediction model evaluated by three different judge models (Gemini-3-Pro, GPT-5.1, Claude-4.5-Sonnet) and their average score.

Model	Tokens	Time	Length
Gemini-3-Pro	1318.44	14.33	698.01
GPT-5.1	227.71	2.81	1006.82
Claude-4.5-Sonnet	210.95	5.34	896.54

Table 4: Performance comparison. **Tokens:** Avg. completion tokens; **Time:** Avg. latency (s); **Length:** Avg. output characters.

verdict, without providing detailed rationales or supporting evidence. In FACTS, this approach had previously been preferred for Claude-3.5-Sonnet due to its lower complexity in structured output generation. However, Claude-4.5-Sonnet exhibited substantially improved JSON generation reliability, allowing us to use the more informative template format (see Table 15).

**Results** Gemini-3-Pro achieved the highest grounding score (0.897), notably outperforming GPT-5.1 (0.680) and Claude-4.5-Sonnet (0.767), with consistently stronger performance across all three evaluators (see Table 3). We additionally recorded completion tokens, response time, and response length for all generated responses. As shown in Table 4, while Gemini-3-Pro required significantly higher completion tokens (1,318 avg. vs. ~220 for others) and generation time (14.33s vs. ~4s), it produced the most concise final answers (698 characters vs. ~950). This suggests that Gemini performs a deeper internal reasoning process to synthesize grounded, compact responses; this behavior is critical for trustworthy clinical QA. Consequently, we selected Gemini-3-Pro as our QA backbone.

## 6 User Study

To assess the practical utility of CLINICALTRIAL-SHUB, we conducted an initial user study with seven medical professionals: hematologists, pathologists, dentists, clinical statisticians, pharmacists, pharmaceutical scientists, and nursing researchers.

Participants explored the system for tasks corresponding to their typical use of PM and CTG. [Table 5](#) summarizes key findings. Detailed study materials are provided in [Appendix D](#).

Metric	Score
<i>Search Feature Experience (0–5):</i>	
Query generation rating	4.50
Filtering rating	4.40
Eligibility check rating	4.50
Combined search utility	4.14
<i>Information Extraction Accuracy (0–5):</i>	
Study Overview	5.0
Study Plan	5.0
Participation Requirements	4.83
Baseline Characteristics	4.83
Outcome Measures	4.83
Participant Flow	4.7
Adverse Events	4.6
Chatbot answer quality (0–5)	4.86

Table 5: User study results summary. Scores reflect participant ratings of ClinicalTrialsHub functionality.

**Search Stage.** Six of seven participants found 6 or more relevant studies with CLINICALTRIALSHUB compared to 5/7 for PubMed and 3/7 for CTG among the top 30 results, demonstrating improved relevance through unified search with BM25 reranking. Query generation (4.50), filtering (4.40), and eligibility checking (4.50) features received high ratings. The overall combined search capability received strong approval (4.14), validating the value of eliminating cross-platform navigation. Time-saving received moderate rating (3.71), potentially reflecting initial learning curve.

**Review Stage.** Six of seven participants reviewed the structured extraction from PubMed. Perceived accuracy was consistently high: Study Overview and Study Plan received perfect scores (5.0), while Participation Requirements, Baseline Characteristics, and Outcome Measures scored 4.83. Results-oriented sections (Participant Flow: 4.7, Adverse Events: 4.6) also rated highly. All seven participants used the chatbot, rating answer quality (4.86) and overall detail page efficiency (4.86) very highly.

## 7 Related Work

Prior work has aimed to improve clinical trial information access and evidence synthesis in various ways. Trialstreamer ([Marshall et al., 2020](#)) structures PubMed articles for rapid evidence browsing. RobotReviewer ([Marshall et al., 2016](#)) auto-

mates trial data extraction and risk-of-bias evaluation specifically for systematic reviews. LinkedCT ([Hassanzadeh et al., 2009](#)) transforms ClinicalTrials.gov data into structured linked data though it does not integrate literature sources, unlike our work. Prior work studied the use of LLMs this space as well—for patient-trial matching, clinical trial design, and participant recruitment ([Wang et al., 2024](#)). Most recently, TrialPanorama ([Wang et al., 2025a](#)) established a large-scale database and benchmark for these tasks, while LEADS ([Wang et al., 2025b](#)) introduced a foundation model specifically designed to enhance human-AI collaboration in medical literature mining. CLINICALTRIALSHUB extends these approaches into a unified platform, bridging registry and literature silos. Unlike Trialstreamer or LinkedCT, it integrates both data sources, and compared to RobotReviewer’s narrow focus on systematic reviews, it supports broader interactive exploration and structured retrieval for diverse clinical and research tasks.

## 8 Conclusion and Future Work

We presented CLINICALTRIALSHUB, a unified platform that integrates structured trial registry data from ClinicalTrials.gov with structured information extracted from unstructured PubMed publications using LLMs. Our system enhances access to comprehensive clinical trial information by enabling unified search, structured information extraction and attributed question answering, supporting the diverse needs of patients, clinicians, and researchers.

Several directions remain for future work. Our extraction evaluation currently covers 21 protocol-level fields; extending this to results-related fields with validated benchmarks is an important next step. We also plan to systematically compare our BM25-based retrieval with dense retrieval methods, conduct larger-scale user studies, and improve extraction accuracy through domain-adapted fine-tuning and error analysis across field types.

## References

- Anthropic. 2025. [Claude 4.5 sonnet system card](#). Technical report, Anthropic.
- Gemini Team, Google. 2025. [Gemini 3: A multimodal family of capable models](#). Technical report, Google DeepMind.
- Oktie Hassanzadeh, Anastasios Kementsietsidis, Lipyeow Lim, Renée J Miller, and Min Wang. 2009. [Linkedct: A linked data space for clinical trials](#). Technical Report CSRG-596, University of Toronto.
- Alon Jacovi, Andrew Wang, Chris Alberti, Connie Tao, Jon Lipovetz, Kate Olszewska, Lukas Haas, Michelle Liu, Nate Keating, Adam Bloniarz, Carl Saroufim, Corey Fry, Dror Marcus, Doron Kukliansky, Gaurav Singh Tomar, James Swirhun, Jinwei Xing, Lily Wang, Madhu Gurusurthy, and 7 others. 2025. [The facts grounding leaderboard: Benchmarking llms’ ability to ground responses to long-form input](#). *Preprint*, arXiv:2501.03200.
- Harold W. Kuhn. 1955. [The hungarian method for the assignment problem](#). *Naval Research Logistics Quarterly*, 2(1–2):83–97.
- Iain J Marshall, Joël Kuiper, and Byron C Wallace. 2016. [Robotreviewer: evaluation of a system for automatically assessing bias in clinical trials](#). *Journal of the American Medical Informatics Association*, 23(1):193–201. Originally published online 2015.
- Iain J Marshall, Benjamin Nye, Joël Kuiper, Anna Noel-Storr, Rachel Marshall, Rory Maclean, Frank Soboczenski, Ani Nenkova, James Thomas, and Byron C Wallace. 2020. [Trialstreamer: A living, automatically updated database of clinical trial reports](#). *Journal of the American Medical Informatics Association*, 27(12):1903–1912.
- OpenAI. 2025. [Gpt-5.1 system card](#). Technical report, OpenAI.
- Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR’94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University*, pages 232–241. Springer.
- Zifeng Wang, Lang Cao, Benjamin Danek, Qiao Jin, Zhiyong Lu, and Jimeng Sun. 2024. [Accelerating clinical evidence synthesis with large language models](#). *Preprint*, arXiv:2406.17755.
- Zifeng Wang, Qiao Jin, Jiacheng Lin, Junyi Gao, Jathurshan Pradeepkumar, Pengcheng Jiang, Benjamin Danek, Zhiyong Lu, and Jimeng Sun. 2025a. [Tri-alpanorama: Database and benchmark for systematic review and design of clinical trials](#). *Preprint*, arXiv:2505.16097.

Zifeng Wang, Qiao Jin, Weicheng Ma, and 1 others. 2025b. [Leads: a foundation model for medical literature mining](#). *Nature Communications*, 16(1):1–16.

Nancy L Wilczynski, Kathleen A McKibbin, and Robert B Haynes. 2011. [Sensitive clinical queries retrieved relevant systematic reviews as well as primary studies: an analytic survey](#). *Journal of Clinical Epidemiology*, 64(12):1341–1349.

## A Additional Statistics Details

Our data expansion is calculated as:

$$\text{Expansion}(\%) = \frac{\text{PMC trials without CTG}}{\text{Total PMC trials}} \times 100 \quad (1)$$

This yields expansion rates of 74.2% (specificity: 102,608/138,279) to 93.4% (sensitivity: 898,846/962,774), averaging 83.8%.

## B Reranking and Result Integration Details

We implement BM25-based reranking mechanisms (Robertson and Walker, 1994) to individual search results, followed by deduplication to merge related publications and trials.

**Corpus Construction:** For PubMed documents, we construct searchable text by concatenating: (1) article titles, (2) abstracts (structured or unstructured), (3) author-provided keywords, (4) Medical Subject Headings (MeSH) descriptors, and (5) journal names. For CTG documents, we concatenate: (1) trial titles, (2) condition specifications, (3) brief summaries, and (4) keywords.

**BM25 Scoring:** During integration, BM25 scores are computed for each source separately, then min-max normalized to the  $[0, 1]$  range and combined with a position-based bonus from the original API ranking. For each document  $d_i$ , the score is computed as:

$$\frac{\text{BM25}(d_i) - \text{BM25}_{\min}}{\text{BM25}_{\max} - \text{BM25}_{\min}} + 0.2 \cdot \frac{N - i}{N}$$

where  $\text{BM25}(d_i)$  represents the raw BM25 score,  $\text{BM25}_{\min}$  and  $\text{BM25}_{\max}$  are the minimum and maximum scores for normalization,  $N$  is the total number of results, and  $i$  is the zero-indexed position in the original API ranking. The constant 0.2 represents the maximum position bonus. When  $\text{BM25}_{\max} = \text{BM25}_{\min}$  (indicating identical scores), the BM25 component is set to 0.0 for all documents, and the final ranking relies entirely on the original API position.

**Bidirectional Deduplication:** We merge PubMed and CTG entries that reference each other.<sup>3</sup> Merged pairs receive a relevance bonus.<sup>4</sup> After rescoring and merging, all entries are sorted by their final scores and shown to the user.

## C Information Extraction Details

### C.1 Modular Prompt System

The extraction framework is organized into 3 main sections corresponding to the CTG API v2.0 structure. We provide the exact prompts in our repository.

- **protocolSection:** Contains basic study information including identification modules (NCT ID, title, sponsors), status modules (enrollment, dates, phases), design modules (study type, allocation, masking), arms and interventions, eligibility criteria, contacts, and locations.
- **resultsSection:** Encompasses study outcomes and results when available, including participant flow, baseline characteristics, outcome measures, adverse events, and statistical analyses.
- **derivedSection:** Includes system-generated data elements such as condition browse modules, intervention browse modules, and MeSH term mappings.

### C.2 Data Validation and Schema Conformance

During extraction, we asynchronously operate a validation pipeline.

- **Schema-level validation** enforces structural integrity across all five data sections. Each extracted field undergoes type checking against CTG's built-in types, including character limits for text fields (briefTitle: 300 chars, officialTitle: 600 chars, eligibilityCriteria: 20,000 chars), ISO 8601 date formatting, and proper nesting for complex structures like baseline characteristics and outcome measures.
- **Enumerated value validation** ensures all categorical fields contain only permitted values. The system validates against ClinicalTrials.gov's comprehensive enum definitions for critical fields

<sup>3</sup>This merging criterion is designed for high precision. We acknowledge that might miss certain trials that may have unidirectional references but do not reference each other (for instance, trials that were registered before the publication but never updated)

<sup>4</sup>We empirically determine this bonus to 0.3.

including studyType, allocation, intervention-Model, phases, sex, and standardized age groups. Non-conforming values trigger immediate correction through fuzzy matching against allowed values.

- **Clinical terminology verification** integrates with the NCBI MeSH API to validate medical subject headings in real-time. The system performs fuzzy matching for condition and intervention terms, automatically suggesting and applying standardized MeSH descriptors. This prevents terminology drift and ensures compatibility with biomedical databases.
- **Statistical consistency checking** validates the coherence of quantitative data, verifying group-measure associations in baseline characteristics, confirming unit consistency across measurements, and validating statistical parameters (means, standard deviations, confidence intervals) for mathematical soundness.

Missing information is handled through deliberate omission rather than placeholder generation, preserving data integrity for downstream analysis.

### C.3 Sample ID Filter Query - PubMed

```
(( "randomized controlled trial"[
  Publication Type]
OR "controlled clinical trial"[
  Publication Type]
OR "randomized"[Title/Abstract]
OR "placebo"[Title/Abstract]
OR "clinical trials as topic"[
  MeSH Terms:noexp]
OR "randomly"[Title/Abstract]
OR "trial"[Title])
NOT ("animals"[MeSH Terms] NOT "
humans"[MeSH Terms])
AND ("english"[Language] OR "
English"[lang])
AND "pubmed pmc open access"[
  Filter]
AND "clinicaltrials gov"[
  Secondary Source ID]
AND "2021/01/01"[Date -
  Publication] : "3000"[Date -
  Publication]
```

### C.4 Semantic Similarity Prompt

```
Compare these clinical trial
field values for '{
```

```
field_explanation}'. Return
only a number 0-1 for semantic
similarity.
Text1: {reference_value}
Text2: {predicted_value}
```

## D User Study Details

To assess the practical utility of ClinicalTrialsHub in realistic workflows, we conducted an initial user study with seven medical professionals. Participants were recruited across diverse roles: hematology, pathology, dentistry, clinical statistics, pharmacy, pharmaceutical science (PhD student), and nursing (PhD candidate). The participant pool comprised a balanced distribution of clinicians (3/7), clinical researchers (2/7), and students or trainees (2/7), thereby representing key stakeholder groups who utilize clinical trial data. Each participant received an explanation of the system's features and was asked to use CLINICALTRIALSHUB for tasks that correspond to their typical use of PM and CTG. The contexts in which participants explored baseline systems (PM/CTG) and CLINICALTRIALSHUB, along with their stated purposes, are described in Table 7, while the actual queries they input for each platform are shown in Table 9.

### D.1 Search Stage Evaluation

**Query Generation** Although only two participants explicitly used the natural language query generation feature, both rated it highly (4.50). This low usage rate is unsurprising given that most participants were already proficient in constructing Boolean queries for PubMed and CTG. Table 8 shows their natural language inputs and the system-generated structured queries. These results demonstrate that the system accurately decomposes complex natural language requests into structured parameters, particularly benefiting less experienced users or those working outside their primary domain.

**Filtering** Five participants used the filtering capabilities and rated this feature highly (4.40). The filters applied spanned temporal constraints (publication date ranges from 5-10 years), study design parameters (study type—interventional/observational, RCT status, data availability requirements (PMC Open Access, CTG with results posted), and population specifications (humans, age restrictions, completion status). Table 10 documents the spe-

cific filter combinations each participant employed. The CLINICALTRIALSHUB's filtering mechanism enabled users to impose identical selection criteria on heterogeneous data sources, reducing the need to mentally translate filter semantics between platform-specific interfaces.

**Eligibility Check** Four participants utilized the eligibility criteria specification feature, which also received a high average rating (4.50). This feature allowed users to define detailed inclusion and exclusion criteria that went beyond simple filter combinations. For instance, one participant specified inclusion criteria requiring "randomized controlled or single-arm registrational trials with  $\geq 50$  patients" that "reported at least one primary outcome," while explicitly excluding "studies including other myeloproliferative neoplasms without fibrosis." A diabetes researcher defined even more granular criteria spanning *intervention type* ("dyadic or family-based behavioral, psychoeducational, or self-management intervention"), *population characteristics* ("adults  $\geq 18$  years with type 2 diabetes"), *outcome requirements* (diabetes self-efficacy, self-management behaviors, dyadic processes, or HbA1c), and *temporal constraints* ("follow-up of at least 3 months"). The complete set of criteria specifications is provided in Table 11. The high rating reflects users' need for nuanced eligibility assessment that cannot be captured through simple keyword filtering alone—particularly for systematic review preparation and evidence synthesis where precise population and study design specifications are critical.

**Overall System Utility** Table 6 shows the distribution of relevant studies participants identified among the top 30 results. 6 of 7 participants found 6 or more relevant studies with CLINICALTRIALSHUB, compared to 5 of 7 for PubMed and only 3 of 7 for CTG. Although CTH demonstrated improved relevance, the maximum response option of '11+' prevents precise quantification of the magnitude of improvement. Nevertheless, our unified search with BM25 reranking successfully integrated heterogeneous data sources into a single ranking system without degrading user satisfaction. This suggests that reducing manual screening effort through cross-source integration is achievable even when reconciling disparate ranking algorithms. The overall combined search capability received strong approval (4.14), validating the value proposition of eliminating cross-platform

Platform	0	1–5	6–10	11+
PubMed	0	2	2	3
ClinicalTrials.gov	0	4	2	1
ClinicalTrialsHub	0	1	3	3

Table 6: Distribution of relevant studies identified among top 30 results.

navigation and manual deduplication. However, the time-saving metric received a more moderate rating (3.71). This may reflect the initial learning curve associated with a new interface or, alternatively, that experienced users already possess efficient workflows for rapidly applying trial filters on baseline systems. We view these responses as an opportunity to identify superior interaction patterns from existing systems and either integrate them into ClinicalTrialsHub or better expose our system’s capabilities to accelerate clinical trial research activities.

## D.2 Review Stage Evaluation

In the review stage, we focused on how well the detail page supported close reading of individual trials. Participants were first asked to rate the accuracy of seven structured sections distilled from PMC (Study Overview, Participation Requirements, Study Plan, Participant Flow, Baseline Characteristics, Outcome Measures, Adverse Events) on a 0–5 scale, then to judge whether this representation and the integrated chatbot detail page helped them interpret studies more efficiently and save time.

**Information Extraction** 6 of 7 participants reported reviewing the structured extraction from PubMed within CLINICALTRIALSHUB. Across those respondents, perceived accuracy was consistently high: both the *Study Overview* and *Study Plan* modules received perfect mean scores of 5.0/5.0, while *Participation Requirements*, *Baseline Characteristics*, and *Outcome Measures* clustered tightly of 4.83. Even the more detailed results-oriented sections, *Participant Flow* (4.7) and *Adverse Events* (4.6, with one non-response), were rated near the top of the scale. These scores suggest that, for the protocol-level fields we benchmark in §4, clinicians also subjectively experience the structured representations as faithful to the source.

**QA** All seven participants used the chatbot on the detail page, and they were asked to provide up

to three concrete examples (paper ID, question, and answer) from their own sessions. Their questions illustrate how the assistant is used as an interpretive layer rather than a generic Q&A tool. Several participants asked design and endpoint focused questions, such as summarizing how the sample size was determined, clarifying what the primary endpoint was and whether it was met, or checking whether specific biomarkers were collected. Others queried safety and practical implications, including requests for the most common adverse events, recommended safest order for clinical steps, or whether particular clinical outcomes showed improvement. Some participants targeted fields in protocol section, such as total enrollment numbers, completion dates, or the geographic distribution of trial sites. Across these diverse uses, participants rated the chatbot’s answers as highly accurate and relevant (4.86) and also agreed that the combined detail page—full text, structured view, and chatbot—helped them interpret study information efficiently (4.86).

Table 7: Participant occupations and task contexts

<b>ID</b>	<b>Occupation</b>	<b>Task for This Evaluation</b>
1	MD, Specializing in Hematology	I am the page editor for the Myelofibrosis evidence page on a hematology platform that aggregates and curates high-quality data from recent hematologic trials. My responsibility is to regularly review newly published clinical trials and update the page with evidence that meets my inclusion criteria. For this evaluation, I need to identify and organize the conclusions of recent myelofibrosis trials that should be added to the Myelofibrosis evidence page.
2	Pharmacist	I am looking for dosing decision evidence for a specific patient who is receiving cefepime and continuous renal replacement therapy (CRRT) in the setting of acute kidney failure, to ensure that the cefepime dose remains below the toxic range while still being effective.
3	PhD Student – pharmaceutical science	I am preparing my candidacy proposal and reviewing preliminary studies exploring aspirin for the prevention of pre-eclampsia in high-risk pregnancies. I would like to examine how prior studies have evaluated placental biomarkers in this context.
4	Statistician	I am designing a clinical trial for a new EGFR-targeted therapy in non-small cell lung cancer. For this evaluation, I need to review existing evidence on EGFR-targeted therapies, including their efficacy and safety profiles, to support the study design.
5	Nurse (PhD Candidate)	I am conducting a systematic review of dyadic and family-based interventions for adults with type 2 diabetes, focusing on diabetes self-management, self-efficacy, and related psychosocial outcomes. My responsibility is to identify recent randomized and registrational trials, screen them against predefined inclusion criteria, and organize the conclusions with respect to dyadic or family-based diabetes self-management interventions.
6	Dentist; General Dentistry	I am preparing a paper that discusses the need to standardize IV moderate sedation training and competency assessment in dental residency programs. For this evaluation, I want to identify clinical trials and observational studies on IV moderate sedation in dentistry, including evidence on patient selection, resident training, and the appropriate order and quantity of medications for successful and safe IV moderate sedation.
7	Pathologist (Dermatopathology / Gynecologic Pathology)	I recently reviewed a case involving a metastatic cutaneous melanoma with a BRAFV600E mutation, confirmed through our molecular service. The patient has completed wide excision and sentinel lymph node evaluation, and the oncology team is now considering adjuvant systemic therapy options. For this case, I would like to review the clinical evidence supporting commonly used adjuvant therapies for metastatic melanoma so I can provide informed context when discussing the pathology findings with the treating oncologists. In particular, I want to find studies evaluating adjuvant systemic treatments used after complete resection of metastatic melanoma, including outcomes such as recurrence risk and treatment-related toxicity.

Table 8: Participant Natural Language Queries and Generated Structured Queries

<b>ID</b>	<b>Natural Language Query</b>	<b>Structured Query (CTH Output)</b>
2	Randomized controlled trials involving continuous renal replacement therapy and cefepime in patients with acute kidney injury	Condition: acute kidney injury Intervention: continuous renal replacement therapy OR cefepime Other terms: randomized controlled trials
7	I want clinical trials evaluating chemotherapy in high-stage melanoma, particularly those reporting BRAFV600E mutation status.	Condition: high-stage melanoma OR BRAFV600E mutation Intervention: chemotherapy

Table 9: Participant Search Queries Across Systems

ID	PubMed Query	ClinicalTrials.gov Query	CTH Query
1	myelofibrosis OR polycythemia vera OR essential thrombocythemia OR myeloproliferative	Condition: myelofibrosis OR polycythemia vera OR essential thrombocythemia OR myeloproliferative	Condition: myelofibrosis OR polycythemia vera OR essential thrombocythemia OR myeloproliferative
2	cefepime and crrt	Condition: Acute Kidney Injury (AKI) Intervention: Continuous Renal Replacement Therapy (CRRT) Other terms: cefepime	–
3	Aspirin AND Pre-eclampsia AND biomarkers	Condition: Pre-eclampsia Intervention: aspirin Other terms: biomarkers	Common: pregnancy Condition: pre-eclampsia Intervention: aspirin Other terms: biomarkers, placental biomarkers
4	EGFR-targeted therapies non small cell lung cancer	Condition: non small cell lung cancer Intervention: EGFR-targeted therapies	Common: “EGFR-targeted therapies non small cell lung cancer” Intervention: EGFR-targeted therapies
5	“type 2 diabetes” OR “type 2 diabetes mellitus” OR T2DM AND (dyadic OR family-based OR spouse OR partner OR dyadic OR caregiver OR self-management OR education)	Condition: Type 2 Diabetes Mellitus (T2DM) Other terms: dyadic / family-based interventions, self-management outcomes	Condition: type 2 diabetes mellitus Other terms: dyadic, family-based, spouse, partner, caregiver, self-management, education
6	IV moderate sedation AND midazolam AND fentanyl	Other terms: IV moderate sedation AND midazolam AND fentanyl	Other terms: Dental Anxiety AND Dental procedures AND IV moderate sedation guidelines
7	metastatic melanoma and BRAFV600E mutation and chemotherapy	Condition: metastatic melanoma Intervention: chemotherapy Other terms: BRAF V600E mutation positive	Common: cutaneous melanoma OR adjuvant melanoma OR resected melanoma OR “stage III melanoma” OR “BRAF V600E melanoma”

Table 10: Participant Filters Used Across Systems

ID	PubMed Filters	CTG Filters	CTH Filters
1	Publication Date (Last 5 years) Phase (III/IV) Article Type (Clinical Trial, RCT) Access (PMC OA) Species (Humans)	Completion Date (Last 5 years) Phase (III/IV) Study Type (Interventional) Has Results (True) Status (Completed)	Date (Last 5 years) Phase (III/IV) Study Type (Clinical Trial, RCT) PubMed (PMC OA, Humans) CTG (With Results, Completed)
2	Age Publication Date	Age Completion Date	Age Date
3	Publication Date	Has Results (True)	Phase Study Type
4	Publication Date (Last 10 years) Article Type (Clinical Trial)	Completion Date (From:11/17/2015)	Date (Last 10 years) Study Type (Clinical Trial)
5	Publication Date (Last 10 years) Article Type (Clinical Trial, RCT) Species (Humans) Age (Adults) Language (English)	Completion Date (Last 10 years) Study Type (Interventional, Randomized) Status (Completed) Has Results (True) Age (Adults)	Date (Last 10 years) Study Type (Clinical Trial, RCT) Age (Adults) PubMed (PMC OA, Humans) CTG (With Results, Completed)

Table 11: Participant Inclusion and Exclusion Criteria

<b>ID</b>	<b>Inclusion Criteria</b>	<b>Exclusion Criteria</b>
1	<p>Randomized controlled or single-arm registrational trials with <math>\geq 50</math> patients</p> <hr/> <p>Reported at least one primary outcome</p>	<p>Other myeloproliferative neoplasms without fibrosis (e.g., PV or ET without MF)</p> <hr/> <p>Prefibrotic MF or early myeloproliferative disease without confirmed fibrosis</p>
2	<p>AKI</p> <hr/> <p>CRRT</p>	
5	<p>Randomized controlled trials or single-arm registrational trials with approximately <math>\geq 50</math> participants</p> <hr/> <p>Adults (<math>\geq 18</math> years) with type 2 diabetes mellitus</p> <hr/> <p>Dyadic or family-based behavioral, psychoeducational, or self-management intervention involving partner/spouse/family</p> <hr/> <p>Reported at least one primary or key outcome (self-efficacy, self-management behaviors, dyadic/family processes, HbA1c, etc.)</p> <hr/> <p>Follow-up <math>\geq 3</math> months after baseline</p>	<p>Studies focusing only on type 1 diabetes, gestational diabetes, or prediabetes</p> <hr/> <p>Purely pharmacologic comparisons without behavioral/educational/family components</p> <hr/> <p>Pediatric/adolescent-only populations</p> <hr/> <p>Case reports, very small pilot studies, qualitative-only studies without quantitative outcomes</p>
6	Dental	General Anesthesia

Table 12: Participant Chatbot Interactions

ID	Chatbot Interaction
1	<p><b>PMC12092244   What was the exact pelabresib dose and schedule?</b> Pelabresib was administered orally at a starting dose of 125 mg once daily for 14 consecutive days followed by a 7-day break, forming a 21-day cycle. Dose modifications between 50–175 mg once daily were permitted based on adverse events and treatment response.</p>
	<p><b>PMC12092244   What was the primary endpoint, and was it met?</b> The primary endpoint was achieving a <math>\geq 35\%</math> reduction in spleen volume from baseline at week 24. This endpoint was met, with a significantly greater proportion of patients in the pelabresib–ruxolitinib arm achieving this reduction compared to the placebo–ruxolitinib arm.</p>
	<p><b>PMC12092244   What were the most common adverse events?</b> The most common treatment-emergent adverse events were hematological—primarily thrombocytopenia and anemia. Among non-hematological events, diarrhea was most frequent in the pelabresib–ruxolitinib arm, whereas constipation was most frequent in the placebo–ruxolitinib arm.</p>
2	<p><b>PMC10576861   Since ACORN showed higher neurotoxicity with Cefepime, should we reconsider de-escalating from Piperacillin-Tazobactam to Cefepime solely based on renal concerns?</b> Based on the ACORN trial findings, reconsidering the practice of de-escalating from Piperacillin-Tazobactam to Cefepime solely for renal protection is warranted. The study demonstrated that there was no significant difference in the highest stage of acute kidney injury or death between the two treatment groups. Furthermore, the trial concluded that Piperacillin-Tazobactam did not increase the incidence of acute kidney injury, whereas Cefepime was associated with a higher incidence of neurological dysfunction.</p>
	<p><b>PMC9463324   The study showed creatinine elevation but no difference in dialysis or mortality. Does this mean the current KDIGO criteria are too sensitive for antibiotic-induced AKI?</b> The study suggests that the creatinine-based KDIGO criteria may identify "pseudotoxicity" rather than true renal injury in the context of Vancomycin and Piperacillin-Tazobactam therapy. The authors propose that the observed creatinine elevation is likely due to the inhibition of tubular secretion rather than damage to the kidney parenchyma, a conclusion supported by the lack of corresponding changes in alternative biomarkers (cystatin C and BUN) and the absence of increased risk for dialysis or mortality.</p>
3	<p><b>NCT03893630   Did the study report any differences in placental biomarker levels between the 81mg and 162mg aspirin groups?</b> Yes, the study reported specific mean changes in placental biomarker levels for the 81mg and 162mg aspirin groups, although statistical significance of the differences was not explicitly detailed in the provided text. For Placental Growth Factor (PIGF), the mean change from baseline to Event 2 was 568.3 pg/ml for the 81mg group and 509.7 pg/ml for the 162mg group. For soluble Intercellular Adhesion Molecule (s-ICAM), the mean change from baseline to Event 2 was 1.6 ng/ml for the 81mg group and 7.9 ng/ml for the 162mg group. Data for IL-6 and TNF were not reported due to insufficient serum samples.</p>
4	<p><b>PMC12139608   Can you summarize how the sample size is determined in this study?</b> The study powered for <math>\geq 85\%</math> to detect a 25% difference in ORR using one-sided <math>\alpha=0.025</math>; planned enrollment was 126 patients.</p>
	<p><b>PMC10318754   Can you summarize the statistical analysis methods used?</b> Analyses included chi-square/Fisher tests and QoL correlations with PFS/OS and sponsorship.</p>
	<p><b>PMC9509093   Summarize tabulated results for all endpoints.</b> Aumolertinib improved PFS (19.3 vs 9.9 mo); ORR 73.8% vs 72.1%; DCR 93.0% vs 96.7%; DoR 18.1 vs 8.3 mo; fewer rash/diarrhea events.</p>
5	<p><b>PMC8266295   Were any clinical outcomes such as HbA1c reported, and if so, what were the main results?</b> The paper does not report clinical outcomes such as HbA1c changes as part of the results. The study focused on family supporter involvement, support techniques, and caregiving experience rather than direct clinical outcomes like HbA1c.</p>
6	<p><b>PMC5858011   What do they recommend as the safest order?</b> Midazolam–first caused fewer desaturation events (90–92%).</p>
7	<p><b>NCT01682083   If OS wasn't clearly significant early on, is the decision to use dabrafenib/trametinib driven mostly by the strong RFS benefit?</b> Yes, the provided results indicate that the primary outcome, Relapse-free Survival (RFS), demonstrated a statistically significant benefit with a hazard ratio of 0.47 and a p-value of <math>&lt; 0.0001</math>. In contrast, the Overall Survival (OS) analysis resulted in a p-value of 0.006, which did not meet the pre-specified significance threshold of 0.000019 for the first interim analysis.</p>

## E Extraction Schema Fields

### E.1 Extraction Schema Evaluation Fields

Piece	# Ref	# EV	Gemini	GPT	Claude
DesignObservationalModel	1	7.67	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
<b>DetailedDescription</b>	<b>64</b>	<b>31.67</b>	<b>1.00</b>	<b>1.00</b>	0.86
StudyType	100	0.00	<b>1.00</b>	<b>1.00</b>	0.85
DesignAllocation	99	0.67	0.99	<b>1.00</b>	0.85
BriefSummary	100	0.00	<b>1.00</b>	<b>1.00</b>	0.83
Condition	100	0.00	<b>1.00</b>	<b>1.00</b>	0.83
DesignInterventionModel	99	0.33	0.98	<b>1.00</b>	0.85
Phase	99	0.67	0.98	<b>1.00</b>	0.85
EnrollmentCount	100	0.00	0.99	0.99	0.84
DesignMasking	99	0.67	0.98	0.99	0.85
NCTId	100	0.00	0.99	<b>1.00</b>	0.82
<b>Keyword</b>	<b>52</b>	<b>44.00</b>	<b>1.00</b>	<b>1.00</b>	0.80
MinimumAge	99	0.67	0.99	0.99	0.77
EligibilityCriteria	100	0.00	0.99	0.99	0.77
HealthyVolunteers	99	1.00	0.99	0.99	0.76
Sex	100	0.00	0.99	0.99	0.77
<b>MaximumAge</b>	<b>54</b>	<b>39.00</b>	0.99	<b>1.00</b>	0.74
StdAge	100	0.00	0.98	0.95	0.76
DesignWhoMasked	59	3.67	0.88	<b>0.94</b>	0.77
BriefTitle	100	0.00	0.71	<b>1.00</b>	0.84
OfficialTitle	100	0.00	0.71	0.81	0.82

Table 13: Average Extra-Valid (EV) case analysis and information extraction performance for each field. **# Ref**: number of evaluation cases (out of 100) where the field appears in CTG ground-truth. **# EV**: average cases across models where the field was Extra-Valid (present in PMC but missing from CTG). Values  $\geq 10$  are bolded. Highest F1 per row is bolded. Rows ordered by average F1.

### E.2 Extraction Schema All Fields

Table 14: Comprehensive list of all fields included in our information-extraction schema.

Index	Piece	Field Index
<b>protocolSection.identificationModule</b>		
1	NCTId	nctId
2	OrgStudyId	orgStudyIdInfo.id
3	OrgStudyIdType	orgStudyIdInfo.type
4	OrgStudyIdLink	orgStudyIdInfo.link
5	SecondaryId	secondaryIdInfos.id
6	SecondaryIdType	secondaryIdInfos.type
7	SecondaryIdLink	secondaryIdInfos.link
8	BriefTitle	briefTitle
9	OfficialTitle	officialTitle
10	Acronym	acronym
11	OrgFullName	organization.fullName
12	OrgClass	organization.class
<b>protocolSection.descriptionModule</b>		
13	BriefSummary	briefSummary

Continued on next page

Table 14 – Continued from previous page

Index	Piece	Field Index
14	DetailedDescription	detailedDescription
<b>protocolSection.conditionsModule</b>		
15	Condition	conditions
16	Keyword	keywords
<b>protocolSection.designModule</b>		
17	StudyType	studyType
18	PatientRegistry	patientRegistry
19	TargetDuration	targetDuration
20	Phase	phases
21	DesignAllocation	designInfo.allocation
22	DesignInterventionModel	designInfo.interventionModel
23	DesignInterventionModelDescription	designInfo.interventionModelDescription
24	DesignPrimaryPurpose	designInfo.primaryPurpose
25	DesignObservationalModel	designInfo.observationalModel
26	DesignTimePerspective	designInfo.timePerspective
27	DesignMasking	designInfo.maskingInfo.masking
28	DesignMaskingDescription	designInfo.maskingInfo.maskingDescription
29	DesignWhoMasked	designInfo.maskingInfo.whoMasked
30	EnrollmentCount	enrollmentInfo.count
31	EnrollmentType	enrollmentInfo.type
<b>protocolSection.armsInterventionsModule</b>		
32	ArmGroupLabel	armGroups.label
33	ArmGroupType	armGroups.type
34	ArmGroupDescription	armGroups.description
35	ArmGroupInterventionName	armGroups.interventionNames
36	InterventionType	interventions.type
37	InterventionName	interventions.name
38	InterventionDescription	interventions.description
39	InterventionArmGroupLabel	interventions.armGroupLabels
<b>protocolSection.outcomesModule</b>		
40	PrimaryOutcomeMeasure	primaryOutcomes.measure
41	PrimaryOutcomeDescription	primaryOutcomes.description
42	PrimaryOutcomeTimeFrame	primaryOutcomes.timeFrame
43	SecondaryOutcomeMeasure	secondaryOutcomes.measure
44	SecondaryOutcomeDescription	secondaryOutcomes.description
45	SecondaryOutcomeTimeFrame	secondaryOutcomes.timeFrame
46	OtherOutcomeMeasure	otherOutcomes.measure
47	OtherOutcomeDescription	otherOutcomes.description
48	OtherOutcomeTimeFrame	otherOutcomes.timeFrame
<b>protocolSection.eligibilityModule</b>		
49	EligibilityCriteria	eligibilityCriteria
50	HealthyVolunteers	healthyVolunteers
51	Sex	sex
52	MinimumAge	minimumAge
53	MaximumAge	maximumAge
54	StdAge	stdAges
55	StudyPopulation	studyPopulation
56	SamplingMethod	samplingMethod

Continued on next page

Table 14 – Continued from previous page

Index	Piece	Field Index
<b>resultsSection.participantFlowModule</b>		
57	FlowPreAssignmentDetails	preAssignmentDetails
58	FlowRecruitmentDetails	recruitmentDetails
59	FlowTypeUnitsAnalyzed	typeUnitsAnalyzed
60	FlowGroupId	groups.id
61	FlowGroupTitle	groups.title
62	FlowGroupDescription	groups.description
63	FlowPeriodTitle	periods.title
64	FlowMilestoneType	periods.milestones.type
65	FlowMilestoneComment	periods.milestones.comment
66	FlowAchievementGroupId	periods.milestones.achievements.groupId
67	FlowAchievementComment	periods.milestones.achievements.comment
68	FlowAchievementNumSubjects	periods.milestones.achievements.numSubjects
69	FlowAchievementNumUnits	periods.milestones.achievements.numUnits
70	FlowDropWithdrawType	periods.dropWithdraws.type
71	FlowDropWithdrawComment	periods.dropWithdraws.comment
72	FlowReasonGroupId	periods.dropWithdraws.reasons.groupId
73	FlowReasonComment	periods.dropWithdraws.reasons.comment
74	FlowReasonNumSubjects	periods.dropWithdraws.reasons.numSubjects
<b>resultsSection.baselineCharacteristicsModule</b>		
75	BaselinePopulationDescription	populationDescription
76	BaselineTypeUnitsAnalyzed	typeUnitsAnalyzed
77	BaselineGroupId	groups.id
78	BaselineGroupTitle	groups.title
79	BaselineGroupDescription	groups.description
80	BaselineDenomUnits	denoms.units
81	BaselineDenomCountGroupId	denoms.counts.groupId
82	BaselineDenomCountValue	denoms.counts.value
83	BaselineMeasureTitle	measures.title
84	BaselineMeasureDescription	measures.description
85	BaselineMeasurePopulationDescription	measures.populationDescription
86	BaselineMeasureParamType	measures.paramType
87	BaselineMeasureDispersionType	measures.dispersionType
88	BaselineMeasureUnitOfMeasure	measures.unitOfMeasure
89	BaselineMeasureDenomUnits	measures.denoms.units
90	BaselineMeasureDenomCountGroupId	measures.denoms.counts.groupId
91	BaselineMeasureDenomCountValue	measures.denoms.counts.value
92	BaselineClassTitle	measures.classes.title
93	BaselineClassDenomUnits	measures.classes.denoms.units
94	BaselineClassDenomCountGroupId	measures.classes.denoms.counts.groupId
95	BaselineClassDenomCountValue	measures.classes.denoms.counts.value
96	BaselineCategoryTitle	measures.classes.categories.title
97	BaselineMeasurementGroupId	measures.classes.categories.measurements.groupId
98	BaselineMeasurementValue	measures.classes.categories.measurements.value
99	BaselineMeasurementSpread	measures.classes.categories.measurements.spread
100	BaselineMeasurementLowerLimit	measures.classes.categories.measurements.lowerLimit
101	BaselineMeasurementUpperLimit	measures.classes.categories.measurements.upperLimit
<b>resultsSection.outcomeMeasuresModule</b>		
102	OutcomeMeasureType	outcomeMeasures.type
103	OutcomeMeasureTitle	outcomeMeasures.title
104	OutcomeMeasureDescription	outcomeMeasures.description

Continued on next page

Table 14 – Continued from previous page

Index	Piece	Field Index
105	OutcomeMeasurePopulationDescription	outcomeMeasures.populationDescription
106	OutcomeMeasureReportingStatus	outcomeMeasures.reportingStatus
107	OutcomeMeasureAnticipatedPostingDate	outcomeMeasures.anticipatedPostingDate
108	OutcomeMeasureParamType	outcomeMeasures.paramType
109	OutcomeMeasureDispersionType	outcomeMeasures.dispersionType
110	OutcomeMeasureUnitOfMeasure	outcomeMeasures.unitOfMeasure
111	OutcomeMeasureCalculatePct	outcomeMeasures.calculatePct
112	OutcomeMeasureTimeFrame	outcomeMeasures.timeFrame
113	OutcomeMeasureTypeUnitsAnalyzed	outcomeMeasures.typeUnitsAnalyzed
114	OutcomeMeasureDenomUnitsSelected	outcomeMeasures.denomUnitsSelected
115	OutcomeGroupId	outcomeMeasures.groups.id
116	OutcomeGroupTitle	outcomeMeasures.groups.title
117	OutcomeGroupDescription	outcomeMeasures.groups.description
118	OutcomeDenomUnits	outcomeMeasures.denoms.units
119	OutcomeDenomCountGroupId	outcomeMeasures.denoms.counts.groupId
120	OutcomeDenomCountValue	outcomeMeasures.denoms.counts.value
121	OutcomeClassTitle	outcomeMeasures.classes.title
122	OutcomeClassDenomUnits	outcomeMeasures.classes.denoms.units
123	OutcomeClassDenomCountGroupId	outcomeMeasures.classes.denoms.counts.groupId
124	OutcomeClassDenomCountValue	outcomeMeasures.classes.denoms.counts.value
125	OutcomeCategoryTitle	outcomeMeasures.classes.categories.title
126	OutcomeMeasurementGroupId	outcomeMeasures.classes.categories.measurements.groupId
127	OutcomeMeasurementValue	outcomeMeasures.classes.categories.measurements.value
128	OutcomeMeasurementSpread	outcomeMeasures.classes.categories.measurements.spread
129	OutcomeMeasurementLowerLimit	outcomeMeasures.classes.categories.measurements.lowerLimit
130	OutcomeMeasurementUpperLimit	outcomeMeasures.classes.categories.measurements.upperLimit
131	OutcomeMeasurementComment	outcomeMeasures.classes.categories.measurements.comment
132	OutcomeAnalysisParamType	outcomeMeasures.analysises.paramType
133	OutcomeAnalysisParamValue	outcomeMeasures.analysises.paramValue
134	OutcomeAnalysisDispersionType	outcomeMeasures.analysises.dispersionType
135	OutcomeAnalysisDispersionValue	outcomeMeasures.analysises.dispersionValue
136	OutcomeAnalysisStatisticalMethod	outcomeMeasures.analysises.statisticalMethod
137	OutcomeAnalysisStatisticalComment	outcomeMeasures.analysises.statisticalComment
138	OutcomeAnalysisPValue	outcomeMeasures.analysises.pValue
139	OutcomeAnalysisPValueComment	outcomeMeasures.analysises.pValueComment
140	OutcomeAnalysisCINumSides	outcomeMeasures.analysises.ciNumSides
141	OutcomeAnalysisCIPctValue	outcomeMeasures.analysises.ciPctValue
142	OutcomeAnalysisCILowerLimit	outcomeMeasures.analysises.ciLowerLimit
143	OutcomeAnalysisCIUpperLimit	outcomeMeasures.analysises.ciUpperLimit
144	OutcomeAnalysisCILowerLimitComment	outcomeMeasures.analysises.ciLowerLimitComment
145	OutcomeAnalysisCIUpperLimitComment	outcomeMeasures.analysises.ciUpperLimitComment
146	OutcomeAnalysisEstimateComment	outcomeMeasures.analysises.estimateComment
147	OutcomeAnalysisTestedNonInferiority	outcomeMeasures.analysises.testedNonInferiority
148	OutcomeAnalysisNonInferiorityType	outcomeMeasures.analysises.nonInferiorityType
149	OutcomeAnalysisNonInferiorityComment	outcomeMeasures.analysises.nonInferiorityComment
150	OutcomeAnalysisOtherAnalysisDescription	outcomeMeasures.analysises.otherAnalysisDescription
151	OutcomeAnalysisGroupDescription	outcomeMeasures.analysises.groupDescription
152	OutcomeAnalysisGroupId	outcomeMeasures.analysises.groupIds
<b>resultsSection.adverseEventsModule</b>		
153	EventsFrequencyThreshold	frequencyThreshold
154	EventsTimeFrame	timeFrame
155	EventsDescription	description

Continued on next page

Table 14 – Continued from previous page

Index	Piece	Field Index
156	EventsAllCauseMortalityComment	allCauseMortalityComment
157	EventGroupId	eventGroups.id
158	EventGroupTitle	eventGroups.title
159	EventGroupDescription	eventGroups.description
160	EventGroupDeathsNumAffected	eventGroups.deathsNumAffected
161	EventGroupDeathsNumAtRisk	eventGroups.deathsNumAtRisk
162	EventGroupSeriousNumAffected	eventGroups.seriousNumAffected
163	EventGroupSeriousNumAtRisk	eventGroups.seriousNumAtRisk
164	EventGroupOtherNumAffected	eventGroups.otherNumAffected
165	EventGroupOtherNumAtRisk	eventGroups.otherNumAtRisk
166	SeriousEventTerm	seriousEvents.term
167	SeriousEventOrganSystem	seriousEvents.organSystem
168	SeriousEventSourceVocabulary	seriousEvents.sourceVocabulary
169	SeriousEventAssessmentType	seriousEvents.assessmentType
170	SeriousEventNotes	seriousEvents.notes
171	SeriousEventStatsGroupId	seriousEvents.stats.groupId
172	SeriousEventStatsNumEvents	seriousEvents.stats.numEvents
173	SeriousEventStatsNumAffected	seriousEvents.stats.numAffected
174	SeriousEventStatsNumAtRisk	seriousEvents.stats.numAtRisk
175	OtherEventTerm	otherEvents.term
176	OtherEventOrganSystem	otherEvents.organSystem
177	OtherEventSourceVocabulary	otherEvents.sourceVocabulary
178	OtherEventAssessmentType	otherEvents.assessmentType
179	OtherEventNotes	otherEvents.notes
180	OtherEventStatsGroupId	otherEvents.stats.groupId
181	OtherEventStatsNumEvents	otherEvents.stats.numEvents
182	OtherEventStatsNumAffected	otherEvents.stats.numAffected
183	OtherEventStatsNumAtRisk	otherEvents.stats.numAtRisk
<b>resultsSection.moreInfoModule</b>		
184	LimitationsAndCaveatsDescription	limitationsAndCaveats.description
185	AgreementPISponsorEmployee	certainAgreement.piSponsorEmployee
186	AgreementRestrictionType	certainAgreement.restrictionType
187	AgreementRestrictiveAgreement	certainAgreement.restrictiveAgreement
188	AgreementOtherDetails	certainAgreement.otherDetails
189	PointOfContactTitle	pointOfContact.title
190	PointOfContactOrganization	pointOfContact.organization
191	PointOfContactEMail	pointOfContact.email
192	PointOfContactPhone	pointOfContact.phone
193	PointOfContactPhoneExt	pointOfContact.phoneExt
<b>derivedSection.conditionBrowseModule</b>		
194	ConditionMeshId	meshes.id
195	ConditionMeshTerm	meshes.term
196	ConditionAncestorId	ancestors.id
197	ConditionAncestorTerm	ancestors.term
198	ConditionBrowseLeafId	browseLeaves.id
199	ConditionBrowseLeafName	browseLeaves.name
200	ConditionBrowseLeafRelevance	browseLeaves.relevance
201	ConditionBrowseBranchAbbrev	browseBranches.abbrev
202	ConditionBrowseBranchName	browseBranches.name
<b>derivedSection.interventionBrowseModule</b>		

Continued on next page

Table 14 – *Continued from previous page*

Index	Piece	Field Index
203	InterventionMeshId	meshes.id
204	InterventionMeshTerm	meshes.term
205	InterventionAncestorId	ancestors.id
206	InterventionAncestorTerm	ancestors.term
207	InterventionBrowseLeafId	browseLeaves.id
208	InterventionBrowseLeafName	browseLeaves.name
209	InterventionBrowseLeafRelevance	browseLeaves.relevance
210	InterventionBrowseBranchAbbrev	browseBranches.abbrev
211	InterventionBrowseBranchName	browseBranches.name

## F QA Evaluation Details

Judge Model	Prompt Template	Macro-F1	Acc.	FPR	FNR	F1 (+)	F1 (-)
Claude 3.5 Sonnet	Span-level	68.85	77.83	20.97	22.38	85.58	52.13
	<b>Implicit span-level</b>	<b>70.24</b>	83.50	45.16	11.34	90.10	50.37
	Response-level	61.88	83.25	72.58	6.69	90.42	33.33
	JSON	56.04	64.78	33.87	35.47	75.64	36.44
	JSON (alt)	55.37	66.75	46.77	30.81	77.91	32.84
	JSON w. double-check	49.50	54.68	25.81	48.84	65.67	33.33
	SimpleQA template	55.39	85.22	88.71	1.45	91.87	18.92
Gemini 1.5 Pro	Span-level	55.84	79.31	79.03	10.17	88.03	23.64
	Implicit span-level	56.66	85.47	87.10	1.45	91.99	21.33
	Response-level	48.82	82.02	95.16	4.07	90.04	7.59
	<b>JSON</b>	<b>71.47</b>	86.95	56.45	5.23	92.48	50.47
	JSON (alt)	66.03	85.96	69.35	4.07	92.05	40.00
	JSON w. double-check	64.89	76.35	37.10	21.22	84.95	44.83
	SimpleQA template	51.54	84.73	93.55	1.16	91.64	11.43
GPT-4o	Span-level	63.08	81.53	64.52	10.17	89.18	36.97
	Implicit span-level	55.43	83.99	87.10	3.20	91.11	19.75
	Response-level	51.54	84.73	93.55	1.16	91.64	11.43
	<b>JSON</b>	<b>69.68</b>	80.54	32.26	17.15	87.83	51.53
	JSON (alt)	66.78	82.02	53.23	11.63	89.28	44.27
	JSON w. double-check	57.62	64.04	17.74	39.24	74.11	41.13
	SimpleQA template	47.04	83.74	98.39	1.45	91.13	2.94

Table 15: Evaluation of judge models and prompt templates, reproduced from the **FACTS Grounding Benchmark** paper (Table 2).

## G User Interface - Search Page

Home About

### Clinical Trials Hub

type2 diabetes insulin child

Condition: Diabetes Intervention: Insulin Other Terms: RCT

PubMed Query: Myelofibrosis AND Randomized Controlled Trial[PT] CTG Query: Hypertension AND ("Cleveland, Ohio") AND recruiting

Select Browsing Mode: Expert ✓ Patient

#### Search Results

PubMed Search Query: ((type 2 diabetes) AND (insulin) AND (child)) AND pubmed pmc open access[Filter]

ClinicalTrials.gov Search Query: Condition: type 2 diabetes | Intervention: insulin | Other terms: child

Total: 1285 Merged: 2 PubMed-only: 998 CTG-only: 285

Figure 4: Search panel

Home About

### Advanced Search Builder

Search Term: Enter term or select from below...

Query Box: (type 2 diabetes OR diabetes mellitus type 2 OR type II diabetes) AND insulin AND (child OR pediatric OR children OR childhood OR minor)

Results Preview: PubMed: 996 CTG: 964 Merged: 4 Total: 1964

#### Refined Condition

Total: 1964 | PM: 997 | CTG: 964 | Merged: 3

type 2 diabetes + diabetes mellitus type 2 + adult-onset diabetes + non-insulin-dependent diabetes mellitus + type II diabetes

(type 2 diabetes OR diabetes mellitus type 2 OR adult-onset diabetes OR non-insulin-dependent diabetes mellitus OR type II diabetes) AND (insulin) AND (child)

#### Refined Intervention

Total: 1961 | PM: 995 | CTG: 961 | Merged: 5

insulin + recombinant insulin + human insulin + insulin therapy + insulin analogs

(type 2 diabetes) AND (insulin OR recombinant insulin OR human insulin OR insulin therapy OR insulin analogs) AND (child)

#### Refined Other Term

Total: 1963 | PM: 996 | CTG: 963 | Merged: 4

child + pediatric + children + childhood + minor

(type 2 diabetes) AND (insulin) AND (child OR pediatric OR children OR childhood OR minor)

Figure 5: Advanced search panel

Home About

Total: 29 Merged: 1 PubMed-only: 3 CTG-only: 25

Merged Result View: PubMed ClinicalTrials.gov Clear Selections (0) Select All Download Selected

PubMed Journal Article Research Support, N.I.H., Extramural Randomized Controlled Trial +1 more

**The obesity paradox: Retinopathy, obesity, and circulating risk markers in youth with type 2 diabetes in the TODAY Study.**

2022 Nov United States Journal of diabetes and its complications

Lynne L Levitsky, Kimberly L Drews, Morey Haymond +7 more

type: NA year: 2022

36150365 | PMC12396272 | NCT00081328

---

MERGED ClinicalTrials.gov INTERVENTIONAL

**A Trial Investigating the Efficacy and Safety of Insulin Detemir Versus Insulin NPH in Combination With Metformin and Diet/Exercise in Children and Adolescents With Type 2 Diabetes Insufficiently Controlled on Metformin With or Without Other Oral Antidiabetic Drug(s) With or Without Basal Insulin**

2014 Jun - Jun 2016 Argentina, Brazil, Croatia +21 more Novo Nordisk A/S

status: TERMINATED type: INTERVENTIONAL phase: PHASE3 has result: Yes

NCT02191272 | 30014302

---

PubMed Journal Article Randomized Controlled Trial Clinical Trial, Phase III +1 more

**Safety, Growth, and Development After Dapagliflozin or Saxagliptin in Children With Type 2 Diabetes (T2NOW Follow-Up).**

2025 May 19 United States The Journal of clinical endocrinology and metabolism

Naim Shehadeh, Pietro Galassetti, Nayyar Iqbal +5 more

type: NA year: 2025

39446459 | PMC12086385

---

ClinicalTrials.gov INTERVENTIONAL

**Low Glycemic Load Diets in Latino Children at Risk for Type 2 Diabetes**

2003 Oct - May 2010 United States Children's National Research Institute

status: COMPLETED type: INTERVENTIONAL phase: PHASE2/PHASE3 has result: No

NCT01068197 | 21309658 | 23255569

Figure 6: Search results

Home About

Search Filters

Merged Result View: PubMed ClinicalTrials.gov Clear Selections (0) Select All Download Selected

MERGED ClinicalTrials.gov INTERVENTIONAL

**Surgical or Medical Treatment**

2019 Dec - Nov 2026 United States Children's Hospital Medical Center, Cincinnati

status: RECRUITING type: INTERVENTIONAL phase: PHASE4 has result: No

NCT04128995 | 34389568

---

ClinicalTrials.gov INTERVENTIONAL

**A Study to Evaluate Tirzepatide (LY3298176) in Pediatric and Adolescent Participants With Type 2 Diabetes Mellitus Inadequately Controlled With Metformin or Basal Insulin or Both**

2022 Apr - Jan 2025 Australia, Brazil, France +6 more Eli Lilly and Company

status: COMPLETED type: INTERVENTIONAL phase: PHASE3 has result: No

NCT05260021

---

ClinicalTrials.gov INTERVENTIONAL

**Rosiglitazone and Insulin in T1DM Adolescents**

2003 Aug - Sep 2005 Australia The University of New South Wales

status: COMPLETED type: INTERVENTIONAL phase: PHASE4 has result: No

NCT00372086

---

PubMed Systematic Review Journal Article

**The relationship between glucose and the liver-alpha cell axis - A systematic review.**

2022 Switzerland Frontiers in endocrinology

Thomas Pixner, Nathalie Stummer, Anna Maria Schneider +11 more

type: NA year: 2022

36686477 | PMC9849557

---

PubMed Journal Article Review Research Support, Non-U.S. Gov't

**Do patients with Prader-Willi syndrome have favorable glucose metabolism?**

2022 May 07 England Orphanet journal of rare diseases

Figure 7: Filtering sidebar

Home About

Search Results

Total: 1285 Merged: 2 PubMed-only: 998 CTG-only: 285

Merged Result View: PubMed ClinicalTrials.gov Clear Selections (0) Select All Download Selected

MERGED ClinicalTrials.gov INTERVENTIONAL

**Surgical or Medical Treatment**

Dec 2019 - Nov 2026 United States Children's Hospital Medical Center, Cincinnati  
 status: RECRUITING type: INTERVENTIONAL phase: PHASE4 has result: No  
 NCT04128995 | 34389568 [View](#)

ClinicalTrials.gov INTERVENTIONAL

**A Study to Evaluate Tirzepatide (LY3298176) in Pediatric and Adolescent Participants With Type 2 Diabetes Mellitus Inadequately Controlled With Metformin or Basal Insulin or Both**

Apr 2022 - Jan 2025 Australia, Brazil, France +6 more Eli Lilly and Company  
 status: COMPLETED type: INTERVENTIONAL phase: PHASE3 has result: No  
 NCT05260021 [View](#)

ClinicalTrials.gov INTERVENTIONAL

**Rosiglitazone and Insulin in T1DM Adolescents**

Aug 2003 - Sep 2005 Australia The University of New South Wales  
 status: COMPLETED type: INTERVENTIONAL phase: PHASE4 has result: No  
 NCT00372086 [View](#)

PubMed Systematic Review Journal Article

**The relationship between glucose and the liver-alpha cell axis - A systematic review.**

2022 Switzerland Frontiers in endocrinology  
 Thomas Pixner, Nathalie Stummer, Anna Maria Schneider +11 more  
 type: NA year: 2022  
 36686477 | PMC9849557 [View](#)

ClinicalTrials.gov Preview

NCTID NCT02131272  
 Related PMIDs 30014302

Study Design

Study Type INTERVENTIONAL  
 Phase PHASE3  
 Enrollment 42 (ACTUAL)  
 Conditions Diabetes Diabetes Mellitus, Type 2  
 Interventions Diet/exercise Insulin detemir  
 Insulin NPH  
 Arm Groups EXPERIMENTAL: Insulin detemir and diet/exercise  
 ACTIVE COMPARATOR: Insulin NPH and diet/exercise

Eligibility Criteria

Age Groups CHILD, ADULT  
 Sexes ALL  
 Age Range 10 Years to 18 Years  
 Healthy Volunteers No

Primary Outcomes

- Change in HbA1c (Glycosylated Haemoglobin)

Secondary Outcomes

- Change in Body Weight Standard Deviation Score (SDS)
- Proportion of Subjects Achieving HbA1c Below 7.0%, Who Have Not Experienced Any Treatment Emergent Severe Hypoglycaemic Episodes Within the Last 14 Weeks of Treatment.
- Proportion of Subjects Achieving HbA1c Below 7.5%, Who Have Not Experienced Any Treatment Emergent Severe Hypoglycaemic Episodes Within the Last 14 Weeks of Treatment

Figure 8: Preview sidebar

Home About

Search Filters

Data Source  
 PubMed 4  
 ClinicalTrials.gov 26

Publication Date  
 1 year 1 4  
 5 years 2 11  
 10 years 4 19  
 Custom Range

Phase  
 Phase I 0  
 Phase II 4  
 Phase III 4 26  
 Phase IV 0

Study Type  
 Clinical Trial 4 26  
 Randomized Controlled Trial 4 26  
 Observational Study 0

Age  
 Child: 0-18 years 4 20  
 Adult: 19+ years 1 4  
 Aged: 65+ years 1 1

PubMed Only  
 PMC Open Access

Article Types

PubMed Clinical Trial, Phase III Journal Article Multicenter Study +1 more

**Rationale, design, and methods for the Medical Optimization and Management of Pregnancies with Overt Type 2 Diabetes (MOMPOD) study.**

2018 Dec 12 England BMC pregnancy and childbirth  
 Diane C Berry, Sonia Davis Thomas, Karen F Dorman +4 more  
 type: NA year: 2018  
 30541506 | PMC6292086 | NCT02932475 [View](#)

MERGED ClinicalTrials.gov INTERVENTIONAL

**A Trial Investigating the Efficacy and Safety of Insulin Detemir Versus Insulin NPH in Combination With Metformin and Diet/Exercise in Children and Adolescents With Type 2 Diabetes Insufficiently Controlled on Metformin With or Without Other Oral Antidiabetic Drug(s) With or Without Basal Insulin**

Jun 2014 - Jun 2016 Argentina, Brazil, Croatia +21 more Novo Nordisk A/S  
 status: TERMINATED type: INTERVENTIONAL phase: PHASE3 has result: Yes  
 NCT02131272 | 30014302 [View](#)

PubMed Journal Article Randomized Controlled Trial Clinical Trial, Phase III +1 more

**Safety, Growth, and Development After Dapagliflozin or Saxagliptin in Children With Type 2 Diabetes (T2NOW Follow-Up).**

2025 May 19 United States The Journal of clinical endocrinology and metabolism  
 Naim Shehadeh, Pietro Galassetti, Nayyar Iqbal +5 more  
 type: NA year: 2025  
 39446459 | PMC12086385 [View](#)

ClinicalTrials.gov INTERVENTIONAL

**Study to Evaluate Safety and Efficacy of Dapagliflozin in Patients With Type 2 Diabetes Mellitus Aged 10-24 Years**

Jun 2016 - Apr 2020 Hungary, Israel, Mexico +4 more AstraZeneca  
 status: COMPLETED type: INTERVENTIONAL phase: PHASE3 has result: Yes  
 NCT02725593 | 35378069 [View](#)

ClinicalTrials.gov Preview

Eligibility Assessment  
 Met 2  Not met 1

Inclusion Criteria

Intervention includes ACE inhibitors (e.g., lisinopril, enalapril, ramipril)  Not Met

Truth:  False

Confidence:  90%

Evidence  
 Interventions: DRUG: Dapagliflozin, DRUG: Dapagliflozin placebo

Reasoning  
 The abstract explicitly mentions the interventions as dapagliflozin and its placebo. There is no mention of ACE inhibitors such as lisinopril, enalapril, or ramipril.

Outcomes include HbA1c, fasting glucose, or similar  Met biomarkers

Truth:  True

Confidence:  90%

Evidence  
 Primary Outcomes: Adjusted Change From Baseline in Glycated Haemoglobin (HbA1c) at Week 24

Reasoning  
 The abstract explicitly states that the primary outcome includes HbA1c, which is a biomarker related to diabetes management.

Exclusion Criteria

Non-RCT designs (observational, cohort, etc.)  Not violated

Figure 9: Eligibility check results

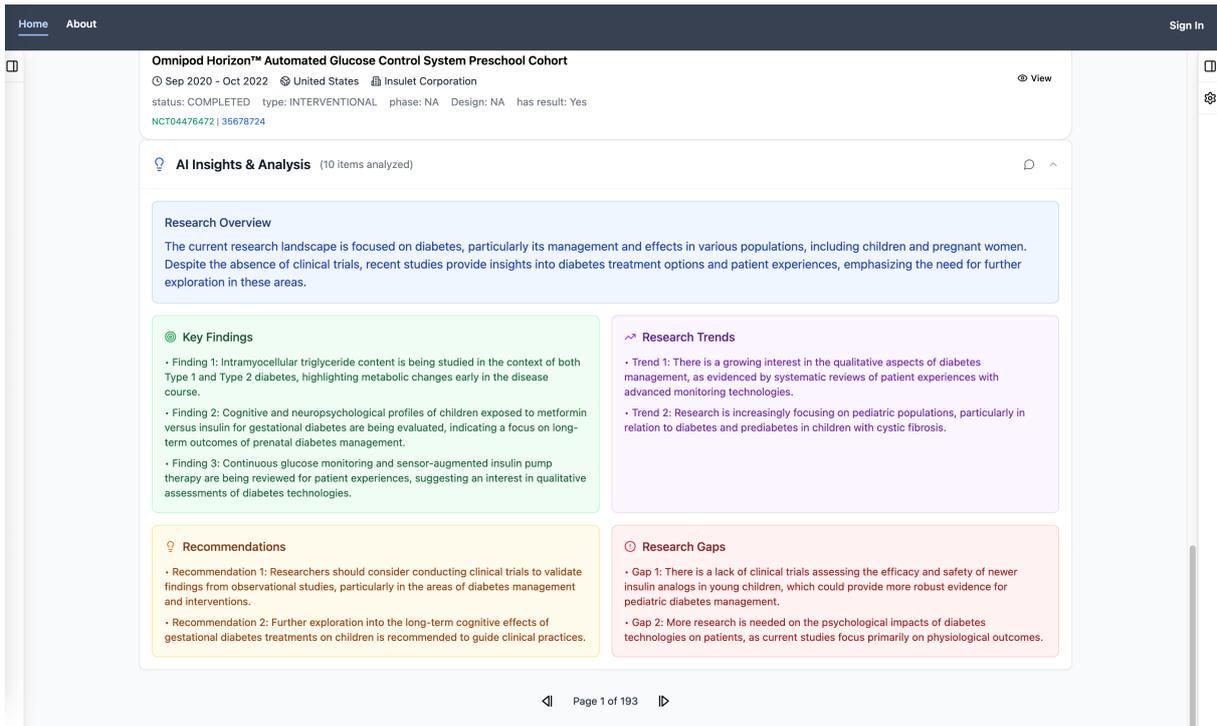


Figure 10: AI Insights

## H User Interface - Detail Page

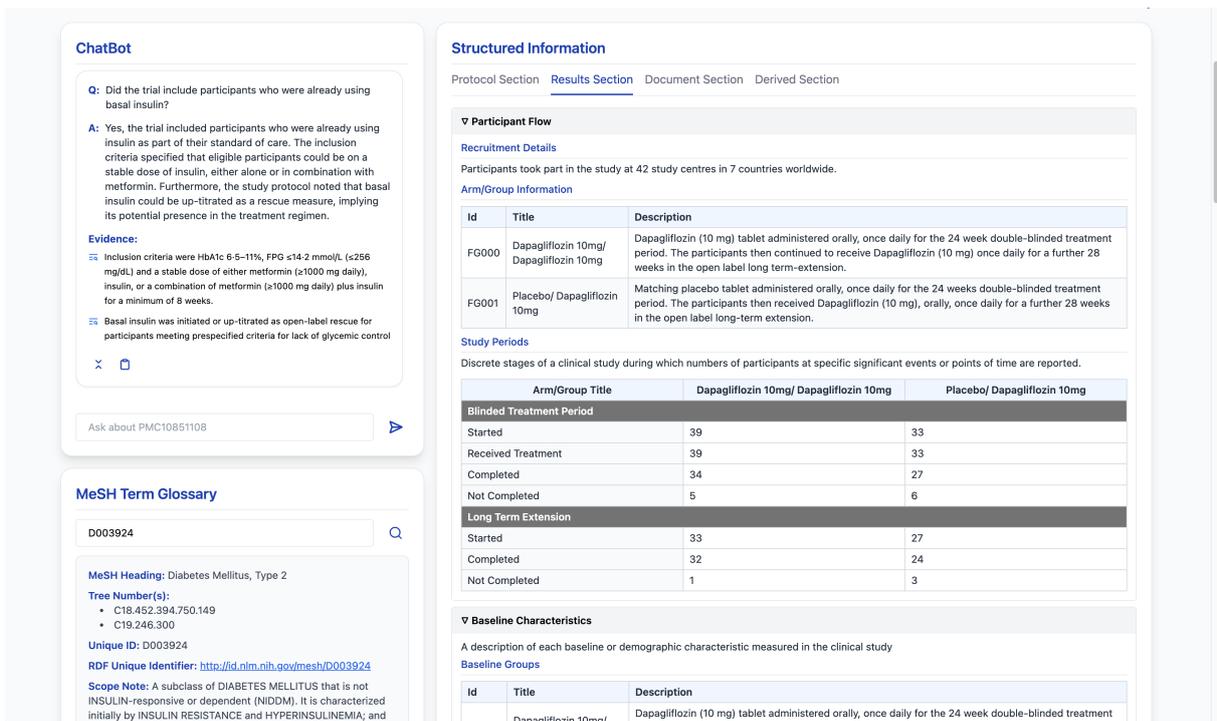


Figure 11: Detail page

## References

### From PubMed

These publications come from PubMed, a public database of scientific and medical articles.

1. Tamborlane WV, Laffel LM, Shehadeh N, Isganaitis E, Van Name M, Ratnayake J, Karlsson C, Norjavaara E. Efficacy and safety of dapagliflozin in children and young adults with type 2 diabetes: a prospective, multicentre, randomised, parallel group, phase 3 study. *Lancet Diabetes Endocrinol.* 2022 May;10(5):341-350. doi: 10.1016/S2213-8587(22)00052-3. Epub 2022 Apr 1.  
PMID: 35378069 | PMCID: PMC10851108  
X Collapse Full Text

### Full Text for Reference 1

#### Randomization and masking

Participants were stratified by sex, age (10–15, >15–<18, ≥18–<25 years) and background medication (metformin alone, insulin alone, or insulin+metformin). *A priori* recruitment of participants aged 18–<25 years was limited to <40% of the total population, while recruitment of participants aged 10–15 years was to comprise ≥20% of the total population. An interactive web/voice response system randomly assigned treatment (placebo or study drug) to each participant. During the 24-week efficacy period, participants and study personnel were blinded to treatment. Treatment during the subsequent 28-week extension period was open-label, although blinding with respect to treatment received in the initial 24-week period was maintained. The sponsor was responsible for randomization and blinding.

#### Procedures

The study treatments were oral, once-daily, dapagliflozin 10 mg or placebo added to standard of care (metformin alone, insulin alone or metformin+insulin). Participants were assessed weekly during the 4-week lead-in period, at baseline, during the double-blind period (Week 1, 2, 4, 8, 10, 12, 16, 20, 24), during the open-label safety extension (Week 28, 32, 36, 40, 46, 52) and 4 weeks post treatment. **Basal insulin was initiated or up-titrated as open-label rescue for participants meeting prespecified criteria for lack of glycemic control:** FPG >13.5 mmol/L (>240 mg/dL) during the double-blind period; FPG >10 mmol/L (>180 mg/dL) or HbA1c >8.0% during the open-label period (Table S2). Rescued participants continued treatment with the study drug and continued in the study. Those receiving insulin (as background or as glycemic rescue) underwent dose adjustments as per investigators' discretion.

#### Outcome measures

##### Efficacy

The primary efficacy outcome was mean change from baseline to 24 weeks in HbA1c with dapagliflozin 10 mg versus placebo. Values after glycemic rescue or permanent discontinuation from study drug were excluded. A prespecified sensitivity analysis was to be performed if >10% of participants in either treatment group had protocol deviations predefined as affecting the primary efficacy results. These predefined protocol deviations are described in Table S3. The primary efficacy endpoint was also described according to subgroups: sex, race (white, non-white), baseline HbA1c (<8%, ≥8%) and background medication (insulin+metformin, metformin only).

Secondary endpoints, in order of hierarchical testing, were mean change from baseline to 24 weeks in FPG, percentage of participants who received glycemic rescue or discontinued study due to lack of glycemic control up to 24 weeks, and percentage of participants with baseline HbA1c ≥7% who achieved HbA1c <7% at 24 weeks.

##### Safety

Safety and tolerability were assessed throughout and included reporting of adverse events (AEs), serious AEs (SAEs), discontinuation due to AEs, hypoglycemia, diabetes ketoacidosis (adjudicated by committee), hepatic laboratory parameters, and vital signs (height, weight, BMI z-score and blood pressure).

Hypoglycemia was defined according to American Diabetes Association (ADA) criteria: severe (required assistance to administer carbohydrate, glucagon or other actions to promote neurological recovery), documented symptoms (typical symptoms and/or plasma glucose [PG] ≤3.9 mmol/L [≤70 mg/dL]), asymptomatic (no symptoms, PG ≤3.9 mmol/L [≤70 mg/dL]), probable symptomatic (typical symptoms without a glucose

Figure 12: QA evidence highlight

## I Search Query Generation Prompt Templates

You are a clinical expert skilled in transforming a user's natural language query into precise search queries for PubMed to identify the most relevant clinical trial papers.

The user's input is provided as key-value pairs in JSON format as follows:

```
{
  "user_query": a free text natural language query,
  "cond": condition/disease terms,
  "intr": intervention/treatment terms,
  "other_term": other terms
}
```

Return your result in JSON format exactly as follows:

```
{
  "cond": refined condition/disease terms,
  "intr": refined intervention/treatment terms,
  "other_term": refined other terms,
  "combined_query": the final combined search query created by concatenating the refined 'cond', 'intr', and 'other_term' using the 'AND' operator with parentheses for non-empty fields.
}
```

Please refine the user's input by following these rules:

- If the user's input is not in English, first translate it into English.
- Retain only the meaningful keywords in each field.
- The "user\_query" field is the primary unstructured input where the user might write, for example, "Find clinical trials for depressive disorder using medication," etc. Extract only the meaningful keywords from "user\_query" and assign them to "cond", "intr", or "other\_term" as appropriate. In other words, merge the information in "user\_query" with any provided in "cond", "intr", and "other\_term" so that the final refined values capture all the important concepts.
- If any keywords in "user\_query" or "other\_term" can be classified as condition or intervention terms, move them to "cond" or "intr" accordingly and remove them from "other\_term."
- For each field, if there are multiple keywords, combine them using appropriate Boolean operators (AND, OR) and group with parentheses.

Example:

- User's input:

```
{
  "user_query": "Find clinical trials for depressive disorder using medication; also interested in ADHD in children, including studies on methylphenidate",
  "cond": "depressive disorder",
  "intr": "medication",
  "other_term": "atomoxetine"
}
```

- Your output should be:

```
{
  "cond": "depressive disorder OR adhd",
  "intr": "Methylphenidate OR Atomoxetine",
  "other_term": "child",
  "combined_query": "(depressive disorder OR adhd) AND (Methylphenidate OR Atomoxetine) AND (child)"
}
```

Note: The final search query (the value of "combined\_query") will be used for the database search, and the content of "user\_query" itself will not be used directly.

Figure 13: Search Query Generation - System Prompt

The user's input for clinical trial search:

```
{{inputData}}
```

Return your result in JSON format as follows:

```
{
  "cond": refined condition/disease terms,
  "intr": refined intervention/treatment terms,
  "other_term": refined other terms,
  "combined_query": final combined search query created by concatenating the refined 'cond', 'intr', and 'other_term'
  using the 'AND' operator with parentheses for non-empty fields.
}
```

Figure 14: Search Query Generation - User Prompt

## J Information Extraction Prompt Templates

```
# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}
# TARGET #
Extract the identificationModule of the study. It includes identifying information such as the trial's unique identifier or title, illustrating "who is conducting which trial and where it is registered."
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{{
  "identificationModule": {
    "nctId": \\ ClinicalTrials.gov Identifier. The format is "NCT" followed by an 8-digit number: TEXT (max 11 chars)
    "orgStudyIdInfo": {
      "id": \\ organization's unique protocol ID: TEXT (max 30 chars)
      "type": \\ Type of organization-issued ID: ENUM (NIH, FDA, VA, CDC, AHRQ, SAMHSA)
      "link": \\ URL link related to OrgStudyId and OrgStudyIdType: TEXT
    },
    "secondaryIdInfos": [ \\ ARRAY of OBJECT
      {
        "id": \\ Secondary identifier for funding or registry: TEXT (max 30 chars)
        "type": \\ Type of secondary ID: ENUM (NIH, FDA, VA, CDC, AHRQ, SAMHSA, OTHER_GRANT, EUDRACT_NUMBER, CTIS, REGISTRY, OTHER)
        "domain": \\ Name of funding organization, registry, or issuer: TEXT (max 119 chars)
        "link": \\ URL link related to SecondaryId and SecondaryIdType: TEXT
      }
    ],
    "organization": {
      "fullName": \\ Name of the sponsoring organization: TEXT
      "class": \\ Organization type: ENUM (NIH, FED, OTHER_GOV, INDIV, INDUSTRY, NETWORK, AMBIG, OTHER, UNKNOWN)
    },
    "briefTitle": \\ Short title of the study: TEXT (max 300 chars)
    "officialTitle": \\ Full official title of the study: TEXT (max 600 chars)
    "acronym": \\ Study acronym: TEXT (max 14 chars)
  }
}}
```
```

Figure 15: Protocol Section - Identification Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the descriptionModule and conditionsModule of the study. DescriptionModule offers a brief introduction or summary of the clinical trial. ConditionsModule specifies the target conditions or topics (keywords), indicating which are being studied.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "descriptionModule": {
    "briefSummary": "\\ Concise summary of the study, including its hypothesis, in layman's terms: TEXT (max 5000 chars)
    "detailedDescription": "\\ Extended study description with technical details, excluding full protocol or duplicate information: TEXT (max 32000 chars)
  },
  "conditionsModule": {
    "conditions": "\\ List of disease(s) or condition(s) studied, preferably using MeSH or SNOMED CT terms: ARRAY of TEXT
    "keywords": "\\ List of descriptive words or phrases related to the study, preferably using MeSH terms: ARRAY of TEXT
  }
}
```

```

Figure 16: Protocol Section - Description Module, Conditions Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the designModule of the study. It defines the overall study design (study type, phase, allocation, masking, number of participants, etc.) in detail.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "designModule": {
    "studyType": "\\ Study classification: ENUM (EXPANDED_ACCESS, INTERVENTIONAL, OBSERVATIONAL)
    "patientRegistry": "\\ Indicates if the study is a patient registry: BOOLEAN
    "targetDuration": "\\ Follow-up duration for observational patient registry studies: TIME
    "phases": "\\ Study phase(s), applicable for drug/biologic trials: ARRAY of ENUM (NA, EARLY_PHASE1, PHASE1, PHASE2, PHASE3, PHASE4)
    "designInfo": {
      "allocation": "\\ Method of assigning participants: ENUM (RANDOMIZED, NON_RANDOMIZED, NA)
      "interventionModel": "\\ Type of intervention design: ENUM (SINGLE_GROUP, PARALLEL, CROSSOVER, FACTORIAL, SEQUENTIAL)
      "interventionModelDescription": "\\ Description of the intervention model: TEXT
      "primaryPurpose": "\\ The main objective of the intervention(s) being evaluated by the clinical trial: ENUM (TREATMENT, PREVENTION, DIAGNOSTIC, EC T, SUPPORTIVE_CARE, SCREENING, HEALTH_SERVICES_RESEARCH, BASIC_SCIENCE, DEVICE_FEASIBILITY, OTHER)
      "observationalModel": "\\ Study model for observational studies: ENUM (COHORT, CASE_CONTROL, CASE_ONLY, CASE_CROSSOVER, ECOLOGIC_OR_COMMUNITY, FAMILY_BASED, DEFINED_POPULATION, NATURAL_HISTORY, OTHER)
      "timePerspective": "\\ Time perspective for observational studies: ENUM (RETROSPECTIVE, PROSPECTIVE, CROSS_SECTIONAL, OTHER)
      "maskingInfo": {
        "masking": "\\ Level of blinding: ENUM (NONE, SINGLE, DOUBLE, TRIPLE, QUADRUPLE)
        "maskingDescription": "\\ Detailed description of masking: TEXT (max 1000 chars)
        "whoMasked": "\\ Groups involved in masking: ARRAY of ENUM (PARTICIPANT, CARE_PROVIDER, INVESTIGATOR, OUTCOMES_ASSESSOR)
      }
    }
  },
  "enrollmentInfo": {
    "count": "\\ Number of participants enrolled: NUMERIC
    "type": "\\ Actual or estimated enrollment: ENUM (ACTUAL, ESTIMATED)
  }
}
```

```

Figure 17: Protocol Section - Design Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the armsInterventionsModule of the study. It details which drugs (or procedures) are administered, to which arms (groups), and how they are applied.
# RESPONSE #
A syntactically correct JSON string:
Format:
```json
{
  "armsInterventionsModule": {
    "armGroups": [ \ \ ARRAY of OBJECT
      {
        "label": \ \ Name of the arm/group: TEXT
        "type": \ \ Type of arm: ENUM (EXPERIMENTAL, ACTIVE_COMPARATOR, PLACEBO_COMPARATOR, SHAM_COMPARATOR, NO_INTERVENTION, OTHER)
        "description": \ \ Description of the arm/group: TEXT
        "interventionNames": \ \ List of interventions used in this arm/group: ARRAY of TEXT
      }
    ],
    "interventions": [ \ \ ARRAY of OBJECT
      {
        "type": \ \ Type of intervention: ENUM (BEHAVIORAL, BIOLOGICAL, COMBINATION_PRODUCT, DEVICE, DIAGNOSTIC_TEST, DIETARY_SUPPLEMENT, DRUG, GENETIC, PROCEDURE, RADIATION, OTHER)
        "name": \ \ Name of the intervention: TEXT
        "description": \ \ Description of the intervention: TEXT
        "armGroupLabels": \ \ List of arm/group labels associated with this intervention: ARRAY of TEXT
      }
    ]
  }
}
```

```

Figure 18: Protocol Section - Arms Interventions Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the outcomesModule of the study. It describes the primary, secondary, and other outcome measures, showing which indicators are used to assess the trial's effectiveness and safety.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type. Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
"outcomesModule": {
  "primaryOutcomes": [ \ \ Required. List of primary outcome measures used to assess the trial's main objectives: ARRAY of OBJECT
    {
      "measure": \ \ Name of the primary outcome measure: TEXT (max 255 chars)
      "description": \ \ Description of the metric used to characterize the primary outcome measure: TEXT (max 999 chars)
      "timeFrame": \ \ Time point(s) at which the outcome is measured: TEXT (max 255 chars)
    }
  ],
  "secondaryOutcomes": [ \ \ Conditional. List of secondary outcome measures for additional study assessments: ARRAY of OBJECT
    {
      "measure": \ \ Name of the secondary outcome measure: TEXT
      "description": \ \ Description of the metric used to characterize the secondary outcome measure: TEXT
      "timeFrame": \ \ Time point(s) at which the outcome is measured: TEXT
    }
  ],
  "otherOutcomes": [ \ \ Optional. List of other pre-specified outcome measures (excluding post-hoc measures): ARRAY of OBJECT
    {
      "measure": \ \ Name of the pre-specified outcome measure: TEXT
      "description": \ \ Description of the metric used to characterize the outcome measure: TEXT
      "timeFrame": \ \ Time point(s) at which the outcome is measured: TEXT
    }
  ]
}
```

```

Figure 19: Protocol Section - Outcomes Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the eligibilityModule of the study. It specifies the eligibility criteria for participating in this clinical trial.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "eligibilityModule": {
    "eligibilityCriteria": "\\ Inclusion and exclusion criteria for participant selection, formatted as a bulleted list under respective headers: TEXT (max 20000 chars)
    "healthyVolunteers": "\\ Indicates if healthy volunteers without the studied condition can participate: BOOLEAN
    "sex": "\\ Eligible participant sex: ENUM (FEMALE, MALE, ALL)
    "minimumAge": "\\ Minimum age required for participation, with unit of time: TEXT (Years, Months, Weeks, Days, Hours, Minutes, N/A)
    "maximumAge": "\\ Maximum age allowed for participation, with unit of time: TEXT (Years, Months, Weeks, Days, Hours, Minutes, N/A)
    "stdAges": "\\ Standardized age categories: ARRAY of ENUM (CHILD, ADULT, OLDER_ADULT)
    "studyPopulation": "\\ (Observational studies only) Description of the population source for cohorts or groups: TEXT (max 1000 chars)
    "samplingMethod": "\\ (Observational studies only) Method used for sampling: ENUM (PROBABILITY_SAMPLE, NON_PROBABILITY_SAMPLE)
  }
}
...

```

Figure 20: Protocol Section - Eligibility Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the participantFlowModule of the study. It describes the flow of participants through each stage of the study, including enrollment, allocation, follow-up, and analysis.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "participantFlowModule": {
    "preAssignmentDetails": "\\ Description of significant events after participant enrollment but before group assignment: TEXT (max 500 chars)
    "recruitmentDetails": "\\ Key information about recruitment process: TEXT (max 500 chars)
    "typeUnitsAnalyzed": "\\ Unit of analysis (e.g., "Participants", "Eyes", "Lesions"): TEXT
    "groups": [ \\ ARRAY of OBJECT - Arms/groups in the study flow
      {
        "id": "\\ Unique group identifier. FG000 is the first group, FG001 is the second, and so on: TEXT
        "title": "\\ Short name of the arm/group: TEXT (max 40 chars)
        "description": "\\ Brief description of the arm/group: TEXT (max 1500 chars)
      }
    ],
    "periods": [ \\ ARRAY of OBJECT - Time periods in the study
      {
        "title": "\\ Period name (e.g., "Overall Study", "Treatment Phase"): TEXT (max 40 chars)
        "milestones": [ \\ ARRAY of OBJECT - Key milestones
          {
            "type": "\\ Milestone name (e.g., "STARTED", "COMPLETED"): TEXT (max 100 chars)
            "comment": "\\ Additional information about the milestone: TEXT (max 500 chars)
            "achievements": [ \\ ARRAY of OBJECT - Numbers for each group
              {
                "groupId": "\\ References a group ID: TEXT (max 500 chars)
                "comment": "\\ Explanation if number differs from expected: TEXT (max 500 chars)
                "numSubjects": "\\ Number of participants: TEXT
                "numUnits": "\\ Number of units if different from participants: TEXT
              }
            ]
          }
        ]
      }
    ],
    "dropWithdraws": [ \\ ARRAY of OBJECT - Reasons for not completing
      {
        "type": "\\ Reason category (e.g., "Adverse Event", "Lost to Follow-up"): TEXT (max 100 chars)
        "comment": "\\ Additional details about the reason: TEXT
        "reasons": [ \\ ARRAY of OBJECT - Numbers for each group
          {
            "groupId": "\\ References a group ID: TEXT
            "comment": "\\ Additional explanation if needed: TEXT
            "numSubjects": "\\ Number of participants: TEXT
          }
        ]
      }
    ]
  }
}
...

```

Figure 21: Results Section - Participant Flow Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding
data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the baselineCharacteristicsModule of the study.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "baselineCharacteristicsModule": { // Baseline demographic and other initial measures, by arm/group.
    "populationDescription": // Brief reason or explanation if baseline participants differ from the assigned groups.
    "typeUnitsAnalyzed": // (Optional) If units are not participants (e.g., eyes, lesions).
    "groups": [ // ARRAY of OBJECT. Arms/groups for baseline assessment. Must include a "Total" group as the last entry.
      {
        "id": // BG000 is the first group, BG001 is the second, and so on.
        "title": // Short label that identifies the group (e.g., "Placebo", "Treatment A", "Total").
        "description": // Brief explanation of the group's characteristics or interventions.
      }
    ],
    "denoms": [ // ARRAY of OBJECT. Structure for Overall Baseline Measure Data (Row).
      {
        "units": // Unit of measure for the data in this row. Default is "Participants".
        "counts": [ // ARRAY of OBJECT. Each object represents a group and its corresponding count in the same order as the "groups" array.
          {
            "groupId": // References an ID from the "groups" array. (e.g., "BG000", "BG001", "BG002").
            "value": // Number of participants in this group.
          }
        ]
      }
    ],
    "measures": [ // ARRAY of OBJECT. Each baseline or demographic characteristic. Required baseline measures include Age, Sex/Gender, Race, Ethnicity (if
    applicable), and any other measures.
      {
        "title": // ENUM - See required/optional values below
        // Required1: Age => ENUM("Age, Continuous", "Age, Categorical", "Age, Customized")
        // Required2: Sex/Gender => ENUM("Sex: Female, Male", "Sex/Gender, Customized")
        // Required3 (if possible): Race and Ethnicity => ENUM("Race (NIH/OMB)", "Ethnicity (NIH/OMB)", "Race/Ethnicity, Customized", "Race and Ethni
        city Not Collected")
        // Required4 (if possible): Region of Enrollment => ENUM("Region of Enrollment")
        // (Optional): Any other measures
        "description": // Additional descriptive information about the baseline measure
        "populationDescription": // (Optional) If the analyzed population differs from the overall baseline population.
        "paramType": // The type of data for the baseline measure. ENUM("COUNT_OF_PARTICIPANTS", "MEAN", "NUMBER", "MEDIAN", "COUNT_OF_UNITS", "GEOMETRIC_
        MEAN", "LEAST_SQUARES_MEAN", "LOG_MEAN", "GEOMETRIC_LEAST_SQUARES_MEAN")
        "dispersionType": // Baseline Measure Dispersion/Precision. ENUM("STANDARD_DEVIATION", "FULL_RANGE", "INTER_QUARTILE_RANGE", "NA", "CONFIDENCE_B
        0", "CONFIDENCE_90", "CONFIDENCE_95", "CONFIDENCE_975", "CONFIDENCE_99", "CONFIDENCE_OTHER", "GEOMETRIC_COEFFICIENT", "STANDARD_ERROR")
        "unitOfMeasure": // e.g., "Participants", "years", "kg", etc.
        "denoms": [ // ARRAY of OBJECT. Same structure as "denoms" above, if needed for measure-specific denominators.
          {
            "units": // TEXT
            "counts": [ // ARRAY of OBJECT.
              {
                "groupId": // TEXT
                "value": // TEXT
              }
            ]
          }
        ]
      }
    ],
    "classes": [ // ARRAY of OBJECT. Within each measure, define rows or classifications.
      {
        "title": // Baseline RowTitle. (e.g., "Sex: Female, Male")
        "denoms": [ // ARRAY of OBJECT. Same structure as "denoms" above, if needed for class-specific counts.
          {
            "units": // TEXT
            "counts": [ // ARRAY of OBJECT.
              {
                "groupId": // TEXT
                "value": // TEXT
              }
            ]
          }
        ]
      }
    ],
    "categories": [ // ARRAY of OBJECT. Each category is essentially a sub-row under the class.
      {
        "title": // e.g., "Female", "Hispanic or Latino"
        "measurements": [ // ARRAY of OBJECT. Data for each group in this category.
          {
            "groupId": // e.g., "BG000"
            "value": // e.g., "53", "64.4"
            "spread": // e.g., "9.4" (for SD), or omitted if not applicable.
            "lowerLimit": // e.g., "5.6" Based on Measure Type and Measure of Dispersion (e.g., lower Limit of Full Range).
            "upperLimit": // e.g., "9.9" Based on Measure Type and Measure of Dispersion (e.g., upper Limit of Full Range)
          }
        ]
      }
    ]
  }
}
...

```

Figure 22: Results Section - Baseline Characteristics Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the outcomeMeasuresModule of the study. It contains the results of primary, secondary, and other outcome measures.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "outcomeMeasuresModule": {
    "outcomeMeasures": [ // ARRAY of OBJECT - All outcome measures with results
      {
        "type": // Outcome type: ENUM ("PRIMARY", "SECONDARY", "OTHER_PRE_SPECIFIED", "POST_HOC")
        "title": // Outcome measure title: TEXT (max 255 chars)
        "description": // Detailed description: TEXT (max 999 chars)
        "populationDescription": // Analysis population if different: TEXT
        "reportingStatus": // Whether data is reported: ENUM ("POSTED", "NOT_POSTED")
        "anticipatedPostingDate": // Expected date if not posted: DATE (format: YYYY or YYYY-MM)
        "paramType": // Type of measure: ENUM ("GEOMETRIC_MEAN", "GEOMETRIC_LEAST_SQUARES_MEAN", "LEAST_SQUARES_MEAN", "LOG_MEAN", "MEAN", "MEDIAN", "NUMBER", "COUNT_OF_PARTICIPANTS", "COUNT_OF_UNITS")
        "dispersionType": // Dispersion/precision type: ENUM ("Not Applicable", "Standard Deviation", "Standard Error", "Inter-Quartile Range", "Full Range", "99% Confidence Interval", "97.5% Confidence Interval", "95% Confidence Interval", "90% Confidence Interval", "80% Confidence Interval", "Other Confidence Interval Level", "Geometric Coefficient of Variation")
        "unitOfMeasure": // Unit of measurement: TEXT
        "calculatePct": // Whether to calculate percentage: BOOLEAN
        "timeFrame": // The description of the time point(s) of assessment must be specific to the outcome measure and is generally the specific duration of time over which each participant is assessed (not the overall duration of the study): TEXT (max 255 chars)
        "typeUnitsAnalyzed": // (Optional) If units are not participants (e.g., eyes, lesions): TEXT (max 40 chars)
        "denomUnitsSelected": // Selected denominator units (e.g., Participants, eyes, lesions): TEXT
        "groups": [ // ARRAY of OBJECT - Study arms/groups
          {
            "id": //Unique group identifier. 0G000 is the first group, 0G001 is the second, and so on: TEXT
            "title": // Group name: TEXT
            "description": // Group description: TEXT
          }
        ],
        "denoms": [ // ARRAY of OBJECT - Denominators
          {
            "units": // Unit type: TEXT
            "counts": [ // ARRAY of OBJECT
              {
                "groupId": // References a group ID: TEXT
                "value": // Count value: TEXT
              }
            ]
          }
        ],
        "classes": [ // ARRAY of OBJECT - Outcome categories/timepoints
          {
            "title": // Category/timepoint name: TEXT
            "denoms": [ // ARRAY of OBJECT. Similar structure as above.
              {
                "units": // TEXT
                "counts": [
                  {
                    "groupId": // References a group ID: TEXT
                    "value": // TEXT
                  }
                ]
              }
            ]
          }
        ],
        "categories": [ // ARRAY of OBJECT - Subcategories
          {
            "title": // Subcategory name: TEXT
            "measurements": [ // ARRAY of OBJECT - Results for each group
              {
                "groupId": // References a group ID: TEXT
                "value": // Result value: TEXT
                "spread": // Spread (e.g., SD, SE): TEXT
                "lowerLimit": // CI lower limit: TEXT
                "upperLimit": // CI upper limit: TEXT
                "comment": // Explanation for NA values: TEXT
              }
            ]
          }
        ]
      }
    ]
  }
}
```

```

Figure 23: Results Section - Outcomes Measures Module 1



```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding
data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the adverseEventsModule of the study. It contains information about adverse events including serious adverse events, other adverse events, and mortal
ity data.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "adverseEventsModule": {
    "frequencyThreshold": // Threshold for reporting other AEs (e.g., "5%"): TEXT
    "timeFrame": // Time period for AE collection: TEXT (max 500 chars)
    "description": // Additional AE collection details: TEXT
    "allCauseMortalityComment": // Explanation about mortality data: TEXT
    "eventGroups": [ // ARRAY of OBJECT - Study arms/groups for AE reporting
      {
        "id": // Unique group identifier. EG000 is the first group, EG001 is the second, and so on: TEXT
        "title": // Group name: TEXT (max 1500 chars)
        "description": // Group description: TEXT
        "deathsNumAffected": // Number affected by all-cause mortality: TEXT
        "deathsNumAtRisk": // Number at risk for mortality: TEXT
        "seriousNumAffected": // Number with any serious AE: TEXT
        "seriousNumAtRisk": // Number at risk for serious AEs: TEXT
        "otherNumAffected": // Number with any other AE: TEXT
        "otherNumAtRisk": // Number at risk for other AEs: TEXT
      }
    ],
    "seriousEvents": [ // ARRAY of OBJECT - Serious adverse events by organ system
      {
        "term": // AE preferred term: TEXT
        "organSystem": // Organ system category: TEXT
        "sourceVocabulary": // Coding dictionary (e.g., "MedDRA 23.0"): TEXT
        "assessmentType": // Collection method: ENUM (NON_SYSTEMATIC_ASSESSMENT, SYSTEMATIC_ASSESSMENT)
        "notes": // Additional description: TEXT
        "stats": [ // ARRAY of OBJECT - Statistics for each group
          {
            "groupId": // References an event group ID: TEXT
            "numEvents": // Total number of events: TEXT
            "numAffected": // Number of participants affected: TEXT
            "numAtRisk": // Number at risk: TEXT
          }
        ]
      }
    ],
    "otherEvents": [ // ARRAY of OBJECT - Other (non-serious) adverse events
      {
        "term": // AE preferred term: TEXT
        "organSystem": // Organ system category: TEXT
        "sourceVocabulary": // Coding dictionary: TEXT
        "assessmentType": // Collection method: ENUM (NON_SYSTEMATIC_ASSESSMENT, SYSTEMATIC_ASSESSMENT)
        "notes": // Additional description: TEXT
        "stats": [ // ARRAY of OBJECT - Statistics for each group
          {
            "groupId": // References an event group ID: TEXT
            "numEvents": // Total number of events: TEXT
            "numAffected": // Number of participants affected: TEXT
            "numAtRisk": // Number at risk: TEXT
          }
        ]
      }
    ]
  }
}
```

```

Figure 25: Results Section - Adverse Events Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding
data does not exist in the target articles. Do not invent or assume information that is not present in the given content.
# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the moreInfoModule of the study. It contains additional information including limitations and caveats, certain agreements, and point of contact.
# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "moreInfoModule": {
    "limitationsAndCaveats": {
      "description": // Study limitations and caveats discussed in the paper: TEXT (max 500 chars)
    },
    "certainAgreement": {
      "piSponsorEmployee": // Whether PIs are sponsor employees: BOOLEAN
      "restrictionType": // Type of disclosure restriction: ENUM ("LTE60", "GT60", "OTHER")
      "restrictiveAgreement": // Whether restrictive agreements exist: BOOLEAN
      "otherDetails": // Details if restriction type is OTHER: TEXT
    },
    "pointOfContact": {
      "title": // Contact person's name or title: TEXT (max 255 chars)
      "organization": // Contact organization: TEXT (max 255 chars)
      "email": // Contact email: TEXT (max 255 chars)
      "phone": // Contact phone: TEXT (max 30 chars)
      "phoneExt": // Phone extension: TEXT
    }
  }
}
```

```

Figure 26: Results Section - More Info Module

```

# CONTEXT #
You are tasked with analyzing clinical trial study reports or papers to extract specific information as structured data. Omit any fields if the corresponding
data does not exist in the target articles. Do not invent or assume information that is not present in the given content.

# PAPER CONTENT #
{{pmc_text}}

# TARGET #
Extract the following modules from the study:
1. conditionBrowseModule: MeSH condition term mappings
2. interventionBrowseModule: MeSH intervention term mappings

# RESPONSE #
The JSON response must adhere to the following structure, where each field is annotated with its expected type.
Ensure that the output is a valid and correctly formatted JSON string.
Format:
```json
{
  "conditionBrowseModule": {
    "meshes": [ \\ Condition MeSH Terms MeSH terms of Condition/Diseases field
      {
        "id": \\ Condition MeSH ID MeSH ID: TEXT,
        "term": \\ Condition MeSH Term MeSH Heading: TEXT
      }
    ],
    "ancestors": [ \\ Ancestors of Condition MeSH Terms Ancestor (higher level and more broad) terms of Condition MeSH terms in MeSH Tree hierarchy
      {
        "id": \\ Condition Ancestor MeSH ID MeSH ID: TEXT,
        "term": \\ Condition Ancestor MeSH Term MeSH Heading: TEXT
      }
    ],
    "browseLeaves": [ \\ Condition Leaf Topics Leaf browsing topics for Condition field
      {
        "id": \\ Condition Leaf Topic ID: TEXT,
        "name": \\ Condition Leaf Topic Name: TEXT,
        "relevance": \\ Relevance to Condition Leaf Topic: ENUM (LOW, HIGH)
      }
    ],
    "browseBranches": [ \\ Condition Branch Topics Branch browsing topics for Condition field
      {
        "abbrev": \\ Condition Branch Topic Short Name: TEXT,
        "name": \\ Condition Branch Topic Name: TEXT
      }
    ]
  },
  "interventionBrowseModule": {
    "meshes": [ \\ Condition MeSH Terms MeSH terms of Condition/Diseases field
      {
        "id": \\ Condition MeSH ID MeSH ID: TEXT,
        "term": \\ Condition MeSH Term MeSH Heading: TEXT
      }
    ],
    "ancestors": [ \\ Ancestors of Condition MeSH Terms Ancestor (higher level and more broad) terms of Condition MeSH terms in MeSH Tree hierarchy
      {
        "id": \\ Condition Ancestor MeSH ID MeSH ID: TEXT,
        "term": \\ Condition Ancestor MeSH Term MeSH Heading: TEXT
      }
    ],
    "browseLeaves": [ \\ Condition Leaf Topics Leaf browsing topics for Condition field
      {
        "id": \\ Condition Leaf Topic ID: TEXT,
        "name": \\ Condition Leaf Topic Name: TEXT,
        "relevance": \\ Relevance to Condition Leaf Topic: ENUM (LOW, HIGH)
      }
    ],
    "browseBranches": [ \\ Condition Branch Topics Branch browsing topics for Condition field
      {
        "abbrev": \\ Condition Branch Topic Short Name: TEXT,
        "name": \\ Condition Branch Topic Name: TEXT
      }
    ]
  }
}
```

```

Figure 27: Results Section - Condition Browse Module, Intervention Browse Module

# SciTrue: Evidence-Grounded Claim Verification in Science

Neşet Özkan TAN

Minghao Li

Mark Gahegan

Department of Computer Science

University of Auckland

Auckland, New Zealand

naset.tan@auckland.ac.nz, minghao.lee2017@outlook.com,

m.gahegan@auckland.ac.nz

## Abstract

Large language models (LLMs) have expanded the potential for AI-assisted scientific claim verification, yet existing systems often exhibit unverifiable attributions, shallow evidence mapping, and hallucinated citations. We present SciTrue, a claim verification system providing source-level accountability and evidence traceability. SciTrue links each claim component to explicit, verifiable scientific sources, enabling users to inspect and challenge model inferences, addressing limitations of both general-purpose and search-augmented LLMs. In a human evaluation of 300 attributions, SciTrue achieves high fidelity in summary traceability, attribution accuracy, and context alignment, substantially outperforming RAG-based baselines such as GPT-4o-search-preview and Perplexity Sonar Pro. These results underscore the importance of principled attribution and context-aware reasoning in AI-assisted scientific verification. A system demo is available at [www.scitruer.org](http://www.scitruer.org).

## 1 Introduction

The increasing adoption of large language models (LLMs) such as GPT-4 (OpenAI, 2023), Gemini 2.5 (Google DeepMind, 2024), and Llama-3 (Meta, 2025) has transformed the landscape of information access, reading comprehension, and scientific literature summarization. Despite their impressive capabilities, LLMs exhibit a significant tendency to generate content that is inconsistent with established knowledge or the provided input context, a phenomenon widely referred to as hallucination (Mittelstadt et al., 2023). This issue manifests in various ways within the scientific domain, including the fabrication of scientific references and the generation of seemingly accurate citations that, upon closer inspection, do not support the claims being made (Tilwani et al., 2024). This is particularly problematic in high-stakes domains such as biomedicine or policy, where unverified or misat-

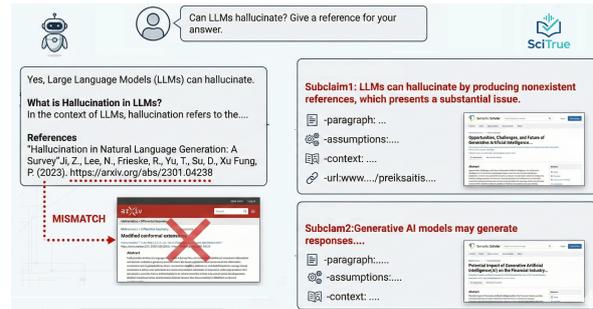


Figure 1: Scientific claim verification with web-access LLMs versus SciTrue on real-world examples. Web-access LLMs often respond with overconfident answers and cite incorrect references (e.g., linking to a differential geometry paper for a claim about LLMs). In contrast, SciTrue provides end-to-end scientific traceability for each attribution. To view the full hallucination example, see: <https://chatgpt.com/share/6867bbc2-73e0-8002-9bda-4eb1223041b0>

tributed evidence can propagate misinformation or erode trust (Chen et al., 2025).

Conventional “research assistant” models, such as OpenAI’s Deep Research<sup>1</sup> and Gemini Deep Research<sup>2</sup>, integrate access to the web and scholarly databases to enhance factual grounding. Nonetheless, use cases reveal persistent deficiencies. These systems often take more than five minutes to process a single claim, ranging from 5 to 30 minutes according to the official website<sup>1</sup>, when reviewing multiple documents, making them impractical for time-sensitive applications such as medicine or journalism. Their outputs tend to be excessively long, which hinders the extraction of concise, actionable insights. Attribution is often shallow: although paragraphs are linked to sources, the synthesis does not clearly indicate which information comes from which source. Moreover, many of these systems follow rigid, predefined formats, such as systematic reviews, that limit their

<sup>1</sup><https://openai.com/index/introducing-deep-research/>

<sup>2</sup><https://gemini.google/overview/deep-research/?hl=en>

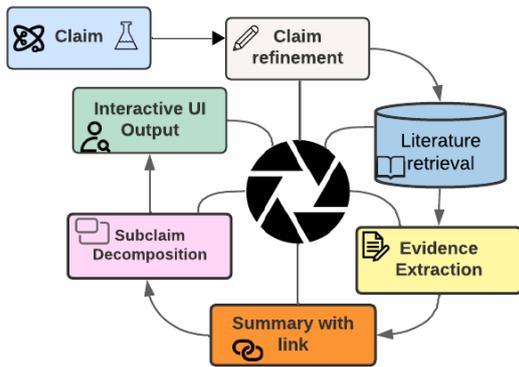


Figure 2: Workflow of the SciTrue agentic pipeline. A user submits a scientific claim, which is then refined. Relevant articles are retrieved, and supporting evidence is extracted. The system generates a summary, decomposes it into subclaims, and presents the results in an interactive interface.

adaptability to diverse, real-world claim verification tasks.

In science, transparency and accountability are foundational principles (Binz et al., 2025). Accurately interpreting and verifying scientific claims necessitates a thorough understanding of the contextual factors and underlying assumptions that frame these claims (Kanavouras and Coutelieris, 2020). Scientific findings are inherently contingent upon specific experimental conditions, population characteristics, and methodological choices (Bottesini et al., 2023). Neglecting such context can lead to overgeneralization or misapplication of results, thereby compromising the validity and utility of the claim. Failure to properly attribute claims can lead to serious issues, including plagiarism, the spread of misinformation, and an erosion of user trust.

To address these challenges, we present **SciTrue**, a scientific claim verification system specifically designed for domains where auditability is critical, such as scientific research. SciTrue provides end-to-end traceability by explicitly attributing verifiable sources, along with the associated context, assumptions, and credibility. Human evaluations demonstrate that SciTrue’s attributions significantly outperform those of leading systems in terms of summary traceability, attribution accuracy, contextual and assumption alignment, and scholarly credibility.

<sup>1</sup><https://openai.com/index/introducing-deep-research/>

## 2 Related Work

The automation of claim verification has received increasing attention due to the rise of misinformation and the need for auditability in scholarly communication (Wadden et al., 2020; Augenstein et al., 2019; Thorne et al., 2018). Benchmarks such as SciFact (Wadden et al., 2020), MultiFC (Augenstein et al., 2019), and FEVER (Thorne et al., 2018) have advanced the development of systems for evidence retrieval and claim stance classification, typically pairing passage retrieval with textual entailment models. However, most approaches focus on closed-domain settings and do not provide sufficient granularity, such as context or assumptions, in associating specific components of a claim with the underlying supporting text.

In the domain of document summarization, a significant line of research has explored multi-document summarization in the medical field (DeYoung et al., 2021; Wallace et al., 2021). More recently, the biomedical community has investigated zero-shot summarization methods for both single- and multi-document settings (Shaib et al., 2023). However, existing techniques generally lack fine-grained provenance for subclaims, such as their source documents, surrounding context, and reasoning trail and often operate at the abstract level (Tan et al., 2024). While commercial tools like (Elicit) and (Consensus) appear to retrieve relevant literature in response to scientific queries, their development processes are opaque. As of 2025, they typically attribute claims only at the document or paragraph level, without clear mappings to the context or assumptions underlying each individual claim.

A promising avenue of research focuses on developing source-aware training methods for LLMs. This approach aims to imbue LLMs with the ability to intrinsically link the knowledge they possess to the specific source documents from their pre-training data (Khalifa et al., 2024). This method of intrinsic source citation offers a potential way to trace the origins of information directly within the model’s parameters, providing an alternative to relying solely on external retrieval mechanisms. However, it faces challenges related to the practicalities of collecting and managing the vast amounts of source data and the specificity required for document identifiers (Khalifa et al., 2024).

| Feature        | GPT-4.1   | Gemini 2.5 | LLaMA 3-70B |
|----------------|-----------|------------|-------------|
| Link generated | ✓ Yes     | ✓ Yes      | ✓ Yes       |
| Link accurate  | ✗ Halluc. | ✗ Halluc.  | ✗ Halluc.   |

Table 1: Link generation and accuracy for leading LLMs. “Halluc.” indicates that the title or content is often hallucinated.

**Positioning SciTrue** A range of research prototypes and commercial products now provide scientific summarisation, literature review, or evidence retrieval, often integrating large language models (LLMs) to improve fluency and coverage. However, most outputs remain overly verbose, inflexible (typically adhering to systematic review templates), and lack precise mappings between synthesized claims and underlying evidence. Attribution is frequently surface-level or error-prone, and users are given limited support to verify individual inference steps. In contrast, **SciTrue** implements principled, context-sensitive, source-level attribution for each decomposed claim component. It explicitly links each subclaim to a specific, verifiable passage in the scientific literature, enabling user-centric audit and challenge. To our knowledge, SciTrue is among the first openly available systems to operationalize fine-grained attribution through an interface purpose-built for scientific accountability and transparent verification.

### 3 SciTrue: Transparent Scientific Claim Verification

SciTrue is an interactive, end-to-end platform for verifying scientific claims by combining large language models (LLMs) with retrieval from scientific literature (see the high-level workflow in Figure 2). This section describes the system’s data flow, user experience, and transparency features.

#### 3.1 Motivation and Overview

SciTrue is designed to assist scientists, journalists, students, and the general public in appraising scientific claims by grounding each claim in the scientific literature and making all evidence and reasoning steps explicit. Users interact through a web-based interface that allows them to submit claims and receive synthesized, well-attributed explanations directly linked to the original scientific sources.

**Claim Refinement:** A user begins by entering a scientific claim along with the desired number

of scientific articles for analysis. If the query is unclear or contains abbreviations, SciTrue’s query refinement agent transforms it into a well-formed scientific statement. If refinement is not possible, SciTrue prompts the user to revise their query, ensuring high-quality retrieval and synthesis in subsequent stages.

**Literature Retrieval:** Leveraging the Semantic Scholar API, which indexes over 214 million scientific papers across all fields of science, SciTrue retrieves articles relevant to each claim from a large database of scientific publications. For each selected article, SciTrue agents extract metadata (e.g., title, authors, journal, year, citation count) and identify candidate abstracts and paragraphs. If the article is deemed relevant, SciTrue extracts the most pertinent sentence and analyzes it using the rubric from Wei (2023), which categorizes evidence as follows:

- Declares something is better.
- Proposes something new.
- Describes a new finding or cause-effect relationship.

A SciTrue agent then assesses whether the evidence fully or conditionally supports or refutes the claim, identifies key assumptions or conditions underlying this support or refutation (such as a limited sample size, the study being conducted on mice rather than humans, or findings specific to a particular demographic), and determines its rhetorical role (i.e., the context in which the evidence is used in the article—such as a main finding, background information, or limitation). It also evaluates the relationship between the claim and the retrieved article by considering the article’s title, abstract, and relevant paragraphs, categorizing the strength of this relationship as strong, medium, or weak. This enables SciTrue to create context-aware, interpretable, and fine-grained links between scientific claims and supporting literature.

**Summary and Verdict Generation :** A SciTrue agent synthesises all supporting and contradictory evidence into a coherent, concise summary by considering all evidence parameters identified in the previous step, including underlying assumptions. Each statement in the summary is precisely linked to a unique source article via a clickable citation,

| Scientific Claim                                                         | Domain                  |
|--------------------------------------------------------------------------|-------------------------|
| Artificial sweeteners are healthier than sugar.                          | Health and Nutrition    |
| Vitamin D supplements prevent respiratory infections.                    | Health and Nutrition    |
| Nuclear waste can't be made safe for the long term.                      | Environment and Climate |
| Renewable energy deployment requires more mining overall.                | Environment and Climate |
| Universal basic income will eliminate poverty.                           | Social Science          |
| School uniforms reduce bullying and misbehaviour.                        | Social Science          |
| Artificial intelligence will inevitably lead to widespread unemployment. | Technology and Policy   |
| AI-based hiring tools reduce human biases in recruitment.                | Technology and Policy   |

Table 2: Examples of scientific claims used in our study and their associated domains.

ensuring transparency and traceability. The system also generates a structured verdict and succinct justification, indicating whether the claim is fully, mostly, partially, or not supported by the aggregated evidence.

**Subclaim and Evidence Interface:** Alongside the summary and verdict, a SciTrue agent automatically decomposes the main claim into underlying subclaims. For each subclaim, the agent retrieves and aligns relevant supporting or refuting passages from primary sources. Where available, citation metrics and journal impact factors are included to contextualize the strength of the evidence. For each subclaim, the system provides the following information:

- **Subclaim expression:** The specific subclaim derived from the summary.
- **Relevant sentence:** The sentence from the source that supports or refutes the subclaim.
- **Exact paragraph:** The paragraph in which the subclaim appears.
- **Contribution label:** The extent to which the subclaim corroborates or contrasts with the main claim.
- **Supporting and refuting assumptions:** Key assumptions that underpin support for or refutation of the subclaim.
- **Rich metadata:** Authors, year, title, journal/venue, section, and a clickable source link.
- **Explicit evidence label:** Indicates whether the evidence fully or conditionally supports or refutes the claim.
- **Context:** The role the evidence plays within the article (e.g., main finding, background information, or limitation).

- **Claim-article relationship:** The relationship between the claim and the retrieved article, categorizing the strength of this relationship as strong, medium, or weak.

- **Quantitative indicators (where available):** Citation counts and impact metrics.

**Interactive Exploration and History:** Users can browse previously analyzed claims in a personalized history, inspect executive summaries, subclaims, and all evidence details with expandable views, and revisit or reevaluate prior results. This workflow supports longitudinal and collaborative usage.

#### Language Model, Retrieval, and UI Choices

While our methodology is compatible with any retrieval system or large language model, our current implementation utilizes the Semantic Scholar API<sup>3</sup> for its comprehensive coverage. For text generation, we employ GPT-4o<sup>4</sup>, chosen after evaluating alternatives such as GPT-4.1, Gemini 2.5, and LLaMA 3 (70B). We selected GPT-4o due to its favorable balance of parsing accuracy, response speed, and lower API costs in small-scale experiments. To further manage costs, we limited the number of processed articles to a maximum of 15. For the user interface, we combined the Streamlit<sup>5</sup> library with front-end technologies including CSS, JavaScript, and HTML.

## 4 Evaluation

### 4.1 Dataset

Existing datasets in the scientific fact-checking domain, such as SciFact (Wadden et al., 2020) and Multi2Claim (Tan et al., 2023), are primarily designed for closed-domain setups. In contrast, SciTrue is designed for an open-domain setting, enabling it to track and incorporate the most recent

<sup>3</sup><https://www.semanticscholar.org/product/api>

<sup>4</sup><https://openai.com/index/gpt-4o-system-card/>

<sup>5</sup><https://streamlit.io/>

| Feature                         | GPT-4o-search-preview | Perplexity | SciTrue |
|---------------------------------|-----------------------|------------|---------|
| Summary Traceability            | 86.36%                | 85.71%     | 98.50%  |
| Overall Verdict                 | 88.00%                | 85.71%     | 97.80%  |
| Title Information               | 93.00%                | 86.00%     | 99.00%  |
| Author Information              | 68.51%                | 67.78%     | 96.50%  |
| Scientific Validity             | 48.18%                | 43.81%     | 94.20%  |
| Factual Accuracy of Attribution | 74.80%                | 80.95%     | 96.70%  |
| Context & Assumptions           | 54.80%                | 51.43%     | 95.30%  |
| Contribution Label Attribution  | 54.80%                | 71.43%     | 94.00%  |
| Source Credibility              | 5.6%                  | 7.5%       | 92.80%  |

Table 3: Human evaluation results comparing SciTrue with GPT-4o-Search-Preview and Perplexity Sonar Pro across key scientific attribution measures. The table reflects inter-annotator agreement, and the scores represent accuracy, as all evaluation questions were binary (Yes/No).

developments in scientific research. This open-domain approach not only allows for greater adaptability but also facilitates the identification and synthesis of 'grey areas' in science, namely, contested claims where multiple perspectives may coexist. Addressing such areas is essential, as the existence of debate, uncertainty, and evolving viewpoints is a fundamental aspect of the scientific process (Kuhn, 1962).

To this end, we tasked large language models (LLMs) with generating scientific claims across four key areas: health and nutrition, environment and climate, social science, and technology and policy. Representative examples of these claims are provided in Table 2, which also contains the curated subset used for our evaluation. The full dataset, comprising more than 3,000 scientific claims, is available in the project repository for community use.

## 4.2 Human Evaluation

We conducted a human evaluation comparing with two leading search-augmented LLMs: GPT-4o-Search-Preview and Perplexity Sonar Pro. Two independent annotators, each holding at least a bachelor's degree in science, assessed approximately 300 attributions, analyzing over 900 full-text sources (300 per model across three models), generated from 60 scientific claims. The claims were evenly distributed across four domains: health and nutrition, environment and climate, social science, and technology and policy. Five articles were retrieved per claim (60 claims  $\times$  5 articles = 300 attributions per model), a limit imposed due to current constraints in the capabilities of search-augmented LLMs (see an example at <https://chatgpt.com/share/6858d793-00fc-8002-9333-434f10b41166>).

In summary, we evaluated the three systems: SciTrue, GPT-4o-search-preview, and Perplexity Sonar Pro—across the following perspectives:

- **Summary Traceability:** Whether the summary is traceable to a specific, verifiable source and accurately reflects that source.
- **Overall Verdict:** Whether the model's overall conclusion and reasoning correctly reflect the summary information.
- **Title and Author Information:** Whether the cited source exists and the title and the listed authors matches the linked content.
- **Scientific Validity:** Whether the cited source is an academic paper, including key metadata such as title, abstract, authors, and publication year.
- **Factual Accuracy of Attribution:** Whether the cited sentence, or a semantically equivalent statement, appears in the referenced article.
- **Context and Assumption Awareness:** Whether the system detects when a source is used out of context or based on misleading assumptions.
- **Contribution Label Attribution:** Whether the assigned contribution type (e.g., corroborating, contrasting) is accurate within the sub-claim's context.
- **Source Credibility:** Whether the cited source's credibility (e.g., citation, impact factor) is accurately reflected.

Annotators followed a detailed rubric (see Appendix A.1). A screenshot of the full evaluation

interface is shown in Figure 3 in the appendix and is also available via the live demo for voluntary evaluation by future users. Inter-annotator agreement was 90%. Any remaining discrepancies were resolved through discussion.

## 5 Results and Discussion

SciTrue consistently outperformed across every perspective assessed. A closer analysis revealed key shortcomings of the competing systems. Notably, GPT-4o-search-preview and Perplexity Sonar Pro frequently retrieved blog posts, news sites, and general web pages as evidence rather than authoritative scientific articles; both systems included less than fifty percent scientific articles among their sources. These less rigorous sources often lack transparent vetting and replicability, severely limiting their usefulness for scientific verification. Moreover, despite relying on these sources, the models frequently hallucinated citations or misattributed findings to unreliable content, further undermining trust and underscoring the critical need for careful source selection in scientific fact-checking.

In contrast, SciTrue consistently prioritized primary scientific literature from reputable publication venues, ensuring that every claim was linked to a verifiable, high-credibility source. Its transparent audit trail not only links claims to sources, but also explicitly documents the attribution process—empowering users to replicate or challenge system decisions. This level of accountability is a foundational requirement for trustworthy AI in both research and journalism.

The results underscore key challenges faced by retrieval augmented LLM-based scientific assistants, particularly in guaranteeing the credibility and appropriateness of retrieved evidence. SciTrue’s ability to outperform mainstream models highlights the need for rigorous source filtering and transparent attribution mechanisms. Looking forward, future research should explore source-aware training methods via prevention of non-scholarly source retrieval, robustness to adversarial or ambiguous cases, and interfaces that support community-driven auditing and feedback. Such advances are critical for the safe and effective integration of AI into scientific, journalistic, and public knowledge environments.

## 6 Conclusion

SciTrue addresses critical gaps in scientific claim verification by emphasizing evidence traceability, rigorous attribution, and context-aware reasoning. This is vital for claims in gray areas of science, where demographic factors, time period, and location shape interpretation. By incorporating these contexts, SciTrue provides nuanced, reliable assessments.

Targeted at researchers, journalists, and policymakers who require trust and auditability, SciTrue offers transparent, source-linked outputs that enable informed debate and scrutiny foundations of trustworthy science and public discourse. It overcomes limitations of both standalone and retrieval-augmented LLMs through granular source mapping.

Following a stepwise pipeline of claim refinement, literature retrieval, evidence extraction, linked summarization, subclaim decomposition, and interactive UI output, SciTrue outperforms retrieval-augmented LLM baselines in human evaluations across multiple metrics, delivering accurate, verifiable, and contextually aware scientific claim assessments.

## 7 Limitations and Ethical Considerations

**Limitations** Our approach depends on the coverage and metadata quality of the Semantic Scholar API. Although it is a broad and openly accessible resource, representation varies across disciplines, and some relevant publications may be missing, paywalled, or not digitally available. As a result, SciTrue is currently limited to the literature indexed by Semantic Scholar, and its effectiveness may be reduced in areas where key source documents are not accessible. These constraints reflect the current state of scholarly data infrastructure, and future integrations with additional corpora may help broaden SciTrue’s reach.

**Ethical Considerations** SciTrue aims to reduce hallucinated attributions and improve verifiability, but automated systems can still inherit biases or overlook contextual nuances present in the underlying data. For this reason, SciTrue should be viewed as an assistive tool, supporting rather than replacing expert judgment, close reading, and peer review. Responsible use involves interpreting its recommendations critically and in conjunction with domain expertise.

## References

- Isabelle Augenstein, Sebastian Riedel, Pontus Stenetorp, Leon Derczynski, Georgios Spithourakis, Alex Dutton, and Yang Ji. 2019. Multifc: A real-world multi-domain dataset for evidence-based fact checking of claims. In *NAACL-HLT*.
- Marcel Binz, Stephan Alaniz, Adina Roskies, Balazs Aczel, Carl T Bergstrom, Colin Allen, Daniel Schad, Dirk Wulff, Jevin D West, Qiong Zhang, Richard M Shiffrin, Samuel J Gershman, Vencislav Popov, Emily M Bender, Marco Marelli, Matthew M Botvinick, Zeynep Akata, and Eric Schulz. 2025. How should the advancement of large language models affect the practice of science? *Proc. Natl. Acad. Sci. U. S. A.*, 122(5):e2401227121.
- Julia G Bottesini, Christie Aschwanden, Mijke Rhemtulla, and Simine Vazire. 2023. How do science journalists evaluate psychology research? *Adv. Methods Pract. Psychol. Sci.*, 6(3).
- Lihu Chen, Shuojie Fu, Gabriel Freedman, Cemre Zor, Guy Martin, James Kinross, Uddhav Vaghela, Ovidiu Serban, and Francesca Toni. 2025. Pub-guard-LLM: Detecting fraudulent biomedical articles with reliable explanations. *arXiv [cs.CL]*.
- Consensus. <https://consensus.app/>.
- Jay DeYoung, Iz Beltagy, Madeleine van Zuylen, Bailey Kuehl, and Lucy Lu Wang. 2021. Ms<sup>2</sup>: Multi-document summarization of medical studies. In *EMNLP*.
- Elicit. <https://elicit.org/>.
- Google DeepMind. 2024. [Gemini 1.5 technical report](#). Accessed June 2025.
- Antonios Kanavouras and Frank Coutelieres. 2020. Similarity among physical phenomena recognized on the basis of the classification of existing knowledge. *Journal of Mathematical Sciences and Modelling*, 3(2):47–54.
- Muhammad Khalifa, David Wadden, Emma Strubell, Honglak Lee, Lu Wang, Iz Beltagy, and Hao Peng. 2024. Source-aware training enables knowledge attribution in language models. *arXiv [cs.CL]*.
- Thomas S. Kuhn. 1962. *The Structure of Scientific Revolutions*. University of Chicago Press.
- Meta. 2025. [Gemini 2.5 report](#).
- Brent Mittelstadt, Sandra Wachter, and Chris Russell. 2023. To protect science, we must use LLMs as zero-shot translators. *Nat. Hum. Behav.*, 7(11):1830–1832.
- OpenAI. 2023. [Gpt-4 technical report](#). *arXiv preprint arXiv:2303.08774*.
- Chantal Shaib, Millicent Li, Sebastian Joseph, Iain James Marshall, Junyi Jessy Li, and Byron Wallace. 2023. Summarizing, simplifying, and synthesizing medical evidence using gpt-3 (with varying success). *ArXiv*, abs/2305.06299.
- Neset Tan, Trung Nguyen, Josh Bensemann, Alex Peng, Qiming Bao, Yang Chen, Mark Gahegan, and Michael Witbrock. 2023. Multi2Claim: Generating scientific claims from multi-choice questions for scientific fact-checking. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2652–2664, Dubrovnik, Croatia. Association for Computational Linguistics.
- Neşet Özkan Tan, Niket Tandon, David Wadden, Oyvind Tafjord, Mark Gahegan, and Michael Witbrock. 2024. Faithful reasoning over scientific claims. *Proceedings of the AAAI Symposium Series*, 3(1):263–272.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Majid Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. In *NAACL-HLT*, pages 809–819.
- Deepa Tilwani, Revathy Venkataramanan, and Amit P Sheth. 2024. Neurosymbolic AI approach to attribution in large language models. *arXiv [cs.CL]*.
- David Wadden, Shanchan Lin, Kyle Lo Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7534–7546.
- Byron C. Wallace, Sayantani Saha, Frank Soboczenski, and Iain James Marshall. 2021. Generating (Factual?) Narrative Summaries of RCTs: Experiments with Neural Multi-Document Summarization. In *AMIA Summits on Translational Science Proceedings*.
- Xin Wei. 2023. [ClaimDistiller: Scientific claim extraction with supervised contrastive learning](#). pages 65–77.

## A Appendix

### A.1 Instructions for Human Evaluators

This section describes the human evaluation interface and protocol used for claim and summary evaluations, as implemented in our Streamlit annotation tool.

**User Identification** Before beginning any evaluation, you must enter your **username** in the sidebar. This is required to personalize and track your progress.

**System Selection & Claim Loading** Using the sidebar, select the system output you wish to evaluate (e.g., Model1, Model2, or Model3). The interface will then display a list of available claims, marking those you have already evaluated. Select a claim to open its evaluation form.

**Summary and Claim Presentation** The main panel presents:

- The **Claim** (as provided by the system).
- The system-generated **Summary**.
- **Q1:** *Is the summary traceable to the cited sources (i.e., when you click sources, do you get the relevant page)?*
- Choose “Yes” or “No.” Optionally, provide a brief justification.
- If you select “No”, you **do not need to complete the rest of the form for this claim**; clicking Save will record all subclaim fields as “No” and end the evaluation for this entry.

#### Overall Verdict Consistency (Q2)

- **Q2:** *Does the system’s overall verdict/label correspond correctly to the summary?*
- Select “Yes” or “No”, and optionally state your reasoning.

**Subclaims Evaluation** For each listed subclaim, you are presented with detailed information such as citation, title, journal, authors, abstract, section, contribution, relevant sentence(s), and credibility scores if provided. For each subclaim, you must answer:

1. **Does the given URL open an article page?**  
(Select “Yes” or “No”)
2. **Does the title align with the URL?**  
 (“Yes” or “No”)
3. **Do the authors align with the URL?**  
 (“Yes” or “No”)
4. **Is the cited source an academic paper?**  
(i.e., does it have plausible title, abstract, authors, and year? “Yes”/“No”)
5. **Factuality:**  
*Does the subclaim (or its semantic equivalent) appear in both the summary and the cited*

*article?* (“Yes”/“No”)

*Note:* If you select “No” here, the following three questions will be automatically locked to “No”:

6. **Contribution Label:** Is the contribution label accurate in context? (“Yes”/“No”)
7. **Context & Assumption:** Does the model correctly detect context and assumptions? (“Yes”/“No”)
8. **Credibility/Impact Representation:** Is the cited source’s impact accurately represented ( $\pm 1$  for impact factor,  $\pm 10$  for citation count)? (“Yes”/“No”)

You may optionally provide reasoning for each subclaim evaluation.

**Saving and Error Checking** Before saving, the interface checks that all required fields are filled. Save your evaluation when complete. If this is the second time you are evaluating the same claim, your previous answers will be updated.

#### Notes and Best Practices

- **Carefully check** each source, title, and author match. Do not assume correctness based on citation style alone.
- **For factuality:** If the subclaim is not justified by the cited article, answer “No” even if it seems plausible.
- Use the “Show Paragraph” and “Show Abstract” toggles to see more context if needed.

**Support** For ambiguities or technical issues, contact the lead researcher.

#### A.2 Evaluation User Interface

##### A.3 LLM-as-judge evaluation

In addition to human annotation, we conducted an LLM-as-judge evaluation on a small subset of scientific claims (10 in total). For each claim, a state-of-the-art language model (GPT-4.1) was prompted with the same annotation criteria used in the human evaluation and asked to assess system outputs based on the provided evidence. The results of this evaluation showed only weak correlation with human annotators’ judgments, as accessing source articles often requires additional clicks, and evaluating information from full articles remains beyond the current capabilities

**User/Login**

Enter your username (required):

**Select JSON Source**

Choose a system output to evaluate:

Model1

Which claim would you like to evaluate?

Claim 2: High-protein diets harm kidney health.

This claim has NOT yet been evaluated by you.

## Summary and Claim Evaluation

Evaluating data from: [source1](#) | User: [user1](#)

**Claim:** High-protein diets harm kidney health.

**Generalized Summary:** Evidence regarding the impact of high-protein diets on kidney health is mixed. Some studies suggest that such diets may lead to glomerular hyperfiltration and potential kidney damage, particularly in individuals. For instance, a study in female Sprague–Dawley rats found that a diet with 35% of energy from protein led to kidney damage, including glomerular injury and renal hypertrophy. ([cambridge.org](#)) Similarly, a review highlighted that high protein intake, potentially resulting in kidney hyperfiltration and proteinuria. ([pubmed.ncbi.nlm.nih.gov](#)) However, other research indicates that in healthy individuals, increased protein intake does not adversely affect kidney function. protein intake within recommended ranges is consistent with normal kidney function in healthy adults. ([pubmed.ncbi.nlm.nih.gov](#)) Additionally, a study involving pre-diabetic older adults found that a higher protein intake was not associated with kidney function decline over one year. ([mdpi.com](#)) Therefore, while high-protein diets may pose risks for individuals with existing kidney issues, they do not appear to harm kidney health in healthy individuals.

**Q1: Is the summary traceable to the sources (i.e., when you click the sources, do you get the relevant page)? (Required)**

Yes  
 No

Reason for your choice (Optional)

**Accuracy / Verdict Label:** Partially True

**General Verdict and Reason:** The claim that high-protein diets harm kidney health is partially true. While such diets may exacerbate kidney issues in individuals with existing conditions, evidence does not consistently show harm in healthy individuals.

**Q2: Does the overall verdict correctly reflect the summary? (Required)**

Yes  
 No

Reason for your verdict (Optional)

### Subclaims Evaluation

**Subclaim 1: A diet with 35% of energy from protein led to kidney damage in female Sprague–Dawley rats.**

**Details**

**URL:** <https://www.cambridge.org/core/journals/british-journal-of-nutrition/article/diet-with-35-of-energy-from-protein-leads-to-kidney-damage-in-female-spraguedawley-rats/9891003F499906219852BA3FC06749D>

**Venue:** British Journal of Nutrition

**Authors:** Jia Y, Hwang SY, House JO, Ogborn MR, Weller HA, O K

**Year:** 2011

**Title:** A diet with 35% of energy from protein leads to kidney damage in female Sprague–Dawley rats

**Section:** Abstract

**Contribution:** completely supports

**Relevant Sentence:** Chronic casein intake in excess of 35 en % results in proteinuria and histological changes in normal and compromised rat kidneys.

Show Paragraph (Subclaim 1)

**Credibility Score:** Impact Factor: 3.334

**Supporting Assumptions:** High-protein intake leads to kidney damage in rats.

**Refuting Assumptions:** Rat models may not directly translate to human physiology.

Does the given URL open an article page?

Does the title align with the URL?

Do the authors align with the URL?

Is the cited source an academic paper (title, abstract, authors, year)?

Does the subclaim, or a semantically equivalent version of it, appear in both the summary and the cited article?

Figure 3: Screenshot of the evaluation interface. This interface is also available via the live demo.

#### **A.4 Automated System Prompt for Claim Evaluation (GPT-4o-search-preview)**

To generate system outputs for benchmark comparison, we used the following prompt and protocol with the GPT-4o-search-preview model. Each claim was presented alongside an explicit instruction to use exactly five scientific articles as evidence. The system was asked to produce a structured JSON object (and no extra text), suitable for downstream human annotation.

# PUCP-Metrix: An Open-source and Comprehensive Toolkit for Linguistic Analysis of Spanish Texts

Javier Alonso Villegas Luis<sup>†</sup> and Marco Antonio Sobrevilla Cabezudo<sup>†‡</sup>

<sup>†</sup>Research Group on Artificial Intelligence, Pontificia Universidad Católica del Perú

<sup>‡</sup>Aveni

{alonso.villegas, msobrevilla}@pucp.edu.pe

## Abstract

Linguistic features remain essential for interpretability and tasks that involve style, structure, and readability, but existing Spanish tools offer limited coverage. We present PUCP-Metrix, an open-source and comprehensive toolkit for linguistic analysis of Spanish texts. PUCP-Metrix includes 182 linguistic metrics spanning lexical diversity, syntactic and semantic complexity, cohesion, psycholinguistics, and readability. It enables fine-grained, interpretable text analysis. We evaluate its usefulness on Automated Readability Assessment and Machine-Generated Text Detection, showing competitive performance compared to an existing repository and strong neural baselines. PUCP-Metrix offers a comprehensive and extensible resource for Spanish, supporting diverse NLP applications.

## 1 Introduction

Linguistic features have gained renewed importance in explainable NLP, particularly for tasks requiring interpretability, stylistic sensitivity, or attention to surface-level properties. Despite advances in end-to-end neural models, recent work shows that handcrafted or derived features remain essential in applications such as AI-generated text detection (Kumarage et al., 2023; Ciccarelli et al., 2024; Petukhova et al., 2024), educational NLP (Mizumoto and Eguchi, 2023; Hou et al., 2025; Atkinson and Palma, 2025), and readability assessment (Zeng et al., 2024; Liu et al., 2025). In automated essay scoring, for instance, models incorporating linguistic features offer more transparent and pedagogically meaningful evaluations (Hou et al., 2025). These trends highlight the need for robust, modular repositories/toolkits that allow to extract linguistic metrics that complement deep models.

Beyond NLP applications, these toolkits also support linguistic research, offering standardized, quantifiable descriptions of texts across genres, reg-

isters, and proficiency levels (Jiang, 2016; Kuiken, 2023). They enable empirical analyses of morphosyntactic variation, cohesion (Gonzalez-Dios and Bengoetxea, 2021), or lexical sophistication (Crossley and Kyle, 2018), and facilitate cross-linguistic studies (Uçar et al., 2024).

Existing tools have demonstrated the value of this approach. For instance, Coh-Metrix (McNamara et al., 2010) provides extensive metrics for English across various linguistic levels. Similar resources include NILC-Metrix for Portuguese (Leal et al., 2023), Coh-Metrix-Esp for Spanish (Quispesaravia et al., 2016), and MultiAzterTest for Spanish, English, and Basque (Bengoetxea and Gonzalez-Dios, 2021). Nevertheless, the latter two exhibit limited metric coverage and present challenges related to installing and inference efficiency.

In this work, we introduce PUCP-Metrix, a new open-source toolkit for extracting linguistic metrics from Spanish texts. It extends coverage to lexical, syntactic, discourse, psycholinguistic, and readability dimensions, providing a total of 182 linguistic metrics. In addition, we demonstrate its utility in two downstream tasks: Automated Readability Assessment and Machine-Generated Text Detection.

Our main contributions are:

- PUCP-Metrix, a comprehensive and extensible open-source toolkit of linguistic metrics for Spanish, featuring metrics not available in existing resources.<sup>1</sup>
- A pip-installable, spaCy-based modular implementation enabling scalable extraction (multi-processing) and easy extension.
- An empirical study evaluating its usefulness in Automated Readability Assessment and Machine-Generated Text Detection.

<sup>1</sup>The code is available at <https://github.com/iapucp/pucp-metrix>.

## 2 PUCP-Metrix

PUCP-Metrix is a modular and extensible linguistic analysis toolkit for Spanish, designed to support both research and large-scale text processing through a Python library. Its architecture emphasizes flexibility and scalability, enabling users to efficiently extract a wide range of linguistic metrics from texts.

To achieve this, we leveraged the widely adopted NLP library spaCy for core processing tasks such as tokenization, part-of-speech tagging, and dependency parsing. Building on spaCy’s modular architecture, we developed custom pipelines that include both general-purpose and category-specific metrics implemented as reusable components. This design allows external users to easily extend or modify the system, ensuring that PUCP-Metrix remains both efficient and adaptable to new linguistic analyses.

### 2.1 Linguistic Categories and Metrics

We employed an open-source Spanish implementation of *Coh-Metrix* (Quispesaravia et al., 2016) to collect initial linguistic metrics and guide our design. To develop additional metrics, we examined the implementations provided by tools such as MultiAzterTest and NILC-Metrix. Overall, we compiled 182 linguistic metrics for Spanish texts. The complete list is available at Appendix A.

- **Descriptives:** 27 indicators that capture general statistics of the text, such as *number of words*, *number of sentences*, *number of paragraphs*, *minimum and maximum length of sentences*, *average word length*.
- **Lexical Diversity:** 22 indicators measure the diversity of the text’s vocabulary, including the *type-token ratio for various word categories (nouns, verbs, etc.)*, *noun density*, *verb density*, *adverb density*, and *adjective density*. Our implementation extends these measures with type-token ratios for additional word categories and their lemmatized forms. Another key indicators include the following:
  - *MTLD (Measure of Textual Lexical Diversity)*: Addresses TTR’s length sensitivity by calculating the average length of sequential word segments that maintain a certain TTR threshold, providing more stable measures across varying text lengths (McCarthy Philip M, 2010).
  - *VOCd (Vocabulary Complexity Diversity)*: Estimates vocabulary richness through curve-fitting techniques on random samples, offering insights into the probability of encountering new word types (McCarthy Philip M, 2010).
  - *Maas Index*: A logarithmic transformation that provides an alternative measure of lexical diversity, particularly useful for comparing texts of different lengths (Mass, 1972).
- **Readability:** 7 indicators that represent how difficult to understand the text is, such as *Flesch Grade Level*, *Brunet Index*, *Gunning Fog Index*, *Honore’s Statistic*, *SMOG Grade*, *The Szigriszt-Pazos Perspicuity Index* and *Readability  $\mu$* . Among the important measures are:
  - *Flesch Grade Level (Fernández-Huertas adaptation)* (Fernández Huerta, 1959): Adapted for Spanish texts, this measure estimates the grade level required for comprehension.
  - *Brunet Index*: A readability measure of lexical richness, where lower values indicate greater vocabulary diversity.
  - *Gunning Fog Index*: Calculates readability by considering both sentence length and complex word percentage, estimating the education level needed for comprehension.
  - *Honore’s Statistic*: Measures vocabulary richness by analyzing hapax legomena (words appearing only once).
  - *SMOG Grade*: Estimates the years of education required to understand a text by analyzing polysyllabic words (3+ syllables).
  - *Szigriszt-Pazos Perspicuity Index*: A Spanish-specific readability measure that evaluates text clarity, offering insights into Spanish text comprehensibility.
  - *Readability  $\mu$* : A statistical measure that evaluates text complexity through letter distribution patterns.
- **Syntactic Complexity:** 12 indicators, reflecting the structural intricacy of text, such as *proportion of sentences with 1-7 clauses*, *minimal edit distances of words*, *POS tags and lemmas*.

Following *Coh-Matrix*, our implementation extends minimal edit distance measures to POS tags and lemmatized forms, providing comprehensive syntactic variation analysis.

- **Psycholinguistics:** 30 indicators, reflecting psycholinguistic properties of words, specifically how they are understood by humans: *concreteness*, *imageability*, *familiarity*, *age of acquisition*, *valence* and *arousal*. These psycholinguistic properties were collected from the EsPal database (Duchon et al., 2013) and works from Stadthagen-Gonzalez et al. (2017):
  - *Concreteness*: Measures the degree to which words refer to tangible, physical objects versus abstract concepts. Higher concreteness values indicate words that are easier to visualize and process cognitively.
  - *Imageability*: Assesses how easily words can evoke mental images. Words with higher imageability are processed more quickly and remembered more easily.
  - *Familiarity*: Evaluates how well-known words are to speakers. Familiar words are processed faster and require less cognitive effort.
  - *Age of Acquisition*: Measures the age at which words are typically learned. Earlier acquired words are processed more automatically and efficiently.
  - *Valence*: Assesses the emotional positivity or negativity of words. Valence influences emotional processing and memory formation.
  - *Arousal*: Measures the emotional intensity or activation level of words. Arousal affects attention and memory consolidation.
- **Word Information:** 24 indicators with more detailed word-level statistics, such as: *number of nouns*, *number of verbs*, *number of adverbs*, *number of adjectives* and *number of content words*.
- **Referential Cohesion:** 12 indicators that serve to measure the interconnections within a text: *noun overlap*, *argument overlap*, *stem overlap*, *content word overlap* and *anaphor overlap*.

- **Textual Simplicity:** 4 indicators, measuring the simplicity of the text using the ratio of short or large sentences, such as: *proportion of short sentences*, *proportion of medium sentences*, *proportion of long sentences*, *proportion of very long sentences*.
- **Semantic Cohesion:** 8 indicators, assessing the degree of semantic relatedness between different parts of the text, such as: *Latent Semantic Analysis (LSA) overlap of adjacent sentences*, *LSA overlap of all sentences*, *LSA overlap of adjacent paragraphs*.
- **Word Frequency:** 16 indicators, various measurements involving the Zipf’s frequency<sup>2</sup> for different kinds of words, such as *rare nouns count*, *rare verbs count*, *rare adverbs count*, *rare content words count*<sup>3</sup> and *mean word frequency*.
- **Syntactic Pattern Density:** 14 indicators, reflecting the density of various syntactic elements, such as: *noun phrase density*, *verb phrase density*, *negative expressions density*, *coordinating conjunctions density* and *subordinating conjunction density*.
- **Connectives:** 6 indicators, measuring the use of words or phrases that establish logical, temporal, or other relationships between different parts of the text, such as: *casual connectives incidence*, *logical connectives incidence*, *adversative connectives incidence*, *temporal connectives incidence*, *additive connectives incidence*, *all connectives incidence*.

## 2.2 Comparison with Existing Tools

Table 1 compares Coh-Matrix-Esp, MultiAzterTest, and PUCP-Matrix across three practical dimensions: ease of installation, processing speed, and metric coverage. Existing tools present complementary strengths but also notable trade-offs. Coh-Matrix-Esp offers a stable implementation but covers a limited number of metrics, while MultiAzterTest provides broader coverage but relies on external dependencies (e.g., syntactic parsers), making installation more complex and increasing inference time, which can hinder scalability in large-scale settings.

<sup>2</sup>Zipf’s frequency is estimate by using the wordfreq tool. It is available at <https://github.com/rspeer/wordfreq/>.

<sup>3</sup>Rare words were defined in a similar way as Bengoetxea and Gonzalez-Dios (2021).

PUCP-Metrix is designed to address these limitations by prioritizing coverage, scalability, and usability. It consolidates 182 linguistic metrics spanning readability, lexical diversity, syntactic complexity, cohesion, and psycholinguistic dimensions, exceeding the coverage of Coh-Metrix-Esp (48 metrics) and MultiAzterTest (141 metrics). The toolkit is implemented as a pure Python library built on spaCy, enabling a modular and extensible architecture that integrates naturally into modern NLP pipelines. In addition, its pip-based distribution and built-in multiprocessing support facilitate reproducible and scalable feature extraction without reliance on heavyweight external toolchains.

While some individual metrics build on prior work, our primary contribution is an integrated and easily deployable framework that emphasizes reproducibility and practical usability. By combining broad metric coverage with lightweight deployment and scalable processing, PUCP-Metrix aims to lower the barrier to large-scale linguistic feature extraction for Spanish NLP research.

|                    | Easy to install | Processing speed | Metric coverage |
|--------------------|-----------------|------------------|-----------------|
| Coh-Metrix-Esp     | ✗               | Fast             | 48              |
| MultiAzterTest     | ✗               | Slow             | 141             |
| PUCP-Metrix (ours) | ✓               | Fast             | 182             |

Table 1: Comparison of PUCP-Metrix with existing tools for Spanish.

Table 2 shows the number of linguistic metrics implemented in Coh-Metrix-Esp, MultiAzterTest and PUCP-Metrix (ours). PUCP-Metrix provides a broader coverage of linguistic metrics compared to Coh-Metrix-Esp and MultiAzterTest, across 13 categories. Notably, PUCP-Metrix includes metrics in categories that are entirely missing or underrepresented in the other tools, such as semantic cohesion, textual simplicity, and psycholinguistics. This way, PUCP-Metrix can capture higher-level discourse, cognitive readability, and psycholinguistic properties.

Furthermore, PUCP-Metrix distributes its metrics more evenly across lexical, syntactic, semantic, and psycholinguistic dimensions. This comprehensive and balanced coverage allows for a more detailed and nuanced characterization of texts, making PUCP-Metrix better suited for in-depth linguistic analysis and a wide range of NLP applications.

### 3 Applications

In order to verify the usefulness of PUCP-Metrix, we use it for two tasks where linguistic metrics have proven to be helpful in past work. In particular, we select Automated Readability Assessment (ARA) and Machine-Generated Text Detection.

#### 3.1 Automated Readability Assessment (ARA)

We adopt an approach similar to that of [Vásquez-Rodríguez et al. \(2022\)](#), who introduced a benchmark for ARA on Spanish texts. Their work unified both public and non-public corpora annotated with language proficiency levels and proposed two- and three-class classification schemas.

In contrast, our study comprises only four publicly available datasets —CAES, Coh-Metrix-Esp, Kwiziq, and HablaCultura— to ensure reproducibility and open accessibility. We adopt the same label mappings described in the paper, adapting all texts to two readability classification schemas: 2-label (simple, complex) and 3-label (basic, intermediate, advanced). The dataset’s descriptions and the labeling strategy can be found in [Appendix B](#).

Overall, the dataset contains 32,167 instances, distributed across the four sources as follows: 31,149 from CAES, 100 from Coh-Metrix-Esp, 206 from Kwiziq, and 713 from HablaCultura.

We experiment with two readability classification schemas mentioned before. All experiments are performed at the document level<sup>4</sup>. The corpus is divided into 80% training, 10% validation, and 10% test sets, stratified by label. We evaluate models using Precision, Recall, and F1-score.

#### 3.2 Machine-Generated Text Detection

We adopt the AuTextification 2023 shared task dataset ([Sarvazyan et al., 2023](#)), which comprises over 160,000 texts in English and Spanish across five domains: tweets, reviews, news, legal, and how-to articles generated by both human and large language models.

For our experiments, we focus on the Machine-generated Text Detection task, which consists of identifying if a text has been created by a human or a machine. The task includes 26,996 human-generated instances and 25,195 machine-generated instances, totaling 52,191 instances. More details about the dataset can be found in [Appendix B](#).

<sup>4</sup>We use the same texts that come from the available resources

| Category                  | Coh-Metrix-Esp | MultiAzterTest | PUCP-Metrix (ours) |
|---------------------------|----------------|----------------|--------------------|
| Descriptive               | 11             | 22             | <b>27</b>          |
| Referential Cohesion      | 12             | 10             | 12                 |
| Lexical Diversity         | 2              | 20             | <b>22</b>          |
| Readability               | 1              | 1              | <b>7</b>           |
| Connectives               | 6              | 12             | 6                  |
| Syntactic Complexity      | 2              | 19             | 12                 |
| Pattern Density           | 3              | 0              | <b>14</b>          |
| Semantic Cohesion         | 0              | 0              | <b>8</b>           |
| Word Information          | 11             | 32             | 24                 |
| Word Frequency            | 0              | 15             | <b>16</b>          |
| Textual Simplicity        | 0              | 0              | <b>4</b>           |
| Psycholinguistics         | 0              | 0              | <b>30</b>          |
| Word Semantic Information | 0              | 4              | 0                  |
| Semantic Overlap          | 0              | 6              | 0                  |
| Total                     | 48             | 141            | 182                |

Table 2: Number of linguistic metrics per category for each tool.

### 3.3 Models

For both tasks, we trained several machine learning models—Logistic Regression (LR), XGBoost (XGB), Support Vector Machines (SVM), and Random Forests (RF)—using all the metrics extracted with both MultiAzterTest and PUCP-Metrix.

Following the AuTextification shared task setup and consistent with the ARA formulation, we use a RoBERTa-based model (Fandiño et al., 2022) (RoBERTa-BNE)<sup>5</sup> as our baseline. We fine-tune this model on the official training splits and evaluate it on the corresponding test splits to ensure comparability. Similarly, we fine-tune RoBERTa-BNE on both the 2-label and 3-label ARA classification schemas.

## 4 Results and Discussion

### 4.1 Automated Readability Assessment

Table 3 reports the results for the 2-label ARA task. PUCP-Metrix slightly outperforms MultiAzter, achieving an overall F1 of 97.46 with XGBoost compared to 97.24. However, this difference is not significant. XGBoost consistently yields the highest scores, followed by Random Forest, while Logistic Regression and SVM lag behind. RoBERTa-BNE achieves the best overall F1 (98.30), indicating that deep contextual models capture subtle semantic patterns beyond what feature-based metrics provide.

Table 4 shows the 3-label ARA results. PUCP-Metrix again slightly surpasses MultiAzter (F1

<sup>5</sup>Model available at <https://huggingface.co/PlanTL-GOB-ES/roberta-base-bne>

| Model       | Simple |       |       | Complex |       |       | F1    |
|-------------|--------|-------|-------|---------|-------|-------|-------|
|             | P      | R     | F1    | P       | R     | F1    |       |
| MultiAzter  |        |       |       |         |       |       |       |
| LR          | 96.42  | 97.27 | 96.85 | 91.75   | 89.37 | 90.54 | 93.70 |
| XGB         | 98.05  | 99.20 | 98.62 | 97.57   | 94.19 | 95.85 | 97.24 |
| SVM         | 96.51  | 97.32 | 96.91 | 91.89   | 89.62 | 90.74 | 93.82 |
| RF          | 97.25  | 99.29 | 98.26 | 97.76   | 91.72 | 94.64 | 96.45 |
| PUCP-Metrix |        |       |       |         |       |       |       |
| LR          | 96.68  | 97.65 | 97.16 | 92.87   | 90.11 | 91.47 | 94.31 |
| XGB         | 98.38  | 99.08 | 98.73 | 97.22   | 95.18 | 96.19 | 97.46 |
| SVM         | 96.60  | 97.69 | 97.14 | 92.97   | 89.86 | 91.39 | 94.27 |
| RF          | 97.45  | 99.20 | 98.32 | 97.52   | 92.34 | 94.86 | 96.59 |
| RoBERTa-BNE | 99.04  | 99.24 | 99.14 | 97.76   | 97.16 | 97.46 | 98.30 |

Table 3: Results on 2-label ARA/Complexity Classification task

96.72 vs. 96.56), with XGBoost as the top-performing classifier. RoBERTa-BNE remains the strongest model, achieving an overall F1 of 98.13 and near-perfect performance on Basic and Intermediate texts.

### 4.2 Machine-Generated Text Detection

Table 5 reports the performance of machine learning models using PUCP-Metrix and MultiAzter metrics, alongside a RoBERTa-BNE model fine-tuned on AuTextification and the shared task’s best-reported results.

PUCP-Metrix consistently outperforms MultiAzter. For human texts, F1 increases from 42–51 (MultiAzter) to 60–66, and for machine texts from 70–73 to 71–76, showing its ability to capture linguistic and structural cues. Tree-based models, especially XGBoost and Random Forest, benefit most, achieving the highest overall F1.

Compared to RoBERTa-BNE, PUCP-Metrix provides more balanced class performance. While RoBERTa-BNE attains high precision on human

| Model       | Basic |       |       | Intermediate |       |       | Advanced |       |       | F1    |
|-------------|-------|-------|-------|--------------|-------|-------|----------|-------|-------|-------|
|             | P     | R     | F1    | P            | R     | F1    | P        | R     | F1    |       |
| Multiazter  |       |       |       |              |       |       |          |       |       |       |
| LR          | 91.43 | 92.56 | 91.99 | 85.66        | 86.30 | 85.97 | 83.00    | 74.20 | 78.36 | 85.44 |
| XGB         | 97.62 | 98.59 | 98.10 | 96.43        | 96.59 | 96.51 | 98.48    | 91.87 | 95.06 | 96.56 |
| SVM         | 90.54 | 93.08 | 91.79 | 85.29        | 85.71 | 85.50 | 84.72    | 68.55 | 75.78 | 84.36 |
| RF          | 96.32 | 98.07 | 97.18 | 94.38        | 94.93 | 94.66 | 98.37    | 85.16 | 91.29 | 94.38 |
| PUCP-Metrix |       |       |       |              |       |       |          |       |       |       |
| LR          | 92.25 | 92.85 | 92.55 | 86.35        | 86.71 | 86.53 | 82.02    | 77.39 | 79.64 | 86.24 |
| XGB         | 97.68 | 98.59 | 98.13 | 97.16        | 96.59 | 96.88 | 96.72    | 93.64 | 95.15 | 96.72 |
| SVM         | 91.10 | 93.55 | 92.31 | 86.06        | 85.63 | 85.85 | 82.72    | 71.02 | 76.43 | 84.86 |
| RF          | 95.55 | 98.18 | 96.85 | 95.11        | 93.77 | 94.44 | 97.63    | 87.28 | 92.16 | 94.48 |
| RoBERTa-BNE | 99.30 | 99.24 | 99.27 | 98.83        | 98.42 | 98.63 | 95.50    | 97.53 | 96.50 | 98.13 |

Table 4: Results on 3-label ARA/Complexity Classification task

| Model          | Human |       |              | Machine |       |       | F1           |
|----------------|-------|-------|--------------|---------|-------|-------|--------------|
|                | P     | R     | F1           | P       | R     | F1    |              |
| Multiazter     |       |       |              |         |       |       |              |
| LR             | 70.52 | 30.28 | 42.37        | 61.84   | 89.93 | 73.29 | 57.83        |
| XGB            | 68.10 | 39.73 | 50.18        | 63.98   | 85.19 | 73.08 | 61.63        |
| SVM            | 70.43 | 30.74 | 42.80        | 61.95   | 89.73 | 73.30 | 58.05        |
| RF             | 62.08 | 43.98 | 51.49        | 63.82   | 78.62 | 70.45 | 60.97        |
| PUCP-Metrix    |       |       |              |         |       |       |              |
| LR             | 71.09 | 55.93 | 62.61        | 70.02   | 81.90 | 75.49 | 69.05        |
| XGB            | 71.34 | 61.36 | <b>65.97</b> | 72.33   | 80.38 | 76.14 | <b>71.06</b> |
| SVM            | 71.04 | 56.05 | 62.66        | 70.06   | 81.82 | 75.48 | 69.07        |
| RF             | 63.57 | 58.24 | 60.79        | 68.85   | 73.44 | 71.07 | 65.93        |
| RoBERTa-BNE    | 86.28 | 46.87 | 60.74        | 68.99   | 94.07 | 79.60 | 70.17        |
| RoBERTa-Autex* | -     | -     | -            | -       | -     | -     | 68.52        |
| Best model*    | -     | -     | -            | -       | -     | -     | 70.77        |

Table 5: Results on AuTexTification. \*The authors of the shared task only provide F1 in the report.

texts (86.28), low recall (46.87) limits its F1, suggesting that contextual embeddings may miss the diversity of human writing. PUCP-Metrix also slightly surpasses the shared task’s best-reported F1 (70.77), indicating that integrating linguistic features with neural models could further improve results.

Finally, we analyze which linguistic metrics contribute most to classification. Overall, detecting machine-generated text depends primarily on features related to frequency, readability, and cohesion, whereas ARA tasks are driven by descriptive, syntactic, and simplicity-related features. Full details of the feature analysis are provided in Appendix C.

## 5 Tool Usage

PUCP-Metrix can be installed via pip:

```
pip install iapucp-metrix
```

To use the library, we need to import the Analyzer class and call `compute_metrics` to compute all metrics. The function supports multiprocessing through spaCy, allowing us to specify the number of workers and the batch size.

```
from iapucp_metrix.analyzer import Analyzer

analyzer = Analyzer()

texts = ["Este es mi ejemplo."]

metrics_list = analyzer.compute_metrics(
    texts,
    workers=4,
    batch_size=2
)

for i, metrics in enumerate(metrics_list):
    print(Readability(Fernández-Huertas):)
    print(f"{metrics['RDFHGL']:.2f}")
```

The output of the code described above is:

```
Readability (Fernández-Huertas):
201.86
```

In addition, PUCP-Metrix supports computing metrics grouped by linguistic categories (via `compute_grouped_metrics`), enabling users to analyze model behavior across dimensions such as lexical, syntactic, and semantic features.

## 6 Conclusion and Future Work

PUCP-Metrix provides a linguistically rich set of 182 metrics for Spanish, offering broader coverage and a larger metric set than previous resources. Empirical evaluations demonstrate its effectiveness in ARA and Machine-generated text detection tasks. Models trained on these metrics match or outperform baseline neural models, underscoring their ability to capture nuanced linguistic information.

Future work includes expanding the metric set to incorporate more discourse and pragmatic metrics, adapting PUCP-Metrix to other Spanish varieties, and integrating these metrics into pre-trained language models or NLP pipelines. Benchmarking on larger and more diverse tasks/datasets will further

validate its robustness and support the development of specialized metric sets.

## Limitations

The current evaluation has several limitations. Although PUCP-Metrix has been tested on multiple datasets, the experiments primarily focus on learner essays, children's texts, and selected AuTextification domains, leaving its performance on other genres and domains uncertain. Additionally, PUCP-Metrix depends heavily on spaCy-based linguistic processing and external lexicons (e.g., psycholinguistic norms), so parsing errors and coverage gaps in these resources can directly affect the reliability of the computed metrics.

## References

- John Atkinson and Diego Palma. 2025. An llm-based hybrid approach for enhanced automated essay scoring. *Nature: Scientific Reports*, 15.
- Francesco Barbieri, Luis Espinosa Anke, and Jose Camacho-Collados. 2022. XLM-T: Multilingual language models in Twitter for sentiment analysis and beyond. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 258–266, Marseille, France. European Language Resources Association.
- Kepa Bengoetxea and Itziar Gonzalez-Dios. 2021. Multiaztertest: A multilingual analyzer on multiple levels of language for readability assessment. *arXiv preprint arXiv:2109.04870*.
- Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. 2021. MultiEURLEX - a multi-lingual and multi-label legal document classification dataset for zero-shot cross-lingual transfer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6974–6996, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Vittorio Ciccarelli, Cornelia Genz, Nele Mastracchio, Wiebke Petersen, Anna Stein, and Hanxin Xia. 2024. Team art-nat-HHU at SemEval-2024 task 8: Stylistically informed fusion model for MGT-detection. In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 1690–1697, Mexico City, Mexico. Association for Computational Linguistics.
- Scott A. Crossley and Kristopher Kyle. 2018. Assessing writing with the tool for the automatic analysis of lexical sophistication (taales). *Assessing Writing*, 38:46–50.
- Andrew Duchon, Manuel Perea, Nuria Sebastián-Gallés, Antonia Martí, and Manuel Carreiras. 2013. Espal: One-stop shopping for spanish word properties. *Behavior Research Methods*, 45(4):1246–1258.
- Asier Gutiérrez Fandiño, Jordi Armengol Estapé, Marc Pàmies, Joan Llop Palao, Joaquin Silveira Ocampo, Casimiro Pio Carrino, Carme Armentano Oller, Carlos Rodriguez Penagos, Aitor Gonzalez Agirre, and Marta Villegas. 2022. Maria: Spanish language models. *Procesamiento del Lenguaje Natural*, 68.
- F Fernández Huerta. 1959. Medidas sencillas de lecturabilidad. *Consigna*, 214:29–32.
- Itziar Gonzalez-Dios and Kepa Bengoetxea. 2021. Multiaztertest@vaxxstance-iberlef 2021: Identifying stances with language models and linguistic features. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2021)*, volume 2943 of *CEUR Workshop Proceedings*, pages 192–201.
- Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. XLsum: Large-scale multilingual abstractive summarization for 44 languages. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4693–4703, Online. Association for Computational Linguistics.
- Zhaoyi Hou, Alejandro Ciuba, and Xiang Li. 2025. Improving llm-based automatic essay scoring with linguistic features. In *Proceedings of the Innovation and Responsibility in AI-Supported Education Workshop*, volume 273 of *Proceedings of Machine Learning Research*, pages 41–65. PMLR.
- Kevin Jiang. 2016. Douglas biber and bethany gray: Grammatical complexity in academic english: Linguistic change in writing. *Applied Linguistics*, 37.
- Folkert Kuiken. 2023. Linguistic complexity in second language acquisition. *Linguistics Vanguard*, 9(s1):83–93.
- Tharindu Kumarage, Joshua Garland, Amrita Bhattacharjee, Kirill Trapeznikov, Scott Ruston, and Huan Liu. 2023. Stylometric detection of ai-generated text in twitter timelines. *arXiv preprint arXiv:2303.03697*.
- Faisal Ladhak, Esin Durmus, Claire Cardie, and Kathleen McKeown. 2020. WikiLingua: A new benchmark dataset for cross-lingual abstractive summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4034–4048, Online. Association for Computational Linguistics.
- Sidney Evaldo Leal, Magali Sanches Duran, Carolina Evaristo Scarton, Nathan Siegle Hartmann, and Sandra Maria Aluísio. 2023. Nilc-metrix: assessing the complexity of written and spoken language in brazilian portuguese. *Language, Resources and Evaluation*, 58(1):73–110.

- Angela Leis, Francesco Ronzano, Miguel A Mayer, Laura I Furlong, and Ferran Sanz. 2019. Detecting signs of depression in tweets in spanish: Behavioral and linguistic analysis. *J Med Internet Res*, 21(6).
- Fengkai Liu, Tan Jin, and John S. Y. Lee. 2025. Automatic readability assessment for sentences: neural, hybrid and large language models. *Language Resources and Evaluation*.
- Heinz-Dieter Mass. 1972. Über den zusammenhang zwischen wortschatzumfang und länge eines textes. *Zeitschrift für Literaturwissenschaft und Linguistik*, 2(8):73.
- Jarvis Scott McCarthy Philip M. 2010. Mtd, vocd-d, and hd-d: A validation study of sophisticated approaches to lexical diversity assessment. In *Behavior Research Methods*, pages 381–392.
- Danielle S. McNamara, Max M. Louwerse, Philip M. McCarthy, and Arthur C. Graesser. 2010. Coh-matrix: Capturing linguistic features of cohesion. *Discourse Processes*, 47(4):292–330.
- Atsushi Mizumoto and Masaki Eguchi. 2023. Exploring the potential of using an ai language model for automated essay scoring. *Research Methods in Applied Linguistics*, 2(2):100050.
- M. Dolores Molina-González, Eugenio Martínez-Cámara, M. Teresa Martín-Valdivia, and L. Alfonso Ureña-López. 2014. Cross-domain sentiment analysis using spanish opinionated words. In *Natural Language Processing and Information Systems*, pages 214–219, Cham. Springer International Publishing.
- Giovanni Parodi. 2015. Corpus de aprendices de español (caes). *Journal of Spanish Language Teaching*, 2(2):194–200.
- Kseniia Petukhova, Roman Kazakov, and Ekaterina Kochmar. 2024. PetKaz at SemEval-2024 task 8: Can linguistics capture the specifics of LLM-generated text? In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 1140–1147, Mexico City, Mexico. Association for Computational Linguistics.
- Andre Quispesaravia, Walter Perez, Marco Sobrevilla Cabezudo, and Fernando Alva-Manchego. 2016. Coh-Matrix-Esp: A complexity analysis tool for documents written in Spanish. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4694–4698, Portorož, Slovenia. European Language Resources Association (ELRA).
- Areg Mikael Sarvazyan, José Ángel González, Marc Franco-Salvador, Francisco Rangel, Berta Chulvi, and Paolo Rosso. 2023. Overview of autextification at iberlef 2023: Detection and attribution of machine-generated text in multiple domains. *arXiv preprint arXiv:2309.11285*.
- Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2020. MLSUM: The multilingual summarization corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8051–8067, Online. Association for Computational Linguistics.
- Hans Stadthagen-Gonzalez, Constance Imbault, Miguel A. Pérez Sánchez, and Marc Brysbaert. 2017. Norms of valence and arousal for 14,031 spanish words. *Behavior Research Methods*, 49(1):111–123.
- Suna-Şeyma Uçar, Itziar Aldabe, Nora Aranberri, and Ana Arruarte. 2024. Exploring automatic readability assessment for science documents within a multilingual educational context. *International Journal of Artificial Intelligence in Education*, 34(4):1417–1459.
- Laura Vázquez-Rodríguez, Pedro-Manuel Cuenca-Jiménez, Sergio Morales-Esquivel, and Fernando Alva-Manchego. 2022. A benchmark for neural readability assessment of texts in Spanish. In *Proceedings of the Workshop on Text Simplification, Accessibility, and Readability (TSAR-2022)*, pages 188–198, Abu Dhabi, United Arab Emirates (Virtual). Association for Computational Linguistics.
- Jinshan Zeng, Xianchao Tong, Xianglong Yu, Wenyan Xiao, and Qing Huang. 2024. Interpretara: Enhancing hybrid automatic readability assessment with linguistic feature interpreter and contrastive learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19497–19505.

## A List of metrics in PUCP-Matrix

## B Datasets for Automated Readability Assessment and Machine-generated Text Detection

### B.1 Automated Readability Assessment

- CAES (*Corpus de Aprendices del Español*)<sup>6</sup> (Parodi, 2015). This corpus consists of essays written by learners of Spanish as a foreign language. Each document is annotated with a CEFR level (A1–C2). Following Vázquez-Rodríguez et al. (2022), we map A1–B1 to "simple" and B2–C2 to "complex" for the 2-label schema, and A1-A2 to "basic", B1-B2 to "intermediate" and C1-C2 to "advanced" for the 3-label schema.
- Coh-Matrix-Esp (Quispesaravia et al., 2016). This dataset is a collection of short Spanish stories that includes children’s tales and texts intended for adults. It provides explicit simple and complex labels, directly aligned to our

<sup>6</sup>Available at <https://galvan.usc.es/caes/>

| Category                                                                 | Metric Description                                                  | Category                                                          | Metric Description                                                      |                                                     |                                        |
|--------------------------------------------------------------------------|---------------------------------------------------------------------|-------------------------------------------------------------------|-------------------------------------------------------------------------|-----------------------------------------------------|----------------------------------------|
| Descriptive Indices                                                      | DESPC: Paragraph count                                              | Syntactic Complexity Indices                                      | SYNNP: Mean number of modifiers per noun phrase                         |                                                     |                                        |
|                                                                          | DESPCi: Paragraph count incidence per 1000 words                    |                                                                   | SYNLE: Mean number of words before main verb                            |                                                     |                                        |
|                                                                          | DESSC: Sentence count                                               |                                                                   | SYNMEDwrd: Minimal edit distance of words between adjacent sentences    |                                                     |                                        |
|                                                                          | DESSCi: Sentence count incidence per 1000 words                     |                                                                   | SYNMEDlem: Minimal edit distance of lemmas between adjacent sentences   |                                                     |                                        |
|                                                                          | DESWC: Word count (alphanumeric words)                              |                                                                   | SYNMEDpos: Minimal edit distance of POS tags between adjacent sentences |                                                     |                                        |
|                                                                          | DESWCU: Unique word count                                           |                                                                   | SYNCLS1: Ratio of sentences with 1 clause                               |                                                     |                                        |
|                                                                          | DESWCUI: Unique word count incidence per 1000 words                 |                                                                   | SYNCLS2: Ratio of sentences with 2 clauses                              |                                                     |                                        |
|                                                                          | DESPLe: Average paragraph length (sentences per paragraph)          |                                                                   | SYNCLS3: Ratio of sentences with 3 clauses                              |                                                     |                                        |
|                                                                          | DESPStd: Standard deviation of paragraph length                     |                                                                   | SYNCLS4: Ratio of sentences with 4 clauses                              |                                                     |                                        |
|                                                                          | DESSL: Average sentence length (words per sentence)                 |                                                                   | SYNCLS5: Ratio of sentences with 5 clauses                              |                                                     |                                        |
|                                                                          | DESSLd: Standard deviation of sentence length                       | SYNCLS6: Ratio of sentences with 6 clauses                        |                                                                         |                                                     |                                        |
|                                                                          | DESSNSL: Average sentence length excluding stopwords                | SYNCLS7: Ratio of sentences with 7 clauses                        |                                                                         |                                                     |                                        |
|                                                                          | DESSNSLd: Standard deviation of sentence length excluding stopwords | Syntactic Pattern Density Indices                                 | DRNP: Noun phrase density per 1000 words                                |                                                     |                                        |
|                                                                          | DESSLmax: Maximum sentence length                                   |                                                                   | DRNPc: Noun phrase count                                                |                                                     |                                        |
|                                                                          | DESSLmin: Minimum sentence length                                   |                                                                   | DRVp: Verb phrase density per 1000 words                                |                                                     |                                        |
|                                                                          | DESWLsy: Average syllables per word                                 |                                                                   | DRVpc: Verb phrase count                                                |                                                     |                                        |
|                                                                          | DESWLsyd: Standard deviation of syllables per word                  |                                                                   | DRNEG: Negation expression density per 1000 words                       |                                                     |                                        |
|                                                                          | DESCWLSy: Average syllables per content word                        |                                                                   | DRNEGc: Negation expression count                                       |                                                     |                                        |
|                                                                          | DESCWLSyd: Standard deviation of syllables per content word         |                                                                   | DRGER: Gerund form density per 1000 words                               |                                                     |                                        |
|                                                                          | DESCWLLt: Average letters per content word                          |                                                                   | DRGERc: Gerund count                                                    |                                                     |                                        |
| DESCWLLtd: Standard deviation of letters per content word                | DRINF: Infinitive form density per 1000 words                       |                                                                   |                                                                         |                                                     |                                        |
| DESWLLt: Average letters per word                                        | DRINFc: Infinitive count                                            |                                                                   |                                                                         |                                                     |                                        |
| DESWLLtd: Standard deviation of letters per word                         | DRCCONJ: Coordinating conjunction density per 1000 words            |                                                                   |                                                                         |                                                     |                                        |
| DESWNSLIt: Average letters per word (excluding stopwords)                | DRCCONJc: Coordinating conjunction count                            |                                                                   |                                                                         |                                                     |                                        |
| DESWNSLItD: Standard deviation of letters per word (excluding stopwords) | DRSCONJ: Subordinating conjunction density per 1000 words           |                                                                   |                                                                         |                                                     |                                        |
| DESLIt: Average letters per lemma                                        | DRSCONJc: Subordinating conjunction count                           |                                                                   |                                                                         |                                                     |                                        |
| DESLItD: Standard deviation of letters per lemma                         | Connective Indices                                                  | CNCAl: All connectives incidence per 1000 words                   |                                                                         |                                                     |                                        |
| Readability Indices                                                      |                                                                     | RDFHGL: Fernández-Huertas Grade Level                             | CNCaus: Causal connectives incidence per 1000 words                     |                                                     |                                        |
|                                                                          |                                                                     | RDSPPP: Szigriszt-Pazos Perspicuity                               | CNCLogic: Logical connectives incidence per 1000 words                  |                                                     |                                        |
|                                                                          |                                                                     | RDMU: Readability $\mu$ index                                     | CNCADC: Adversative connectives incidence per 1000 words                |                                                     |                                        |
|                                                                          |                                                                     | RDSMOG: SMOG index                                                | CNCTemp: Temporal connectives incidence per 1000 words                  |                                                     |                                        |
|                                                                          |                                                                     | RDFOG: Gunning Fog index                                          | CNCAdd: Additive connectives incidence per 1000 words                   |                                                     |                                        |
|                                                                          |                                                                     | RDHS: Honoré Statistic                                            | Word Information Indices                                                | WRDCONT: Content word incidence per 1000 words      |                                        |
|                                                                          |                                                                     | RDBR: Brunet index                                                |                                                                         | WRDCONTc: Content word count                        |                                        |
|                                                                          |                                                                     | Referential Cohesion Indices                                      |                                                                         | CRFNOI: Noun overlap between adjacent sentences     | WRDNOUN: Noun incidence per 1000 words |
|                                                                          |                                                                     |                                                                   |                                                                         | CRFAOI: Argument overlap between adjacent sentences | WRDNOUNc: Noun count                   |
|                                                                          | CRFSOI: Stem overlap between adjacent sentences                     |                                                                   |                                                                         | WRDVERB: Verb incidence per 1000 words              |                                        |
| CRFCWOI: Content word overlap between adjacent sentences (mean)          | WRDVERBc: Verb count                                                |                                                                   |                                                                         |                                                     |                                        |
| CRFCWOIstd: Content word overlap between adjacent sentences (std dev)    | WRDADJ: Adjective incidence per 1000 words                          |                                                                   |                                                                         |                                                     |                                        |
| CRFANPI: Anaphore overlap between adjacent sentences                     | WRDADJc: Adjective count                                            |                                                                   |                                                                         |                                                     |                                        |
| CRFNOa: Noun overlap between all sentences                               | WRDADV: Adverb incidence per 1000 words                             |                                                                   |                                                                         |                                                     |                                        |
| CRFAOa: Argument overlap between all sentences                           | WRDADVc: Adverb count                                               |                                                                   |                                                                         |                                                     |                                        |
| CRFSOa: Stem overlap between all sentences                               | WRDPRO: Personal pronoun incidence per 1000 words                   |                                                                   |                                                                         |                                                     |                                        |
| CRFCWOa: Content word overlap between all sentences (mean)               | WRDPROc: Personal pronoun count                                     |                                                                   |                                                                         |                                                     |                                        |
| CRFCWOad: Content word overlap between all sentences (std dev)           | WRDPRP1s: First person singular pronoun incidence per 1000 words    |                                                                   |                                                                         |                                                     |                                        |
| CRFANPa: Anaphore overlap between all sentences                          | WRDPRP1sc: First person singular pronoun count                      |                                                                   |                                                                         |                                                     |                                        |
| Lexical Diversity Indices                                                | LDITTRa: Type-token ratio for all words                             | WRDPRP1pc: First person plural pronoun incidence per 1000 words   |                                                                         |                                                     |                                        |
|                                                                          | LDITTRcw: Type-token ratio for content words                        | WRDPRP1pc: First person plural pronoun count                      |                                                                         |                                                     |                                        |
|                                                                          | LDITTRno: Type-token ratio for nouns                                | WRDPRP2s: Second person singular pronoun incidence per 1000 words |                                                                         |                                                     |                                        |
|                                                                          | LDITTRvb: Type-token ratio for verbs                                | WRDPRP2sc: Second person singular pronoun count                   |                                                                         |                                                     |                                        |
|                                                                          | LDITTRadv: Type-token ratio for adverbs                             | WRDPRP2p: Second person plural pronoun incidence per 1000 words   |                                                                         |                                                     |                                        |
|                                                                          | LDITTRadj: Type-token ratio for adjectives                          | WRDPRP2pc: Second person plural pronoun count                     |                                                                         |                                                     |                                        |
|                                                                          | LDITTRLa: Type-token ratio for all lemmas                           | WRDPRP3s: Third person singular pronoun incidence per 1000 words  |                                                                         |                                                     |                                        |
|                                                                          | LDITTRLno: Type-token ratio for noun lemmas                         | WRDPRP3sc: Third person singular pronoun count                    |                                                                         |                                                     |                                        |
|                                                                          | LDITTRLvb: Type-token ratio for verb lemmas                         | WRDPRP3p: Third person plural pronoun incidence per 1000 words    |                                                                         |                                                     |                                        |
|                                                                          | LDITTRLadv: Type-token ratio for adverb lemmas                      | WRDPRP3pc: Third person plural pronoun count                      |                                                                         |                                                     |                                        |
|                                                                          | LDITTRLadj: Type-token ratio for adjective lemmas                   | Psycholinguistic Indices                                          | PSYC: Overall concreteness ratio                                        |                                                     |                                        |
|                                                                          | LDITTRLpron: Type-token ratio for pronouns                          |                                                                   | PSYC0: Very low concreteness ratio (1-2.5)                              |                                                     |                                        |
|                                                                          | LDITTRLpron: Type-token ratio for relative pronouns                 |                                                                   | PSYC1: Low concreteness ratio (2.5-4)                                   |                                                     |                                        |
|                                                                          | LDITTRLipron: Type-token ratio for indefinite pronouns              |                                                                   | PSYC2: Medium concreteness ratio (4.5-5)                                |                                                     |                                        |
|                                                                          | LDITTRLifn: Type-token ratio for functional words                   |                                                                   | PSYC3: High concreteness ratio (5.5-7)                                  |                                                     |                                        |
|                                                                          | LDMLTD: Measure of Textual Lexical Diversity (MTLD)                 |                                                                   | PSYIM: Overall imageability ratio                                       |                                                     |                                        |
|                                                                          | LDVOCd: Vocabulary Complexity Diversity (VoCD)                      |                                                                   | PSYIM0: Very low imageability ratio (1-2.5)                             |                                                     |                                        |
|                                                                          | LDMaas: Maas index                                                  |                                                                   | PSYIM1: Low imageability ratio (2.5-4)                                  |                                                     |                                        |
|                                                                          | LDDno: Noun density                                                 |                                                                   | PSYIM2: Medium imageability ratio (4.5-5)                               |                                                     |                                        |
|                                                                          | LDDvb: Verb density                                                 |                                                                   | PSYIM3: High imageability ratio (5.5-7)                                 |                                                     |                                        |
| LDDadv: Adverb density                                                   | PSYFM: Overall familiarity ratio                                    |                                                                   |                                                                         |                                                     |                                        |
| LDDadj: Adjective density                                                | PSYFM0: Very low familiarity ratio (1-2.5)                          |                                                                   |                                                                         |                                                     |                                        |
| Word Frequency Indices                                                   | WFRChn: Rare noun count                                             | PSYFM1: Low familiarity ratio (2.5-4)                             |                                                                         |                                                     |                                        |
|                                                                          | WFRChoi: Rare noun incidence per 1000 words                         | PSYFM2: Medium familiarity ratio (4.5-5)                          |                                                                         |                                                     |                                        |
|                                                                          | WFRChb: Rare verb count                                             | PSYFM3: High familiarity ratio (5.5-7)                            |                                                                         |                                                     |                                        |
|                                                                          | WFRChvi: Rare verb incidence per 1000 words                         | PSYAoa: Overall age of acquisition ratio                          |                                                                         |                                                     |                                        |
|                                                                          | WFRCadj: Rare adjective count                                       | PSYAoa0: Very early acquisition ratio (1-2.5)                     |                                                                         |                                                     |                                        |
|                                                                          | WFRCadji: Rare adjective incidence per 1000 words                   | PSYAoa1: Early acquisition ratio (2.5-4)                          |                                                                         |                                                     |                                        |
|                                                                          | WFRCadv: Rare adverb count                                          | PSYAoa2: Medium acquisition ratio (4.5-5)                         |                                                                         |                                                     |                                        |
|                                                                          | WFRCadvi: Rare adverb incidence per 1000 words                      | PSYAoa3: Late acquisition ratio (5.5-7)                           |                                                                         |                                                     |                                        |
|                                                                          | WFRCCw: Rare content word count                                     | PSYARO: Overall arousal ratio                                     |                                                                         |                                                     |                                        |
|                                                                          | WFRCCwi: Rare content word incidence per 1000 words                 | PSYAR00: Very low arousal ratio (1-3)                             |                                                                         |                                                     |                                        |
| WFRCCwd: Distinct rare content word count                                | PSYAR01: Low arousal ratio (3-5)                                    |                                                                   |                                                                         |                                                     |                                        |
| WFRCCwdi: Distinct rare content word incidence per 1000 words            | PSYAR02: Medium arousal ratio (5-7)                                 |                                                                   |                                                                         |                                                     |                                        |
| WFMew: Mean frequency of content words                                   | PSYAR03: High arousal ratio (7-9)                                   |                                                                   |                                                                         |                                                     |                                        |
| WFMw: Mean frequency of all words                                        | PSYVAL: Overall valence ratio                                       |                                                                   |                                                                         |                                                     |                                        |
| WFMwv: Mean frequency of rarest words per sentence                       | PSYVAL0: Very negative valence ratio (1-4)                          |                                                                   |                                                                         |                                                     |                                        |
| WFMwcv: Mean frequency of rarest content words per sentence              | PSYVAL1: Negative valence ratio (3-5)                               |                                                                   |                                                                         |                                                     |                                        |
| Semantic Cohesion Indices                                                | SECLoSadj: LSA overlap between adjacent sentences (mean)            | PSYVAL2: Positive valence ratio (5-7)                             |                                                                         |                                                     |                                        |
|                                                                          | SECLoSadjd: LSA overlap between adjacent sentences (std dev)        | PSYVAL3: Very positive valence ratio (7-9)                        |                                                                         |                                                     |                                        |
|                                                                          | SECLoSall: LSA overlap between all sentences (mean)                 | Textual Simplicity Indices                                        | TSSRsh: Ratio of short sentences (<11 words)                            |                                                     |                                        |
|                                                                          | SECLoSalld: LSA overlap between all sentences (std dev)             |                                                                   | TSSRmd: Ratio of medium sentences (11-12 words)                         |                                                     |                                        |
|                                                                          | SECLOPadj: LSA overlap between adjacent paragraphs (mean)           |                                                                   | TSSRlg: Ratio of long sentences (13-14 words)                           |                                                     |                                        |
|                                                                          | SECLOPadjd: LSA overlap between adjacent paragraphs (std dev)       |                                                                   | TSSRxl: Ratio of very long sentences ( $\geq 15$ words)                 |                                                     |                                        |
|                                                                          | SECLoSgiv: LSA overlap between given and new sentences (mean)       |                                                                   |                                                                         |                                                     |                                        |
|                                                                          | SECLoSgivd: LSA overlap between given and new sentences (std dev)   |                                                                   |                                                                         |                                                     |                                        |

Table 6: List of linguistic metrics implemented in PUCP-Metrix

2-label schema and to "basic" vs "advanced" in the 3-label schema.

- Kwiziq<sup>7</sup>. Kwiziq is an online language-learning platform that offers graded Spanish readings labeled with CEFR levels. We use the available data proposed by Vázquez-Rodríguez et al. (2022) and map the CEFR annotations to our 2- and 3-label classification schemes using the same criteria.
- HablaCultura. This dataset comprises educational readings sourced from the HablaCultura platform<sup>8</sup>, where each text is labeled by instructors with CEFR levels. We use the same level mappings used by Vázquez-Rodríguez et al. (2022).

## B.2 Machine-generated Text Detection

Human-generated texts in AuTextTification were sourced from publicly available datasets, including MultiEURLEX (Chalkidis et al., 2021) (legal), XLSUM/MLSUM (Hasan et al., 2021; Scialom et al., 2020) (news), COAR/COAH (Molina-González et al., 2014) (reviews), XLM-Tweets (Barbieri et al., 2022) and TSD (Leis et al., 2019) (tweets), and WikiLingua (Ladhak et al., 2020) (how-to articles). Machine-generated texts were produced using six large language models: three from the BLOOM family (BLOOM-1B<sup>9</sup>, BLOOM-3B<sup>10</sup>, BLOOM-7B1<sup>11</sup>) and three from the GPT-3 family (babbage, curie, text-davinci-003).

## C Feature Analysis

We applied Anova over our dataset using all the metrics. We set a p-value of 0.05 and remove the features that do not make contribution for our analysis.

Figure 1 shows a heatmap representing the coverage of linguistic categories along the ranking, i.e., the distribution of linguistic features as more signals are included. Overall, the contribution of linguistic features varies across tasks. For machine-generated content detection, top-ranked signals are dominated by word frequency, readability, and semantic cohesion metrics. In contrast, descriptive

and connective metrics play a more limited role and appear only at later ranks.

For ARA tasks, the importance shifts toward descriptive features, syntactic pattern density, readability, syntactic complexity, and textual simplicity metrics. Conversely, semantic cohesion and connective metrics are comparatively less important.

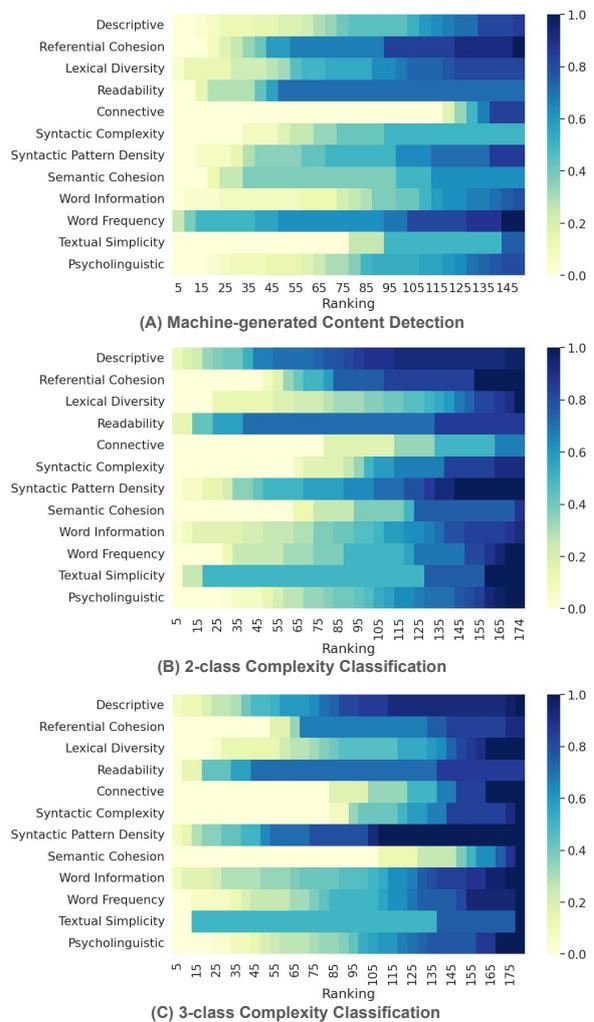


Figure 1: Category coverage along the ranking for PUCP-Metrix

<sup>7</sup>The platform is available at <https://www.kwiziq.com/>

<sup>8</sup>Available at <https://hablacultura.com/>

<sup>9</sup>Available at <https://huggingface.co/bigscience/bloom-1b7>.

<sup>10</sup>Available at <https://huggingface.co/bigscience/bloom-3b>.

<sup>11</sup>Available at <https://huggingface.co/bigscience/bloom-7b1>.

# INTEGRITYSHIELD

## A System for Ethical AI Use and Authorship Transparency in Assessments

Ashish Raj Shekhar\* Shiven Agarwal\* Priyanuj Bordoloi Yash Shah  
Tejas Anvekar Vivek Gupta

Arizona State University

[Project Page](#) [Demo](#) [Video](#) [Code](#)

{ashekha9, sagar147, pbordoloi, yshah124, tanvekar, vgupt140}@asu.edu

### Abstract

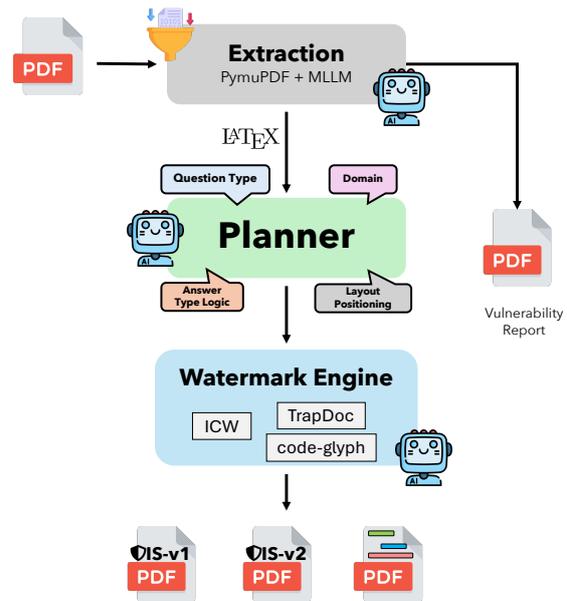
Multimodal Large Language Models (MLLMs) can now solve entire exams directly from uploaded PDF assessments, raising urgent concerns about academic integrity and the reliability of grades and credentials. Existing watermarking techniques either operate at the token level or assume control over the model’s decoding process, making them ineffective when students query proprietary black-box systems using instructor-provided documents. We present INTEGRITYSHIELD, a document-layer watermarking system that embeds schema-aware, item-level watermarks into assessment PDFs while keeping their human-visible appearance unchanged. These watermarks consistently prevent MLLMs from answering shielded exam PDFs and encode stable, item-level signatures that can be reliably recovered from model or student responses. Across 30 question papers spanning STEM, humanities, and medical reasoning, INTEGRITYSHIELD achieves exceptionally high prevention (91-94% exam-level blocking) and strong detection reliability (89-93% signature retrieval) across four commercial MLLMs. Our demo showcases an interactive interface where instructors upload an exam, preview watermark behavior, and inspect pre/post AI performance and authorship evidence.

## 1 Introduction

LLMs and MLLMs can now interpret full PDF assessments, reason over diagrams and tables, and produce fluent step-by-step solutions within seconds. While these capabilities expand access to high-quality assistance, they simultaneously undermine the credibility of homework and online exams by enabling students to outsource entire assessments to AI tools (OpenAI, 2023; Team, 2024; Susnjak, 2022).

Institutions have responded with post-hoc detection (e.g., authorship classifiers (Emi and Spero,

\*contributed equally



Shielded PDFs & Attribution Report

Figure 1: Overview of INTEGRITYSHIELD. The system extracts question structure from an assessment PDF, uses an LLM-based planner to select schema-aware watermarking tactics, and applies document-layer perturbations through the watermark engine. It outputs shielded PDF variants (VIS-v1, VIS-v2) and an attribution report summarizing AI vulnerability along with authorship signals.

2024; Thai et al., 2025)) and surveillance-heavy proctoring (e.g., keystroke, browser, or gaze monitoring (Atoum et al., 2017; Kundu et al., 2024)). However, detectors struggle with short answers, code, paraphrasing, and mixed authorship (Mitchell et al., 2023; Niu et al., 2024), while invasive monitoring raises significant privacy, accessibility, and equity concerns.

These approaches share a fundamental limitation: they analyze the student’s output. In practice, the dominant workflow is the opposite students upload instructor-provided PDFs to black-box AI systems. Existing watermarking methods, which modify generation at the model’s decoder (Kirchen-

bauer et al., 2023; Liu et al., 2025), cannot operate in this setting. This motivates an interesting question: *Can assessments themselves be instrumented so that AI reliance becomes observable, without altering visible exam content or student workflows?*

**From detection to document-level watermarking.** We exploit the render-parse gap in PDFs: what humans see often differs from what AI parsers ingest. By injecting invisible text, glyph remappings, and lightweight overlays, we influence model interpretation while leaving the exam visually unchanged. **INTEGRITYSHIELD** operationalizes this idea as an authorship-aware watermarking system. Rather than asking whether a student cheated, we ask: *to what extent do model-generated responses follow a consistent watermark signature embedded in the exam?* This reframing provides instructors with an interpretable notion of authorship degree while maintaining fairness for honest students. Finally, we summarize our contributions as:

- We introduce **INTEGRITYSHIELD**, a document-layer watermarking system that embeds schema-aware watermarks into assessment PDFs while keeping their human-visible appearance unchanged.
- We develop an LLM-driven planner and PDF watermark engine that adapts tactics to question type, achieving consistently high prevention (91-94% exam-level blocking) and strong detection reliability (89-93% retrieval) across four commercial MLLMs on a thirty-exam benchmark.
- We release an interactive demo that allows instructors to upload exams, preview watermarks, and inspect pre/post AI performance and authorship evidence, enabling ethical and transparent AI use in education.

## 2 Background and Related Work

**AI assistance and mixed authorship.** LLMs increasingly participate in writing and problem-solving tasks, often producing blended human-AI content. Recent work formalizes this as *homogeneous* vs. *heterogeneous* mixed authorship (Thai et al., 2025). Existing detectors including perplexity-based methods (Mitchell et al., 2023), style-based classifiers (Emi and Spero, 2024), and multilingual cheating detectors (Niu et al.,

2024) struggle with short answers, paraphrasing, and multi-author mixtures, and they analyze only the *output*, leaving the assessment itself uninstrumented.

**AI watermarking.** Watermarking embeds provenance signals into generated content, typically by modifying decoding distributions (Kirchenbauer et al., 2023) or via prompt-based in-context cues (Liu et al., 2025). Parallel work explores invisible watermarks for AI-generated writing, designed to survive paraphrasing and editing (Liu and Bu, 2024). These methods assume control over generation, which is infeasible when students query proprietary black-box systems using instructor-provided PDFs.

**Document-layer perturbations.** Recent work shows that perturbing the PDF substrate - via phantom tokens, font CMaps, or off-page text - can induce systematic model errors without affecting human readability (Jin et al., 2025; Xiong et al., 2025). Our work builds on these insights but shifts the objective: rather than deceiving models or detecting cheating, we embed *recoverable watermark signatures* that quantify the extent of AI involvement in solving an exam.

## 3 **INTEGRITYSHIELD** System Architecture and Workflow

**INTEGRITYSHIELD** is designed as a practical tool for instructors who want to harden PDF-based assessments against AI assistants without redesigning their exams or changing grading workflows. The system keeps all human-facing content (layout, typography, pagination, item numbering) unchanged while embedding signals that reliably influence model-side parsing. It adapts watermark tactics to the item schema, treating **MCQ**, **true/false**, and **Long-Form** questions differently, remaining robust across black-box MLLMs; and exposes a lightweight web interface where instructors can upload assessments, preview watermark behavior, and inspect calibration and authorship signals with minimal configuration.

### 3.1 End-to-End Architecture and Workflow

Figure 1 summarizes the end-to-end workflow of **INTEGRITYSHIELD**. A single-page web front end communicates with a stateless backend that operates directly on the PDF substrate. The backend is organized around four logical services: document ingestion, which parses the uploaded PDF

into a structured item schema with stems, options, diagrams, and answer keys; strategy planning, via a lightweight instruction-tuned LLM that assigns a watermark plan to each item based on its type, gold answer, and local layout metadata; PDF rewriting, which applies the plan while enforcing that the rendered appearance of the document remains unchanged; and an authorship and calibration service that runs reference models on original and watermarked PDFs and later scores submitted answers against stored watermark signatures. From an instructor’s perspective, this architecture is depicted in a three-stage interaction flow.

### Stage 1: Upload and Watermark Planning.

The instructor uploads an exam PDF through the browser. The ingestion service segments pages into questions, detects answer options and numbering, and associates inline figures or tables with the relevant items, producing a compact item schema with content spans, page coordinates, and answer keys when available, as illustrated in Figure 2.

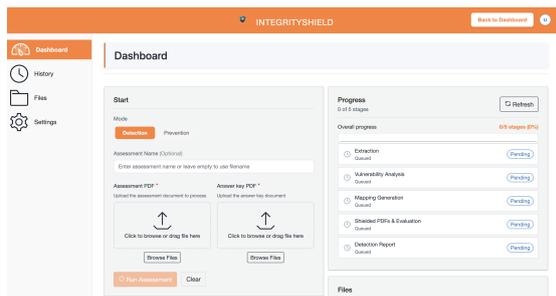


Figure 2: Stage 1: Upload and Watermark Planning. Instructors upload an assessment PDF and answer key, after which the system extracts question structure and previews the planned schema-aware watermarking strategies.

The strategy planner then assigns, for each item, either a *target distractor* (for multiple-choice and true/false questions) or a small set of signature key phrases (for long-form questions) and decides which document-layer mechanisms to apply. The interface presents a split-screen preview of original and watermarked pages with per-question summaries of the chosen strategy, allowing instructors to inspect and optionally disable aggressive tactics (such as strong glyph remapping) before proceeding.

### Stage 2: Watermark Embedding and AI Calibration.

Once the plan is confirmed, the PDF rewriting service applies it directly to the assessment file. It injects invisible text spans anchored near stems and options, applies CMap-based glyph

remapping so that visually identical tokens are parsed differently by models, and, when appropriate, adds clipped or off-page overlays that insert distractor-oriented cues into the parsable stream while keeping them outside the visible canvas, as shown in Figure 3.

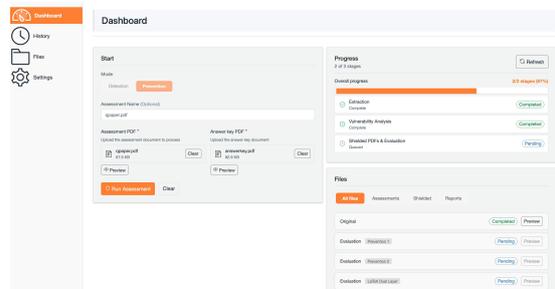


Figure 3: Stage 2: Watermark Embedding and AI Calibration. After planning, the system applies document-layer watermarks to the assessment PDF and evaluates original vs. watermarked versions against multiple MLLMs to generate prevention and detection reports.

We instantiate two watermark configurations, **VIS-v1** and **VIS-v2**, which differ in the density and combination of these mechanisms: **VIS-v1** uses a lighter mix of hidden-text and minimal glyph remapping, whereas **VIS-v2** employs stronger multi-layer perturbations for maximal robustness across parsing pipelines.

After rewriting, the system verifies that the rendered appearance of the PDF matches the original across common viewers (Adobe Reader, Chrome, macOS Preview). In the same stage, the authorship and calibration service evaluates both the original and watermarked versions with a panel of reference models in a simulated “*student uploads the exam*” setting, computing pre- versus post-watermark accuracy, the fraction of incorrect answers that land on intended distractors, and per-item watermark retrieval rates. An interactive report summarizes these statistics and assists instructors in selecting an appropriate watermark “*strength*” preset.

### Stage 3: Authorship Analysis.

After an assessment has been protected with our **VIS** watermarked PDFs, instructors can use it to analyze responses. The interface accepts either raw model outputs (for research) or anonymized student responses exported from a learning management system, depicted in Figure 4.

For each question, the authorship engine checks whether the response follows the stored watermark signature: for objective questions, this reduces to matching the target distractor (or a small tied set);

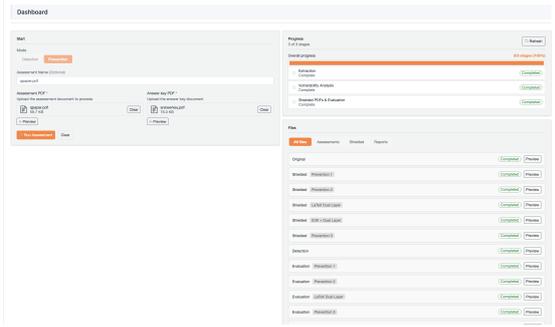


Figure 4: Stage 3: Authorship Analysis. The dashboard displays per-question watermark retrieval, exam-level authorship scores, and previewable shielded PDFs, enabling instructors to inspect AI-reliance signals.

for long-form questions, a judge LLM scores how closely the response tracks the watermark’s key phrases or erroneous reasoning patterns. These per-item scores are aggregated into an exam-level authorship degree and displayed on a dashboard with cohort-level distributions and drill-down views for individual questions. The tool is explicitly positioned as an aid for triage rather than an automatic decision system: high authorship scores are intended to trigger follow-up actions such as brief oral checks or additional written assessments, keeping human judgment in the loop.

## 4 Experiments

**Models and prompts.** We evaluate **INTEGRITYSHIELD** against a panel of four proprietary frontier MLLMs that support direct PDF ingestion: GPT-5, Claude Sonnet-4.5, Grok-4.1, and Gemini-2.5 Flash. All models are treated as black boxes and accessed via their official APIs, with temperature set to 0 and maximum output length sufficient to cover all questions in an exam. For each exam, we use a minimal, instruction-style prompt that (i) asks the model to answer all questions in order, (ii) returns a structured list of answers (e.g., “*Q1*: A, *Q2*: C, ...” for MCQ and T/F; numbered paragraphs for Long-Form (LF)), and (iii) forbids external tools or web browsing. We use the same prompting templates for original and watermarked PDFs; full prompt text for MCQ, T/F, and LF questions appears in Appendix A as **Prompt A**, **Prompt B**, and **Prompt C**.

**Baselines.** We compare our two watermark configurations, **VIS-v1** and **VIS-v2**, against three document- or prompt-level baselines. **ICW** is an in-context watermarking method that attempts to steer

model outputs using prompt-side patterns without modifying the PDF content using an invisible white color small sized font size (0.1-0.5). (Liu et al., 2025). **code-glyph** is a document-layer baseline that manipulates bitcode to glyph mapping on question text to perturb parsing while keeping human readability intact (Xiong et al., 2025). **TRAPDOC** adapts document-layer perturbations that introduce phantom tokens and layout tricks to cause models to produce plausible but incorrect answers without visible changes to the PDF (Jin et al., 2025). In contrast, **VIS-v1** and **VIS-v2** operate directly at the PDF substrate with schema-aware hidden text, glyph remapping, and overlays; **VIS-v1** applies a lighter combination aimed at minimal perturbation, while **VIS-v2** uses denser, multi-layer perturbations for maximal robustness.

**Benchmark Dataset.** To approximate real assessment settings, we compile a diverse benchmark of exam-style PDFs by web-scraping publicly available quizzes, homework sets, and midterm assessments from university course websites (e.g., Stanford and other institutions). The collected material spans mathematics, science, programming, humanities, and medical reasoning, and includes a mix of MCQ, T/F, and long-form questions. From this pool, we sample  $\approx 10\%$  of items to construct our benchmark, selecting documents that (i) contain at least ten questions, (ii) include at least three question formats, and (iii) render cleanly as PDFs. All items and answer keys are qualitatively reviewed by two authors and spot-validated quantitatively (e.g., via official solutions when available) to filter out ambiguous or mislabeled questions.

**Evaluation metrics.** We evaluate **INTEGRITYSHIELD** along two complementary axes: *prevention*, which measures how strongly watermarking disrupts a model’s ability to answer correctly, and *detection*, which captures how reliably watermark signatures can be recovered from model or student responses.

For **prevention**, we simply check whether watermarking causes the model to fail or refuse to answer the exam PDF. For exam  $d$ ,  $\text{Prev}(f, d) = 1[y^{\text{wm}}$  contains no usable answers], and we report the percentage of PDFs where this occurs.

For **detection**, we measure the degree to which model outputs follow the embedded watermark signature. For each item, the authorship engine assigns a retrieval score  $c_i \in [0, 1]$ : for MCQ; T/F,  $c_i = 1$  iff the model selects the target distractor;

for long-form,  $c_i$  is produced by a judge LLM evaluating alignment with watermark keyphrases. The exam-level detection score is

$$\text{Det}(d', y) = \frac{1}{n_{d'}} \sum_{i=1}^{n_{d'}} c_i,$$

representing the proportion of responses that exhibit watermark-consistent behavior. We report detection scores per model and method, with breakdowns by question type.

#### 4.1 Quantitative Analysis on Prevention and Detection

Table 1 summarizes prevention and detection performance for all baselines and our **IS** variants across GPT, Claude, Grok, and Gemini.

| Method                | GPT          | Claude       | Grok         | Gemini       |
|-----------------------|--------------|--------------|--------------|--------------|
| <i>Prevention-ASR</i> |              |              |              |              |
| ICW                   | 07.20        | 05.80        | 04.10        | 03.50        |
| code-glyph            | 86.30        | 84.7         | 83.20        | 81.90        |
| TRAPDOC               | 88.70        | 86.40        | 85.10        | 40.50        |
| <b>IS-v1</b>          | 91.20        | 90.80        | 90.50        | 90.10        |
| <b>IS-v2</b>          | <b>93.60</b> | <b>92.90</b> | <b>92.30</b> | <b>91.70</b> |
| <i>Detection</i>      |              |              |              |              |
| ICW                   | 06.80        | 05.30        | 04.60        | 03.20        |
| code-glyph            | 85.90        | 84.20        | 82.70        | 81.40        |
| TRAPDOC               | 87.90        | 85.80        | 84.60        | 43.40        |
| <b>IS-v1</b>          | 90.70        | 90.30        | 89.90        | 89.50        |
| <b>IS-v2</b>          | <b>92.80</b> | <b>92.10</b> | <b>91.60</b> | <b>91.00</b> |

Table 1: Prevention and detection performance across models. Prevention-ASR is the percentage of exam PDFs on which watermarking causes the model to refuse or fail to produce usable answers. Detection is the percentage of questions whose responses follow the embedded watermark signature. For both metrics, higher is better. ICW: in-context watermarking; code-glyph: glyph perturbation; TRAPDOC: phantom-token PDF attack; IS: **IS** variants.

In the *Prevention-ASR* block, ICW almost never prevents models from answering full exams, with single-digit prevention rates across all models. This confirms that prompt-only steering is ineffective when students upload raw PDFs to black-box MLLMs. Document-layer baselines such as code-glyph and TRAPDOC are substantially stronger on GPT, Claude, and Grok (around 83-89% prevention), but TRAPDOC degrades sharply on Gemini (40.5%), suggesting that its perturbations do not transfer reliably across parsing and model stacks. By contrast, **IS-v1** and **IS-v2** achieve consistently high prevention on *all* models (90-94%), indicating that schema-aware, multi-

layer PDF watermarking can robustly block end-to-end exam solving for contemporary MLLMs.

The *Detection* block shows a similar pattern. ICW again yields negligible detection rates, while code-glyph and TRAPDOC achieve strong detection on GPT, Claude, and Grok (mid-80s); however, TRAPDOC drops to 43.4% on Gemini. In contrast, **IS-v1** and especially **IS-v2** maintain high detection performance across all four models (around 89-93%), meaning that whenever models do attempt to answer on watermarked exams, their outputs follow the embedded watermark signatures in a highly consistent way, enabling reliable authorship attribution.

#### 4.2 **IS** Performance for Question-Category

Table 2 breaks down the impact of **IS** on answer accuracy by question type (MCQ, T/F, LF) and model, comparing performance on original (*w/o*) and watermarked (*w/*) exams. Without watermarking, all four MLLMs attain very high accuracy across categories (typically 94-97%), reflecting their strong baseline performance on our exam-style benchmark.

| Type | GPT        |            | Claude     |            | Grok       |            | Gemini     |            |
|------|------------|------------|------------|------------|------------|------------|------------|------------|
|      | <i>w/o</i> | <i>w/</i>  | <i>w/o</i> | <i>w/</i>  | <i>w/o</i> | <i>w/</i>  | <i>w/o</i> | <i>w/</i>  |
| MCQ  | 96.2       | <b>7.8</b> | 95.8       | <b>6.9</b> | 94.9       | <b>5.7</b> | 94.1       | <b>4.3</b> |
| T/F  | 95.7       | <b>6.5</b> | 95.3       | <b>5.8</b> | 94.6       | <b>4.9</b> | 93.8       | <b>3.6</b> |
| LF   | 96.8       | <b>5.2</b> | 96.4       | <b>4.6</b> | 95.3       | <b>3.8</b> | 94.7       | <b>3.1</b> |

Table 2: Per-question-type answer accuracy without (*w/o*) and with (*w/*) our **IS** watermarks. Values show the residual accuracy of each model on shielded exams; lower *w/* accuracy indicates stronger prevention for that question type.

With **IS** enabled, residual accuracy collapses into the low single digits for every model and question type (3-8%), corresponding to an 85-90 point drop. Long-form questions show the largest reductions for GPT and Claude, while MCQ and T/F items are also heavily disrupted across all models. These results indicate that our document-layer watermarks are effective not only at the exam level, but also uniformly across different assessment formats.

#### 4.3 Utility of **IS**

Beyond aggregate metrics, **IS** provides instructors with actionable, item-level evidence of AI reliance. Table 3 illustrates this with a qualitative example: on an OSI-model question, all baseline attacks (ICW, code-glyph, TRAPDOC)

| Attack Method | GPT | Claude | Grok | Gemini |
|---------------|-----|--------|------|--------|
| ICW           | A   | A      | A    | A      |
| code-glyph    | A   | A      | A    | A      |
| TRAPDOC       | A   | A      | A    | A      |
| 🛡️IS-v1       | B   | B      | B    | B      |
| 🛡️IS-v2       | C   | C      | C    | C      |

Table 3: Model predictions across attack methods for the OSI model question. *Q: Which layer of the OSI model is responsible for routing packets between networks?*  
**Gold Answer: A**

collapse to the same incorrect prediction across models, offering no consistent signal for attribution. In contrast, our schema-aware variants (🛡️IS-v1, 🛡️IS-v2) drive models toward distinct, watermark-aligned distractors (B and C, respectively), enabling clear and separable authorship signatures.

In a *prevention-focused* deployment, the system summarizes where watermarking fully blocks a model from answering an exam, providing a document-level view of which assessments are resilient to AI-based shortcuts. In a *detection-focused* deployment, the system aggregates authorship evidence across questions, showing, for example, that “Q3 follows the 🛡️IS-v2 signature across multiple models”.

These reports are intended as triage tools: instructors can identify items likely influenced by AI, perform brief oral checks or follow-up tasks, and intervene proportionally. By surfacing interpretable authorship signals rather than relying on opaque classifiers or intrusive proctoring, 🛡️INTEGRITYSHIELD enables ethical, transparent, and governance-aligned AI use in educational assessments.

## 5 Conclusion

🛡️INTEGRITYSHIELD A System for Ethical AI Use & Authorship Transparency in Assessments, demonstrates that assessment integrity can be strengthened without invasive monitoring by instrumenting the exam document itself. By operating directly at the PDF substrate, our system embeds schema-aware watermarks that both (i) prevent modern MLLMs from answering shielded exams (91-94% exam-level blocking) and (ii) yield stable, recoverable authorship signatures (89-93% retrieval) when AI is used. These effects hold consistently across question types and four commercial MLLMs, highlighting the robustness of document-layer watermarking as a practical defense.

The demo showcases a complete workflow for

real instructional use: uploading an exam, previewing watermark strategies, generating shielded variants, running automated AI calibration, and inspecting item-level authorship evidence. This combination of prevention and attribution provides instructors and institutions with actionable, interpretable signals supporting fair assessment practices, targeted follow-up, and transparent communication with students.

We hope 🛡️INTEGRITYSHIELD serves as a step toward ethically grounded AI governance in education, enabling institutions to observe AI reliance without resorting to surveillance or restricting access to assistive technologies.

## Limitations

Our evaluation is limited to a thirty-exam benchmark, a fixed set of frontier MLLMs, and simulated usage in which models directly consume instructor PDFs. Real-world deployments may involve broader variation in domains, languages, accessibility workflows (e.g., screen readers), and institution-specific formats. As MLLMs and their PDF-parsing pipelines evolve, watermark robustness may drift, necessitating periodic recalibration.

🛡️INTEGRITYSHIELD is not a definitive detector of misconduct. Authorship scores indicate alignment with watermark signatures not whether a student violated policy and should be used as a triage signal for human follow-up (e.g., brief oral checks), not as automatic evidence for sanctions. We acknowledge that we did not evaluate robustness against rasterization, ghostscript re-encoding, Word export/import, or PDF sanitization tools. A technically sophisticated adversary who transforms the document first falls outside our current scope. Students who instead choose to retype questions, take screenshots, or manually transcribe content introduce meaningful friction that itself serves as a partial deterrent. More importantly, such workarounds would be inconsistent with the typical behavior of students seeking a quick, low-effort solution, which is the primary risk profile the system addresses.

The system is designed for deployability under current parsing pipelines, with the expectation that watermark strategies will need to be updated over time; analogous to how anti-virus signatures or plagiarism detection tools require periodic updates.

Finally, our approach assumes institutional control over assessment PDFs. Similar watermarking

techniques could be misapplied to non-assessment documents, so we explicitly restrict the intended use to formal educational settings with clear governance, transparency, and AI-use policies. PDF is the dominant format for distributed assessments in higher education, which motivated our focus. Extension to other formats (e.g., Images, HTML-based assessments) is an important direction for future work.

## Ethics Statement

This work aims to support ethical and transparent AI use in educational assessment settings. **INTEGRITYSHIELD** operates exclusively on instructor-provided documents and does not monitor students, avoiding surveillance-heavy practices such as keystroke logging, webcam tracking, or device control. The system is designed to keep all responses and analyzes within institutional infrastructure, respecting student privacy and data-governance requirements.

Authorship scores produced by our watermarking framework indicate alignment with embedded watermark signatures; they do *not* constitute evidence of misconduct. We recommend that institutions (i) clearly communicate AI-use policies and the presence of watermarking to students, (ii) treat high authorship scores only as signals for human review (e.g., follow-up questions or oral checks), and (iii) ensure that any use of these signals aligns with local policies, academic integrity guidelines, and privacy regulations. The system is designed as a practical deterrent and authorship signal for institutional triage, not as a cryptographically secure system. We also note that IS-v2's multi-layer approach (CMap remapping + overlays) is harder to isolate than pure hidden-text injection, though we do not claim this is detection-proof.

All experiments were conducted with fixed model parameters (e.g., temperature,  $top_p$ ,  $top_k$ ) to mitigate stochastic variability in black-box LLMs. Models used in this work (e.g., GPT-5, Gemini-2.5 Flash, Grok-4.1, Claude Sonnet-4.5) were accessed in accordance with their respective usage policies. Data labeling and verification were performed by author-annotators, and AI-based tools (e.g., Grammarly, ChatGPT) were used strictly for language refinement. To the best of our knowledge, this study introduces no additional ethical risks beyond those common to LLM evaluation in controlled educational settings.

## Acknowledgements

We thank the Complex Data Analysis and Reasoning Lab at Arizona State University for computational support, and Sandipan De and Eun Woo Im for their helpful reviews of earlier drafts of this work. We are also grateful to Ben Zhou, Pavan Turaga, and Janice Mak for their valuable feedback and suggestions, and to the Arizona State University academic integrity office for their support and feedback.

## References

- Yousef Atoum, Liping Chen, Alex X. Liu, Stephen D. H. Hsu, and Xiaoming Liu. 2017. [Automated Online Exam Proctoring](#). *IEEE Transactions on Multimedia*, 19(7):1609–1624.
- Bradley Emi and Max Spero. 2024. [Technical Report on the Pangram AI-Generated Text Classifier](#). Technical report, Pangram Labs.
- Hyundong Jin, Sicheol Sung, Shinwoo Park, SeungYeop Baik, and Yo-Sub Han. 2025. TRAPDOC: Deceiving LLM Users by Injecting Imperceptible Phantom Tokens into Documents. *EMNLP Findings*. ArXiv:2506.00089.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, and 1 others. 2023. A Watermark for Large Language Models. *arXiv preprint arXiv:2301.10226*.
- Debnath Kundu, Atharva Mehta, Rajesh Kumar, Naman Lal, Avinash Anand, Apoorv Singh, and Rajiv Ratn Shah. 2024. [Keystroke Dynamics Against Academic Dishonesty in the Age of LLMs](#). In *Proceedings of the IEEE International Joint Conference on Biometrics (IJCB)*.
- Yepeng Liu and Yuheng Bu. 2024. [Adaptive Text Watermark for Large Language Models](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 30718–30737. PMLR.
- Yepeng Liu, Xuandong Zhao, Christopher Kruegel, Dawn Song, and Yuheng Bu. 2025. In-Context Watermarks for Large Language Models. *arXiv preprint arXiv:2505.16934*.
- Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. DetectGPT: Zero-shot Machine-Generated Text Detection Using Probability Curvature. *Proceedings of the 40th International Conference on Machine Learning*.
- Chenhao Niu, Kevin P. Yancey, Ruidong Liu, Mirza Basim Baig, André Kenji Horie, and James Sharpnack. 2024. Detecting LLM-Assisted Cheating on Open-Ended Writing Tasks on Language Proficiency Tests. *EMNLP Industry Track*.

OpenAI. 2023. [ChatGPT](#). OpenAI blog.

Teo Susnjak. 2022. [ChatGPT: The End of Online Exam Integrity?](#) *Preprint*, arXiv:2212.09292.

Gemini Team. 2024. Gemini 1.5: Unlocking Multi-modal Understanding Across Millions of Tokens of Context. *arXiv preprint arXiv:2403.05530*.

Katherine Thai, Bradley Emi, Elyas Masrou, and Mohit Iyyer. 2025. [EditLens: Quantifying the Extent of AI Editing in Text](#). *Preprint*, arXiv:2510.03154.

Junjie Xiong, Changjia Zhu, Shuhang Lin, Chong Zhang, Yongfeng Zhang, Yao Liu, and Lingyao Li. 2025. [Invisible Prompts, Visible Threats: Malicious Font Injection in External Resources for Large Language Models](#). *Preprint*, arXiv:2505.16957.

## **A Prompts Details**

## Prompt A: MCQ Perturbation

You are an expert at generating text substitutions for academic multiple-choice questions.

Given:

```
- LaTeX code for the question stem: {latex_stem_text}
- Gold answer: {gold_answer}
- Question type: {question_type}
- Options: {options}
- Strategy: replacement
- Reasoning steps:
{reasoning_steps}
- Copyable text (use this exact text when selecting substrings):
<<<COPY
{copyable_text}
>>>
{prefix_note}{answer_guidance}{retry_instructions}
```

Your task:

Generate {k} valid mappings that satisfy the replacement strategy. Each mapping should:

1. Replace one contiguous substring of the question stem
2. Change the answer from the gold option ({gold\_answer}) to a different incorrect option
3. Ensure the replacement is semantically meaningful and natural
4. Cause a verifiable deviation in the answer

For each mapping, provide:

1. question\_index: The question number ({question\_index})
2. latex\_stem\_text: Exact LaTeX text of the question stem (must match the input exactly)
3. original\_substring: The substring to replace (must be a contiguous substring of latex\_stem\_text)
4. replacement\_substring: The replacement text
5. start\_pos: Start position of original\_substring relative to latex\_stem\_text (0-based index)
6. end\_pos: End position of original\_substring relative to latex\_stem\_text (exclusive, 0-based index)
7. target\_wrong\_answer: The target incorrect option label (e.g., "B", "C", "D")
8. reasoning: Brief explanation of why this mapping satisfies the strategy

IMPORTANT:

- The original\_substring MUST be an exact substring of latex\_stem\_text
- The start\_pos and end\_pos MUST be accurate ( $\text{start\_pos} + \text{len}(\text{original\_substring}) = \text{end\_pos}$ )
- The target\_wrong\_answer MUST be different from the gold answer
- CRITICAL: The replacement\_substring MUST be DIFFERENT from the original\_substring. Do NOT generate mappings where
  - ↪ original\_substring == replacement\_substring (e.g., "power" --> "power" is INVALID). The replacement MUST change the
  - ↪ text to create actual manipulation.
- CRITICAL: Neither original\_substring nor replacement\_substring can be empty strings. Both must contain actual text.
- LENGTH CONSTRAINT: The replacement\_substring MUST be smaller or equal in length to the original\_substring ( $\text{len}(\text{replacement\_substring}) \leq \text{len}(\text{original\_substring})$ ). This is critical for maintaining document layout and preventing
  - ↪ text overflow.
- latex\_stem\_text is provided exactly as it appears in the LaTeX source. Do NOT trim, normalise, or reformat it when
  - ↪ determining positions.
- The latex\_stem\_text may include \item tokens from enumerate environments. Keep the \item token intact and operate on the
  - ↪ descriptive text that follows it whenever possible.
- The replacement should be natural and semantically meaningful

Return as JSON array:

```
[
  {{
    "question_index": {question_index},
    "latex_stem_text": "...",
    "original_substring": "...",
    "replacement_substring": "...",
    "start_pos": 0,
    "end_pos": 5,
    "target_wrong_answer": "B",
    "reasoning": "..."
  }},
  ...
]
```

Return ONLY valid JSON, no markdown or additional text.

## Prompt B: True False Perturbation

You are an expert at generating text substitutions for True/False questions.

Given:

- LaTeX code for the question stem: {latex\_stem\_text}
- Gold answer: {gold\_answer}
- Question type: {question\_type}
- Strategy: replacement
- Reasoning steps: {reasoning\_steps}
- Copyable text (use this exact text when selecting substrings):  
<<<COPY  
{copyable\_text}  
>>>  
{prefix\_note}{answer\_guidance}{retry\_instructions}

Your task:

Generate {k} valid mappings that satisfy the replacement strategy. Each mapping should:

1. Replace one contiguous substring of the question stem
2. Flip the answer from {gold\_answer} to the opposite answer
3. Ensure the replacement is semantically meaningful and natural
4. Cause a verifiable deviation in the answer

For each mapping, provide:

1. question\_index: The question number ({question\_index})
2. latex\_stem\_text: Exact LaTeX text of the question stem (must match the input exactly)
3. original\_substring: The substring to replace (must be a contiguous substring of latex\_stem\_text)
4. replacement\_substring: The replacement text
5. start\_pos: Start position of original\_substring relative to latex\_stem\_text (0-based index)
6. end\_pos: End position of original\_substring relative to latex\_stem\_text (exclusive, 0-based index)
7. target\_wrong\_answer: The opposite answer (e.g., "False" if gold is "True", or "True" if gold is "False")
8. reasoning: Brief explanation of why this mapping satisfies the strategy

IMPORTANT:

- The original\_substring MUST be an exact substring of latex\_stem\_text
- The start\_pos and end\_pos MUST be accurate ( $\text{start\_pos} + \text{len}(\text{original\_substring}) = \text{end\_pos}$ )
- The target\_wrong\_answer MUST be the opposite of the gold answer
- CRITICAL: The replacement\_substring MUST be DIFFERENT from the original\_substring. Do NOT generate mappings where
  - ↪ original\_substring == replacement\_substring (e.g., "force" --> "force" is INVALID). The replacement MUST change the text to create actual manipulation.
- CRITICAL: Neither original\_substring nor replacement\_substring can be empty strings. Both must contain actual text.
- LENGTH CONSTRAINT: The replacement\_substring MUST be smaller or equal in length to the original\_substring ( $\text{len}(\text{replacement\_substring}) \leq \text{len}(\text{original\_substring})$ ). This is critical for maintaining document layout and preventing text overflow.
- latex\_stem\_text is provided exactly as it appears in the LaTeX source. Do NOT trim, normalise, or reformat it when
  - ↪ determining positions.
- The latex\_stem\_text may include \item tokens from enumerate environments. Keep the \item token intact and operate on the
  - ↪ descriptive text that follows it whenever possible.
- The replacement should be natural and semantically meaningful

Return as JSON array:

```
[
  {
    "question_index": {question_index},
    "latex_stem_text": "...",
    "original_substring": "...",
    "replacement_substring": "...",
    "start_pos": 0,
    "end_pos": 5,
    "target_wrong_answer": "False",
    "reasoning": "..."
  },
  ...
]
```

Return ONLY valid JSON, no markdown or additional text.

## Prompt C: LongForm Perturbation

You are an expert at generating text substitutions for long-form questions (essay, short answer, etc.).

Given:

- LaTeX code for the question stem: {latex\_stem\_text}
- Gold answer: {gold\_answer}
- Question type: {question\_type}
- Strategy: replacement
- Reasoning steps: {reasoning\_steps}
- Copyable text (use this exact text when selecting substrings):  
<<<COPY  
{copyable\_text}  
>>>  
{prefix\_note}{answer\_guidance}{retry\_instructions}

Your task:

Generate {k} valid mappings that satisfy the replacement strategy. Each mapping should:

1. Replace one contiguous substring of the question stem
2. Cause a verifiable and detectable deviation from the gold answer
3. Ensure the replacement is semantically meaningful and natural
4. Change the question focus in a way that affects the expected answer

For each mapping, provide:

1. question\_index: The question number ({question\_index})
2. latex\_stem\_text: Exact LaTeX text of the question stem (must match the input exactly)
3. original\_substring: The substring to replace (must be a contiguous substring of latex\_stem\_text)
4. replacement\_substring: The replacement text
5. start\_pos: Start position of original\_substring relative to latex\_stem\_text (0-based index)
6. end\_pos: End position of original\_substring relative to latex\_stem\_text (exclusive, 0-based index)
7. target\_wrong\_answer: Description of how the answer should deviate (e.g., "focuses on different aspect", "changes key  
↪ concept")
8. reasoning: Brief explanation of why this mapping satisfies the strategy and how it causes deviation

IMPORTANT:

- The original\_substring MUST be an exact substring of latex\_stem\_text
- The start\_pos and end\_pos MUST be accurate ( $\text{start\_pos} + \text{len}(\text{original\_substring}) = \text{end\_pos}$ )
- The replacement should cause a verifiable deviation in the answer
- CRITICAL: The replacement\_substring MUST be DIFFERENT from the original\_substring. Do NOT generate mappings where  
↪ original\_substring == replacement\_substring. The replacement MUST change the text to create actual manipulation.
- CRITICAL: Neither original\_substring nor replacement\_substring can be empty strings. Both must contain actual text.
- LENGTH CONSTRAINT: The replacement\_substring MUST be smaller or equal in length to the original\_substring ( $\text{len}(\text{replacement\_substring}) \leq \text{len}(\text{original\_substring})$ ). This is critical for maintaining document layout and preventing  
↪ text overflow.
- latex\_stem\_text is provided exactly as it appears in the LaTeX source. Do NOT trim, normalise, or reformat it when  
↪ determining positions.
- The latex\_stem\_text may include \item tokens from enumerate environments. Keep the \item token intact and operate on the  
↪ descriptive text that follows it whenever possible.
- The replacement should be natural and semantically meaningful

Return as JSON array:

```
[
  {{
    "question_index": {question_index},
    "latex_stem_text": "...",
    "original_substring": "...",
    "replacement_substring": "...",
    "start_pos": 0,
    "end_pos": 5,
    "target_wrong_answer": "focuses on different aspect",
    "reasoning": "..."
  }},
  ...
]
```

Return ONLY valid JSON, no markdown or additional text.

# Using a Human-AI Teaming Approach to Create and Curate Scientific Datasets with the SCILIRE System

Necva Bölücü<sup>1†</sup>, Jessica Irons<sup>1†</sup>, Changhyun Lee<sup>1</sup>, Brian Jin<sup>1</sup>,  
Maciej Rybinski<sup>2\*</sup>, Huichen Yang<sup>1</sup>, Andreas Duenser<sup>1†</sup>, Stephen Wan<sup>1†</sup>

<sup>1</sup>CSIRO, Sydney, Australia

firstname.lastname@csiro.au

<sup>2</sup>ITIS, University of Málaga, Málaga, Spain

maciek.rybinski@uma.es

## Abstract

The rapid growth of scientific literature has made manual extraction of structured knowledge increasingly impractical. To address this challenge, we introduce SCILIRE, a system for creating datasets from scientific literature. SCILIRE has been designed around Human-AI teaming principles centred on workflows for verifying and curating data. It facilitates an iterative workflow in which researchers can review and correct AI outputs. Furthermore, this interaction is used as a feedback signal to improve future LLM-based inference. We evaluate our design using a combination of intrinsic benchmarking outcomes together with real-world case studies across multiple domains. The results demonstrate that SCILIRE improves extraction fidelity and facilitates efficient dataset creation.

## 1 Introduction

The exponential growth of scientific literature, which makes it increasingly challenging for researchers to stay up to date with the latest scientific developments (Cai et al., 2024; Reddy and Shojaee, 2025), represents an opportunity: scientific papers can be mined to generate high-value datasets (Dunn et al., 2022; Jiang et al., 2025; Wei et al., 2025). Such datasets are key in creating Artificial Intelligence (AI) to revolutionise scientific workflows and discovery, an endeavour generally referred to as *AI for Science* (AI4S).

Building on this potential, recent progress in Large Language Models (LLMs) offers powerful new tools, such as Elicit<sup>1</sup> and SciSpace<sup>2</sup>, to assist researchers in navigating and extracting knowledge from the vast scientific literature. However, these systems treat AI-data extraction as a single pass.

\*This work was done when the author was affiliated with CSIRO Data61.

<sup>†</sup>Primary authors for this work.

<sup>1</sup><https://elicit.com/>

<sup>2</sup><https://scispace.com>

Given that AI results are usually not perfect, single-pass tools force users to improve data outside of the tool without AI-assistance, a challenge when working with big datasets. This can limit the adoption of such tools in research workflows where output must conform to a certain standard and where such user processes to validate data manually are often arduous and time-consuming (Rahman and Kandogan, 2022; Pham and Lin, 2025; Schmidt et al., 2025).

Indeed, recent studies highlight risks related to LLM-based extraction: these models may generate hallucinated (confabulated) information, with empirical evaluations showing nontrivial error rates that require human correction and verification (Helms Andersen et al., 2025). Reviews of AI for literature synthesis further highlight ongoing problems with explainability and reliability, showing that generative AI cannot be fully trusted without expert oversight (Bolanos et al., 2024).

We adopt a Human-AI Teaming (HAT) (Berretta et al., 2023) design in SCILIRE, enabling users to curate data (hereafter: **HAT for Data Curation** (HAT-DC)). By combining expert validation with AI-assisted extraction, researchers can correct errors and mitigate hallucinations. Moreover, iterative human feedback helps improve model performance and fosters transparency, accountability, and trust in AI-enabled workflows (Gao et al., 2025; Schroeder et al., 2025).

SCILIRE differs from existing tools in that it supports iterative extraction and correction, enables dynamic sampling using the user’s curation history as an evolving source of examples, and scales to large literature search collections.

We evaluate the HAT-DC design of SCILIRE using public scientific datasets, provide real-world case studies, and report on feedback from users. Our contributions are: (1) evaluation of an HAT-DC workflow; (2) insights into the effectiveness of dynamic sampling.

| Tool          | Dynamic Sampling | Multi-record Support | Provenance Data | Table Export (CSV, JSON, ...) |
|---------------|------------------|----------------------|-----------------|-------------------------------|
| Elicit        | ✗ <sup>‡</sup>   | ✗                    | ✓ <sup>†</sup>  | ✓                             |
| SciSpace      | ✗ <sup>‡</sup>   | ✗                    | ✓ <sup>†</sup>  | ✓                             |
| NotebookLM    | ✗                | ✓                    | ✓ <sup>†</sup>  | ✗                             |
| Claude.ai     | ✗                | ✓                    | ✗               | ✓                             |
| SciLIRE(Ours) | ✓                | ✓                    | ✓ <sup>†</sup>  | ✓                             |

Table 1: Comparison of data curation tools’ functionality, highlighting differences from SciLIRE. <sup>‡</sup> Elicit and SciSpace can be made to accept static examples in column definitions. <sup>†</sup> Elicit, SciSpace and NotebookLM provide paragraph- or sentence-level citations. SciLIRE provides the degree of alignment with the source, along with relevant paragraphs. **Multi-record support** indicates if a tool is designed to produce multiple extracted records per document.

## 2 Related Work

Our work focuses on AI-powered tools leveraging LLMs for data curation, such as Elicit and SciSpace, which extract key information from PDFs. A listing of existing tools that extract some data from PDFs is presented in Table 1. For the HAT-DC, two key features are required: (1) the exporting of tables (in CSV or JSON), and (2) provision of provenance data for data verification and curation – this leaves Elicit and SciSpace as the two most relevant tools to our workflow.

Beyond these tools, there is a growing body of research that looks more broadly at how information from scientific papers can be turned into structured tables. ArxivDIGESTables (Newman et al., 2024) studies cross-paper table generation with LLMs and proposes an automatic evaluation method, while ArXiv2Table (Deng et al., 2024) presents a more comprehensive benchmark in the computer science domain. Several domain-specific efforts have also attempted to extract structured or tabular information directly from research papers in other domains, such as material science, chemistry and food manufacturing (Dunn et al., 2022; Wei et al., 2025; Bölücü et al., 2025). Collectively, these studies highlight growing interest in turning unstructured documents into tabular formats, a goal that AI-powered tools put into practice.

## 3 AI-augmented Curation Workflow

SciLIRE is designed to support the existing data curation workflow, with the use of AI focused on human skill augmentation. Typically, users iterate through possible schema structures (either provided by the user or selected from built-in templates) with

SciLIRE. The HAT-DC workflow is as follows:<sup>3</sup>

- 1. Bibliography Upload and Schema Definition:** Users upload a collection of documents and provide a schema file (e.g., a spreadsheet) specifying the column headers of the target curated table. The schema can be modified at any time. Once the documents are uploaded, SciLIRE automatically triggers the document preprocessing module to provide machine-readable versions of the text.
- 2. Pilot Phase:** Users select a small sample of documents ( $\leq 10$ ), generate data tables, and manually vet and correct the outputs. During this process, they also assess which documents and extracted records are relevant to the target aim. For relevant documents, users verify the generated results, correcting and curating data where necessary. Verified corrections are used by the system for dynamic sampling for subsequent batches. The pilot phase (multiple batches) is repeated for  $\sim 50$ –100 documents or until the user is satisfied.
- 3. Batch Phase:** The remaining documents are processed at scale. The user either continues to check data or exports the data as is.<sup>4</sup>

## 4 System Components

SciLIRE is built as a modular framework for data curation. Its architecture (Figure 1) comprises three modules that handle document preprocessing, LLM-based record generation, and verification support. Together, these modules form a flexible and transparent pipeline designed to support the user in a HAT-DC workflow.

### 4.1 Document Preprocessing

**Parsing.** SciLIRE checks each document for PDF type. If the PDF contains machine-readable text, it is processed using two PDF parsing pipelines (GROBID (Lopez, 2009)<sup>5</sup> and Apache Tika<sup>6</sup>). For PDFs with no machine-readable text (e.g., PDFs with only scanned images), SciLIRE uses the OCRmyPDF<sup>7</sup> library to convert a PDF of page images into a PDF with machine-readable text.

<sup>3</sup>See Appendix E for screenshots and a more detailed system usage description.

<sup>4</sup>The degree of further verification will depend on thresholds for acceptable quality related to the user’s end goal.

<sup>5</sup><https://github.com/kermitt2/grobid>

<sup>6</sup><https://tika.apache.org>

<sup>7</sup><https://github.com/ocrmypdf/OCRmyPDF>

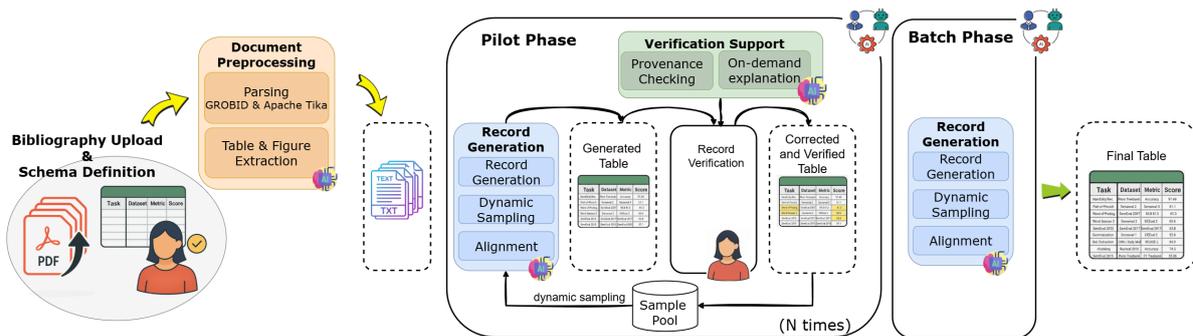


Figure 1: SCILIRE components and AI-augmented curation workflow.

The resulting PDFs are processed using the same pipelines as non-scanned PDFs, ensuring consistency across all documents: (1) The GROBID is the preferred pipeline, providing high-quality text extraction and structured metadata (Meuschke et al., 2023), particularly for scientific literature. The Tika pipeline provides an alternative when GROBID occasionally fails, thereby providing better support of other text genres. Since GROBID does not process figures and provides limited table extraction quality (Meuschke et al., 2023), we develop a custom table and figure extraction module.

**Table & figure extraction.** We implemented a two-stage pipeline (Figure 3, Appendix B) for the task. First, SCILIRE detects pixel regions indicating tables and figures in PDFs.<sup>8</sup> SCILIRE then performs table structure and caption recognition on the detected regions. Finally, the inferred table structure, together with table contents and captions, is rendered in markdown format. This is appended to the text extracted by GROBID based on the position of tables in the PDF. Implementation details of this module are provided in the Appendix B.

**Chunking.** Since LLMs have fixed context windows, long documents are segmented into overlapping chunks. We apply a configurable sliding-window strategy (by characters), preserving local coherence through overlapping spans. The window size and overlapping ratio are configurable parameters in the system (window size=LLM context length, overlap=10%).

## 4.2 Record Generation Module

Given a user’s schema that outlines concepts of interest (e.g., context or variables in an experiment with a measured result), SCILIRE automatically constructs prompts to generate structured records

using an LLM. Initially, the prompt follows a zero-shot baseline approach. If human-corrected data exists, the prompt uses a few-shot “In-Context Learning” (ICL) approach (Ghosh et al., 2024), which dynamically picks an example to include in the prompt.

SCILIRE uses the schema concepts to define a JSON dictionary structure (Oestreich and Müller, 2025) as the desired output format, which is included in the prompt (Appendix A). This structure also houses any ICL examples if required. Two versions of the prompt are then created, using data from the GROBID and Tika pipelines. The two prompts are then sent to the LLM.<sup>9</sup>

**Alignment.** Since the generation phase can yield two record sets (GROBID-based and Tika-based) per PDF, we merge the sets using the Hungarian maximum-matching algorithm (Kuhn, 1955) to identify overlapping records. We compute a similarity matrix by encoding records with sentence embeddings<sup>10</sup> and computing pairwise cosine similarity. The algorithm selects the optimal one-to-one alignment, allowing SCILIRE to suggest alternative records which users can compare during the curation task.

**Dynamic sampling for ICL.** Selecting an effective ICL example is critical in few-shot prompting. Instead of using static examples, SCILIRE retrieves a document-specific ICL example using BM25 (Robertson et al., 1995) from the pool of previously human-corrected documents. The closest match record is used as a 1-shot prototype (Ghosh et al., 2024). This dynamic ICL selection ensures that the LLM receives the most relevant example,

<sup>9</sup>SCILIRE can use any LLM, including closed and open weight models. In the SCILIRE, we use GPT-4o as the model.

<sup>10</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

<sup>8</sup>Currently, we do not further process figures.

improving extraction fidelity. This is a key distinction from other systems (e.g., Elicit and SciSpace).

### 4.3 Verification Support module

To allow experts to verify the generated content, SCILIRE provides several support tools:

**Provenance checking.** We implement a cell-level hallucination checker that compares generated content against the source document using fuzzy string matching fuzzywuzzy library<sup>11</sup>. The system visualises match strength as a graded signal, enabling users to build trust in generated answers by identifying aligned and unaligned answers.

**On-demand explanations and insights.** SCILIRE provides justifications of AI answers by aligning generated text to the source material, implemented as a search function that retrieves the top three supporting paragraphs from the original document using BM25. Aligned words are shown highlighted in bold to the user to assist data verification. We also provide the detected figures and tables from the PDF as additional resources.

SCILIRE also facilitates user requests for LLM-based explanations that identify source paragraphs relevant to a generated answer (Appendix A).<sup>12</sup> Together, these mechanisms support the user’s verification and curation tasks.

## 5 Experiments

Here we report on intrinsic benchmarking experiments to answer the following research questions: **(RQ1)** Can dynamic sampling (with data from a HAT-DC iterative workflow) lead to improved data extraction? **(RQ2)** How does SCILIRE (and the HAT-DC approach) compare to existing tools providing dataset generation capabilities?

### 5.1 Evaluation Framework

#### 5.1.1 Data and Metrics

To evaluate the dataset generation capabilities of SCILIRE, we conduct experiments on 18 datasets from five scientific domains, covering varying levels of granularity in data curation (Appendix C.1). We report on a primary evaluation metric which focuses on the very strict *record-level*  $F_1$  evaluation rather than cell-level evaluation, as typically

<sup>11</sup><https://pypi.org/project/fuzzywuzzy>

<sup>12</sup>Currently, SCILIRE uses GPT-4o as the model, but any model can be used.

| Dataset         | 0-shot | ICL-10       | ICL-50       | ICL-100      | ICL-all      |
|-----------------|--------|--------------|--------------|--------------|--------------|
| TDM5            | 10.14  | 19.01        | 23.01        | 24.54        | <b>25.02</b> |
| SciREX          | 3.66   | 13.51        | 15.49        | 15.71        | <b>18.27</b> |
| MPEA            | 29.23  | 32.52        | 30.39        | <b>30.81</b> | 30.64        |
| Diffusion       | 17.52  | <b>17.99</b> | 17.59        | –            | 17.20        |
| YSHEAY          | 5.34   | 7.87         | <b>8.30</b>  | 8.03         | 7.93         |
| CCRMG           | 1.82   | 2.48         | –            | –            | <b>2.89</b>  |
| Doping          | 7.55   | 13.23        | <b>14.65</b> | –            | 12.95        |
| MMD             | 0.78   | –            | –            | –            | <b>8.54</b>  |
| MRL             | 1.80   | 1.99         | <b>2.04</b>  | –            | 1.82         |
| PNCEExtract     | 31.14  | 40.21        | 42.46        | <b>45.04</b> | 43.59        |
| PolyIE          | 14.29  | 18.58        | 18.64        | –            | 18.34        |
| BRENDA_enzyme   | 26.58  | <b>36.81</b> | 34.33        | 36.14        | 36.59        |
| BRENDA_ribozyme | 11.74  | 17.84        | <b>19.12</b> | 18.69        | 18.48        |
| OPE             | 22.33  | <b>28.91</b> | 28.17        | 23.80        | 22.12        |
| PPE             | 51.67  | <b>67.83</b> | 64.76        | 64.60        | 64.60        |
| SE              | 36.6   | 42.71        | <b>47.00</b> | –            | 46.78        |
| AE              | 15.17  | <b>17.90</b> | 16.14        | –            | 11.53        |
| SuperMat        | 5.66   | <b>19.31</b> | 16.58        | 16.83        | 17.14        |
| AVG.            | 16.28  | 23.45        | 24.92        | <b>28.42</b> | 22.47        |

Table 2:  $F_1$  results across datasets. LLM: GPT-5.  $F_1$  reported with 0–100 scale; best score is **boldfaced**. For the full table, see Table 10.

scientists are compiling a set of scientific findings that comprise several dependent fields.<sup>13</sup>

### 5.2 RQ1. Evaluating a HAT-DC Approach

Using the benchmark data, we evaluate SCILIRE’s effectiveness in supporting data curation through HAT by simulating user corrections over multiple scientific domains. With a random sample of papers as an initial data pool, we use the associated human-authored ground truth data from that pool as a stand-in for the corrected records (see Appendix C.2). We then apply dynamic sampling from that pool to create ICL prompts for use with GPT-5.<sup>14</sup> Table 2 reports the summary results.<sup>15</sup>

In line with prior work (Jiang et al., 2024), we observe that generating accurate records is a hard task. The best reported  $F_1$  score is 67.83 (PPE dataset). The best averaged  $F_1$  score was just 28.42, highlighting the complexity of matching full records. Despite this, we see that the results support the HAT-DC approach: (1) all ICL variants are better than the zero-shot performance, and (2) using a sample pool of  $n = 100$  generally leads to the best performance with marginal gain or even performance degradation beyond that, as the increasing pool size tends to introduce redundancy rather than

<sup>13</sup>[https://github.com/bolucunecva/table\\_generation](https://github.com/bolucunecva/table_generation); See Appendix C.3 for detail of the evaluation metrics considered.

<sup>14</sup>GPT-5 was found to be the best performing LLM overall. See Appendix C.4 for performance of each tested LLM.

<sup>15</sup>Here we report on 1-shot ICL, guided by an engineering trade-off, given finite context, to balance between ICL examples and the flexibility of the system to accept an arbitrarily long list of columns.

| Dataset         | SciSpace |              | Elicit |              | SCILIRE      |              |
|-----------------|----------|--------------|--------|--------------|--------------|--------------|
|                 | 0-shot   | ICL-S        | 0-shot | ICL-S        | 0-shot       | ICL-D        |
| TDMS            | 0.0      | 0.0          | 0.0    | 3.13         | 3.97         | <b>11.76</b> |
| SciREX          | 0.0      | 0.0          | 1.08   | 6.49         | 2.98         | <b>18.22</b> |
| MPEA            | 13.26    | 13.26        | 0.0    | 0.0          | 40.67        | <b>42.27</b> |
| Diffusion       | 0.65     | 0.65         | 0.06   | 0.53         | 6.80         | <b>8.71</b>  |
| YSHEAY          | 0.0      | 0.0          | 2.22   | <b>13.33</b> | 3.29         | 5.54         |
| CCRMG           | 0.0      | 0.0          | 0.0    | <b>22.22</b> | 1.80         | 2.67         |
| Doping          | 0.0      | 0.0          | 3.6    | 9.01         | 5.41         | <b>12.12</b> |
| MMD             | 0.0      | 0.0          | 0.0    | 0.26         | 0.78         | <b>8.65</b>  |
| MRL             | 0.13     | 0.13         | 0.0    | 0.58         | <b>1.75</b>  | 1.57         |
| PNCEXtract      | 2.56     | 2.56         | 5.13   | 5.86         | 29.69        | <b>34.96</b> |
| PolyIE          | 0.0      | 0.0          | 0.0    | 0.74         | 12.05        | <b>18.58</b> |
| BRENDA_enzyme   | 0.05     | 0.33         | 0.42   | 1.35         | 34.34        | <b>47.44</b> |
| BRENDA_ribozyme | 0.0      | 0.0          | 1.96   | 4.34         | 26.03        | <b>30.95</b> |
| OPE             | 16.86    | <b>20.69</b> | 10.73  | 16.86        | 19.37        | 16.25        |
| PPE             | 5.41     | 0.0          | 1.80   | 12.61        | 48.83        | <b>62.69</b> |
| SE              | 0.0      | 0.0          | 0.78   | 0.78         | 34.12        | <b>45.25</b> |
| AE              | 0.0      | 0.0          | 0.0    | 0.0          | <b>21.23</b> | 15.31        |
| SuperMat        | 4.67     | 0.0          | 0.67   | 1.0          | 11.58        | <b>26.80</b> |
| AVG.            | 2.42     | 2.09         | 1.58   | 5.50         | 16.93        | <b>22.76</b> |

Table 3: F<sub>1</sub> results across datasets comparing SCILIRE with other data generation tools. F<sub>1</sub> reported with 0–100 scale; best score is **boldfaced**. SCILIRE results are based on GPT-5. Abbreviations: ICL-S: ICL static, ICL-D: ICL Dynamic ( $n=all$ ). For the full table, see Table 11. Results are shown for 10 randomly selected PDFs.

informative diversity, since dynamic sampling here results in the inclusion of samples that are highly similar.

### 5.3 RQ2. Related Commercial Software

As outlined in the Section 2, SciSpace and Elicit are comparable to SCILIRE as they also allow users to curate datasets. We evaluate the standard version of these commercial tools and our 1-shot usage of these tools, where we co-opt the column header input textbox in the UI to provide a statically chosen prototype example.<sup>16</sup> We use one randomly selected static sample per dataset as the static example.

Elicit, SciSpace, and SCILIRE can process a different number of PDFs, with Elicit providing the lower bound on PDF uploads. Here, we selected a random sample of 10 PDFs from each dataset to create a data subsample (in total 180 PDFs) used with each tool. The results are given in Table 3.<sup>17</sup> SciSpace outperforms Elicit in a zero-shot setting. However, Elicit is able to better utilise static ICL. Ultimately, SCILIRE consistently outperforms both tools, due to its dynamic sampling (ICL) capability. This highlights the benefits of the key differentiator of SCILIRE: the adoption of the HAT-DC approach over a single-pass AI approach.

<sup>16</sup>We use the Extract Data tools of SciSpace and Elicit for comparison.

<sup>17</sup>A full comparison between SciSpace and SCILIRE across all datasets is given in Table 12.

## 6 Case studies

To complement the intrinsic benchmarking results, we provide an overview of real usage. We report on four case studies from different scientific domains, with each study involving one or more researchers engaged in their own data curation tasks.

In each case study, domain expert scientists worked with the system development team to document task goals and evaluate progress. Variables such as the number of documents per pilot phase iteration were determined by the experts, who were able to seek advice from the development team. Users were free to judge when to end the pilot phase and proceed to the batch phase.

### 6.1 Scenarios

**Agriculture** The scientist was interested in extracting a dataset of reported plant-pest interactions (e.g., plant taxon, insect taxon). This research is ongoing with a desire to extend the extraction to thousands of articles. Here, the scientist performed multiple rounds in the pilot phase (10-20 articles), and the batch phase included 100 articles. The results of the batch phase were then manually verified by the scientist.

**Environmental Studies** The scientist was interested in performing a meta-analysis, collating environmental datasets published in academic publications and grey literature. The aim was to survey the field to determine a standard data schema and then to release a harmonised version of the amalgamated data. In the pilot phase of the SCILIRE workflow, the user conducted 4 rounds of validation on 5 documents each, before scaling up to perform the batch phase on the remainder of the dataset (approx. 200 documents). Because the results were intended for publication, high accuracy through human validation of the full table was required.

**Biochemistry** This multi-user team was interested in extracting and classifying bioactivity information on various plant compounds. Their goal in using SCILIRE was to find the “needles in the haystack”: the rare documents describing specific types of bioactivity. Because such information was scarce, the case study involved validating a mostly sparse table. The users performed the pilot phase through 3 iterations with 18 documents, before electing to perform the batch phase workflow on an additional 81 documents.

| Domain   | # Docs | # Records | Edits (0-shot) | Edits (ICL) | Time (0-shot) | Time (ICL) |
|----------|--------|-----------|----------------|-------------|---------------|------------|
| Agri.    | 42     | 96        | 30             | 20          | 6             | 7          |
| Env. St. | 20     | 20        | 35             | 3           | 13            | 3          |
| Biochem. | 18     | 56        | 7              | 5           | 11            | 3          |
| Med Man. | 15     | 15        | 12             | 25          | 6             | 9          |

Table 4: User validation behaviour across case studies at the initial (zero-shot) and final pilot phases (with ICL). Edits: Percentage of values edited by the user. Time: Average time (minutes) curating/correcting each PDF.

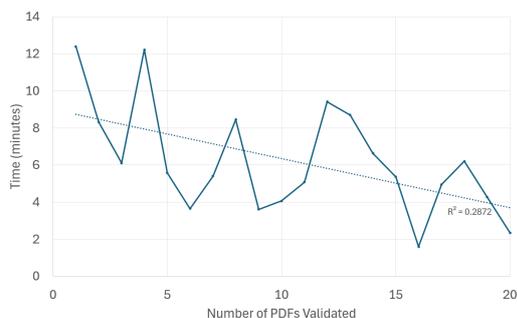


Figure 2: Average time spent validating data from the first 20 papers across the early adopter user cases.

**Medical manufacturing** In this case study, the scientists wished to extract and curate a small dataset (approx. 30 documents) reporting on drug trials. The extracted table included factors related to the drug formulation, experimental design and application area.

## 6.2 Overview of User Interactions

Table 4 presents an overview of four case studies (e.g., number of documents checked; number of records completed). These statistics confirm that a reference article can yield more than one record. In general, more edits are made in the initial pilot phase (a zero-shot setting), and in 3 of the 4 cases, in the dynamic sampling (ICL) setting, both the number of required edits and the time to complete the task drop. We see further evidence of the benefits of the HAT-DC approach in Figure 2, which shows the average time spent to validate data for the first 20 PDFs, averaged over the case studies. While there is a high degree of variability (std.err.: 2.5 mins/paper), we do see a significantly decreasing trend in time ( $p < 0.025$ ), indicating a decreasing validation workload as curation interaction increases.

## 7 Discussion and User Feedback

While quantitative benchmarks validate the performance of SCILIRE, qualitative analysis can provide

a deeper understanding of its practical strengths. This section presents user feedback that shows the real-world impact of the system.

### Verifying the need for data verification tools

As in earlier findings (e.g., Naddaf (2025)), interview data with users indicated that they were not prepared to trust AI-generated results in a fully automated (zero-shot) setting. That is, they wanted to inspect and review the content, with the ability to correct the results. This is consistent with our interaction analysis (see Appendix D), where we see that experts will prefer to verify the data before accepting or rejecting it. We observed that how the user performed this review was idiosyncratic: some preferred to check the source PDFs, while others preferred the supporting paragraphs. *This validates our design decision to include tools to facilitate the curation of AI-generated data.*

**Opportunities in efficiency and scale** Users engaging with SCILIRE saw opportunities on two fronts: (a) **Scale**, the ability to generate datasets at a scale that would not be feasible manually; and (b) **Efficiency**, the ability to create datasets with less manual effort. Scaling emerged as the most frequently reported need from users. For example, the agriculture user reported that, given the need to process over 6,000 documents, they could not attempt the task manually. Another noted that for historical datasets, SCILIRE allowed revisiting documents to add columns on experimentation context, a task that would otherwise not be feasible manually, given the dataset size. Users noted that the dataset compilation task would normally be performed by a team of researchers; with SCILIRE, it could now be managed by one researcher with a modest budget to cover the additional computational costs.

Given the fallibility of AI, users reported that they preferred checking pre-populated fields to manually populating the table from scratch. They felt that having a starting point from which to start the validation process increased their efficiency. This efficiency was recognised by users as time savings. For example, one user estimated that if they were to redo a recent manual systematic review with SCILIRE, they might perform the task twice as fast. Another estimated that SCILIRE reduced the time required from 2-3 months to 1 month. *This validates the design decision to consider data curation AI tools leading to productivity benefits.*

## 8 Conclusion

We presented SCILIRE, a Human-AI Teaming (HAT) system, where AI capabilities and human expert judgements work together to enable effective dataset creation from scientific literature. Our intrinsic benchmarking demonstrates dynamic samples for few-shot learning (from an iterative workflow), guides, and improves AI results. These outcomes mirror findings in our user studies, which reveal that the SCILIRE and the HAT approach lead to new opportunities for scientists. With an AI-enabled data creation and curation workflow, users can work at scales that would be infeasible without AI support, while gaining efficiency in overseeing and validating the AI-generated results.

As future work, we plan to extend SCILIRE with additional analysis capabilities, allowing users to transform curated datasets into actionable insights by facilitating the discovery of trends and patterns.

## Acknowledgments

This work is supported by CSIRO as part of the AI for Missions Program (<https://research.csiro.au/ai4m/>). We acknowledge the CSIRO Language Technology team and Samuel Walker, a former CSIRO employee (UX Engineer), for their contributions to the SCILIRE.

## Ethical Consideration

The public datasets used in quantitative experiments are adopted from existing repositories. Therefore, we do not foresee any serious or harmful issues related to their content. We have collected documents with the given identifiers (e.g., DOI, PubMed ID).

Collection of user data was approved by the CSIRO Social and Interdisciplinary Science Human Research Ethics Committee (090/23), and participants whose data are presented here provided informed written consent.

## References

Sophie Berretta, Alina Tausch, Greta Ontrup, Björn Gilles, Corinna Peifer, and Annette Kluge. 2023. *Defining human-ai teaming the human-centered way: a scoping review and network analysis*. *Frontiers in Artificial Intelligence*, Volume 6 - 2023.

Francisco Bolanos, Angelo Salatino, Francesco Osborne, and Enrico Motta. 2024. *Artificial intelligence for literature reviews: Opportunities and challenges*. *Artificial Intelligence Review*, 57(10):259.

Necva Bölücü, Jordan Pennells, Huichen Yang, Maciej Rybinski, and Stephen Wan. 2025. *An Evaluation of Large Language Models for Supplementing a Food Extrusion Dataset*. *Foods*, 14(8):1355.

Christopher KH Borg, Carolina Frey, Jasper Moh, Tresa M Pollock, Stéphane Gorsse, Daniel B Miracle, Oleg N Senkov, Bryce Meredig, and James E Saal. 2020. *Expanded dataset of mechanical properties and observed phases of multi-principal element alloys*. *Scientific Data*, 7(1):430.

Hengxing Cai, Xiaochen Cai, Junhan Chang, Sihang Li, Lin Yao, Wang Changxin, Zhifeng Gao, Hongshuai Wang, Li Yongge, Mujie Lin, Shuwen Yang, Jiankun Wang, Mingjun Xu, Jin Huang, Xi Fang, Jiaxi Zhuang, Yuqi Yin, Yaqi Li, Changhong Chen, and 4 others. 2025. *SciAssess: Benchmarking LLM Proficiency in Scientific Literature Analysis*. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 2335–2357, Albuquerque, New Mexico. Association for Computational Linguistics.

Hengxing Cai, Xiaochen Cai, Shuwen Yang, Jiankun Wang, Lin Yao, Zhifeng Gao, Junhan Chang, Sihang Li, Mingjun Xu, Changxin Wang, Hongshuai Wang, Yongge Li, Mujie Lin, Yaqi Li, Yuqi Yin, Linfeng Zhang, and Guolin Ke. 2024. *Uni-SMART: Universal Science Multimodal Analysis and Research Transformer*. *CoRR*, abs/2403.10301.

Zhaowei Cai and Nuno Vasconcelos. 2018. *Cascade r-cnn: Delving into high quality object detection*. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162.

Jerry Cheung, Yuchen Zhuang, Yinghao Li, Pranav Shetty, Wantian Zhao, Sanjeev Grampurohit, Rampi Ramprasad, and Chao Zhang. 2024. *POLYIE: A Dataset of Information Extraction from Polymer Material Scientific Literature*. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2370–2385. Association for Computational Linguistics.

Zewen Chi, Heyan Huang, Heng-Da Xu, Houjin Yu, Wanxuan Yin, and Xian-Ling Mao. 2019. *Complicated Table Structure Recognition*.

Filip Darmanović, Allan Hanbury, and Markus Zlabinger. 2023. *SCI-3000: A Dataset for Figure, Table and Caption Extraction from Scientific PDFs*. page 234–251, Berlin, Heidelberg. Springer-Verlag.

Zheyue Deng, Chunkit Chan, Weiqi Wang, Yuxi Sun, Wei Fan, Tianshi Zheng, Yauwai Yim, and Yangqiu Song. 2024. *Text-Tuple-Table: Towards Information Integration in Text-to-Table Generation via Global Tuple Extraction*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9300–9322. Association for Computational Linguistics.

- Alexander Dunn, John Dagdelen, Nicholas Walker, Sanghoon Lee, Andrew S Rosen, Gerbrand Ceder, Kristin Persson, and Anubhav Jain. 2022. [Structured information extraction from complex scientific text with fine-tuned large language models](#). *arXiv preprint arXiv:2212.05238*.
- Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. 2024. [Sciknoweval: Evaluating multi-level scientific knowledge of large language models](#). *arXiv preprint arXiv:2406.09098*.
- Luca Foppiano, Sae Dieb, Akira Suzuki, Pedro Baptista de Castro, Suguru Iwasaki, Azusa Uzuki, Miren Garbine Esparza Echevarria, Yan Meng, Kensei Terashima, Laurent Romary, Yoshihiko Takano, and Masashi Ishii. 2021. [SuperMat: construction of a linked annotated dataset from superconductors-related publications](#). *Science and Technology of Advanced Materials: Methods*, 1(1):34–44.
- Jiayuan Gao, Yingwei Zhang, Yiqiang Chen, Yihan Dong, Yuanzhe Chen, Shuchao Song, Boshi Tang, and Yang Gu. 2025. [Agent-in-the-loop to distill expert knowledge into artificial intelligence models: a survey](#). *Artificial Intelligence Review*, 58(9):266.
- Satanu Ghosh, Neal Brodnik, Carolina Frey, Collin Holgate, Tresa Pollock, Samantha Daly, and Samuel Carton. 2024. [Toward Reliable Ad-hoc Scientific Information Extraction: A Case Study on Two Materials Dataset](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 15109–15123. Association for Computational Linguistics.
- Ross Girshick. 2015. [Fast r-cnn](#). In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. [Mask r-cnn](#). In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969.
- T Helms Andersen, TM Marcussen, AD Termanssen, TWH Lawaetz, and O Nørgaard. 2025. [Using Artificial intelligence tools as second reviewers for data extraction in systematic reviews: a performance comparison of two AI tools against human reviewers](#). *Cochrane Evidence Synthesis and Methods*, 3(4):e70036.
- Yufang Hou, Charles Jochim, Martin Gleize, Francesca Bonin, and Debasis Ganguly. 2019. [Identification of Tasks, Datasets, Evaluation Metrics, and Numeric Scores for Scientific Leaderboards Construction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5203–5213. Association for Computational Linguistics.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2025. [Qwen2.5 technical report](#). *Preprint*, arXiv:2412.15115.
- Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. [SciREX: A Challenge Dataset for Document-Level Information Extraction](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7506–7516. Association for Computational Linguistics.
- Jinling Jiang, Jie Hu, Siwei Xie, Menghao Guo, Yuhang Dong, Shuai Fu, Xianyue Jiang, Zhenlei Yue, Junchao Shi, Xiaoyu Zhang, and 1 others. 2025. [Enzyme Co-Scientist: Harnessing Large Language Models for Enzyme Kinetic Data Extraction from Literature](#). *bioRxiv*, pages 2025–03.
- Peiwen Jiang, Xinbo Lin, Zibo Zhao, Ruhui Ma, Yvonne Jie Chen, and Jinhua Cheng. 2024. [TKGT: Redefinition and A New Way of Text-to-Table Tasks Based on Real World Demands and Knowledge Graphs Augmented LLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16112–16126. Association for Computational Linguistics.
- Ghazal Khalighinejad, Defne Circi, L. Brinson, and Bhuwan Dhingra. 2024. [Extracting Polymer Nanocomposite Samples from Full-Length Documents](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13163–13175. Association for Computational Linguistics.
- Harold W Kuhn. 1955. [The Hungarian method for the assignment problem](#). *Naval research logistics quarterly*, 2(1-2):83–97.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. [Microsoft coco: Common objects in context](#). In *European conference on computer vision*, pages 740–755. Springer.
- Patrice Lopez. 2009. [GROBID: Combining automatic bibliographic data recognition and term extraction for scholarship publications](#). In *Research and Advanced Technology for Digital Libraries: 13th European Conference, ECDL 2009, Corfu, Greece, September 27-October 2, 2009. Proceedings 13*, pages 473–474. Springer.
- Norman Meuschke, Apurva Jagdale, Timo Spinde, Jelena Mitrović, and Bela Gipp. 2023. [A benchmark of pdf information extraction tools using a multi-task and multi-domain evaluation framework for academic documents](#). In *International Conference on Information*, pages 383–405. Springer.
- Miryam Naddaf. 2025. [How are researchers using AI? Survey reveals pros and cons for science](#). *Nature*.
- Benjamin Newman, Yoonjoo Lee, Aakanksha Naik, Pao Siangliulue, Raymond Fok, Juho Kim, Daniel S Weld, Joseph Chee Chang, and Kyle Lo. 2024. [ArxivDIGESTables: Synthesizing Scientific Literature into Tables using Language Models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9612–9631. Association for Computational Linguistics.

- Julian Oestreich and Lydia Müller. 2025. [Evaluating Structured Decoding for Text-to-Table Generation: Evidence from Three Datasets](#). *Preprint*, arXiv:2508.15910.
- ShengYun Peng, Aishwarya Chakravarthy, Seongmin Lee, Xiaojing Wang, Rajarajeswari Balasubramanian, and Duen Horng Chau. 2024. [UniTable: Towards a Unified Framework for Table Recognition via Self-Supervised Pretraining](#). In *NeurIPS 2024 Third Table Representation Learning Workshop*.
- Vy Pham and Fu-ren Lin. 2025. [The Design and Evaluation of the Collaboration between Researchers and Generative AI for Systematic Literature Reviews](#). In *Proceedings of the Hawaii International Conference on System Sciences (HICSS)*.
- Maciej P Polak, Shrey Modi, Anna Latosinska, Jinning Zhang, Ching-Wen Wang, Shaonan Wang, Ayan Deep Hazra, and Dane Morgan. 2024. [Flexible, model-agnostic method for materials data extraction from text using general purpose language models](#). *Digital Discovery*, 3(6):1221–1235.
- Sajjadur Rahman and Eser Kandogan. 2022. [Characterizing Practices, Limitations, and Opportunities Related to Text Information Extraction Workflows: A Human-in-the-loop Perspective](#). In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22. Association for Computing Machinery.
- Chandan K Reddy and Parshin Shojaee. 2025. [Towards scientific discovery with generative ai: Progress, opportunities, and challenges](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 28601–28609.
- Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, and 1 others. 1995. Okapi at TREC-3. *Nist Special Publication Sp*, 109:109.
- Lena Schmidt, Ailbhe N Finnerty Mutlu, Rebecca Elmore, Babatunde K Olorisade, James Thomas, and Julian PT Higgins. 2025. [Data extraction methods for systematic review \(semi\) automation: Update of a living systematic review](#). *F1000Research*, 10:401.
- Noah L. Schroeder, Chris Davis Jaldi, and Shan Zhang. 2025. [Large Language Models with Human-In-The-Loop Validation for Systematic Review Data Extraction](#). *Preprint*, arXiv:2501.11840.
- Ray Smith. 2007. [An overview of the tesseract ocr engine](#). In *Ninth international conference on document analysis and recognition (ICDAR 2007)*, volume 2, pages 629–633. IEEE.
- Roelien C. Timmer, Yufang Hou, and Stephen Wan. 2025. [A Position Paper on the Automatic Generation of Machine Learning Leaderboards](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 30749–30772. Association for Computational Linguistics.
- Galen Wei, Xinchun Ran, Runeem AI-Abssi, and Zhongyue Yang. 2025. [Finding the dark matter: Large language model-based enzyme kinetic data extractor and its validation](#). *Protein Science*, 34(9):e70251.
- Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, Chunyu Kit, Clara Grazian, Wenjie Zhang, and 1 others. 2023. [Darwin series: Domain specific large language models for natural science](#). *arXiv preprint arXiv:2308.13565*.
- Youxue Zhang, Huaiwei Ni, and Yang Chen. 2010. [Diffusion data in silicate melts](#). *Reviews in Mineralogy and Geochemistry*, 72(1):311–408.

## A Prompts

The prompt used in the record generation module:

```
Please, extract ATTRIBUTE_1, ATTRIBUTE_2, ..., ATTRIBUTE_n from the given article.
```

```
For the extracted information, you MUST respond in a list of JSON dictionaries structure with the given Dictionary Key Mapping.
```

```
[Dictionary Key Mapping in your response]
{
  ATTRIBUTE_1: (example: VALUE_1),
  ATTRIBUTE_2: (example: VALUE_2),
  ...
  ATTRIBUTE_n: (example: ATTRIBUTE_n)
}
```

```
[Given Article Start]
ARTICLE CONTENT
[Given Article End]
```

The prompts used in on-demand explanations by LLMs:

```
Please find the relevant paragraph that shows that the ATTRIBUTE is VALUE from given article.
```

```
[Given Article Start]
ARTICLE CONTENT
[Given Article End]
```

## B Table & Figure Extraction Module

Our pipeline (Figure 3) consists of two main stages: (1) table, figure, and caption detection, and (2) table structure recognition (TSR).

### B.1 Stage I: Table, figure, caption detection

**Training.** The detection model architecture is based on Cascade R-CNN (Cai and Vasconcelos, 2018). We fine-tune the pre-trained (on COCO (Lin et al., 2014) dataset) Cascade R-CNN using the

SCI-3000 dataset (Darmanović et al., 2023)<sup>18</sup> to detect tables, figures, and their captions from PDF page images. All the implementations are based on Detectron2<sup>19</sup>. The results are given in Table 5.

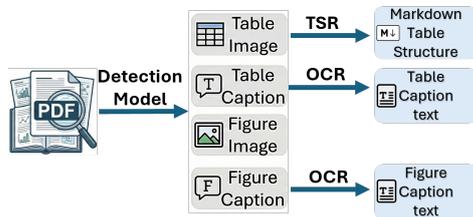


Figure 3: The Table & Figure Extraction module.

| Model                        | mAP@.5       | mAP@.75      | mAP          |
|------------------------------|--------------|--------------|--------------|
| Fast R-CNN (Girshick, 2015)  | 96.18        | 94.20        | 87.36        |
| Mask R-CNN (He et al., 2017) | 96.21        | 94.39        | 87.57        |
| Cascade R-CNN                | <b>97.12</b> | <b>95.40</b> | <b>90.55</b> |

Table 5: Evaluation results comparing Cascade R-CNN with baseline approaches. The best results are **bold-faced**.

| Model                             | TEDS         | mAP@.5       |
|-----------------------------------|--------------|--------------|
| UniTable (Peng et al., 2024)      | 95.23        | 96.23        |
| UniTable (finetuned on SciTSR)    | <b>97.98</b> | <b>96.98</b> |
| Qwen2.5-VL-72B (Hui et al., 2025) | 82.72        | —            |

Table 6: Evaluation results comparing UniTable with and without finetuning. The best results are **boldfaced**.

**Inference.** We convert each PDF page into an image. The detection model then identifies tables, figures, and their captions on each page.

## B.2 Stage II: Table structure recognition

**Training.** Table structure recognition architecture is based on UniTable (Peng et al., 2024). We fine-tune the model on the SciTSR dataset (Chi et al., 2019), which contains 15k table images (12k for training and 3k for testing) and their corresponding structure labels obtained from LaTeX source files. We convert the SciTSR data annotation format to HTML to align with the UniTable training and TEDS evaluation format. The comparison results between the base model and our fine-tuned

<sup>18</sup>The original SCI-3000 dataset uses a single caption label for both tables and figures. For our task, we re-annotate the captions to separate table captions from figure captions to provide finer granularity. The resulting dataset contains four labels: table, table caption, figure, and figure caption.

<sup>19</sup><https://github.com/facebookresearch/detectron2>

model are shown in Table 6. Since the model predicts HTML-formatted table structures, we convert the predicted HTML into markdown as the final output.

**Inference.** The TSR model reads each detected table image and infers the table structure along with its content. It outputs the table in markdown format, the extracted cell content, and the corresponding caption (extracted by Tesseract OCR (Smith, 2007)).

## C Experimental Details

### C.1 Datasets

For the quantitative experiments (Section 5.1), we use datasets from multiple domains spanning machine learning, materials science, chemistry, medicine, and physics<sup>20</sup>. Together, these datasets provide broad, heterogeneous, and real scientific documents for evaluating SciLIRE.

- **Machine Learning:** Leaderboard construction task from NLP papers (TDMS, SciREX).
- **Materials Science:** Extraction of compositions, experimental parameters, and material properties from diverse subfields (MPEA, Diffusion, YSHEAY, CCRMG, MRL, Doping, MMD, PolyIE, PNCEXtract).
- **Chemistry:** Extraction of molecular structures, reaction properties, and material characteristics (BRENDA, OPE, PPE, SE).
- **Medicine:** Affinity extraction involving molecules, SMILES, and bioassay targets (AE).
- **Physics:** Extraction of superconductor materials and properties (SuperMat).

The statistical details of datasets are given in Table 7.

### C.2 Experimental Settings of HAT-DC Workflow

Algorithm 1 outlines the HAT-DC-mimicking workflow. This procedure allows us to simulate iterative human-AI interactions and evaluate the benefits of HAT-DC guidance in a controlled, reproducible manner. As a reminder, the user actively selects samples for a sample pool for dynamic sampling in SciLIRE, rather than random sampling.

<sup>20</sup>You can find more details about datasets (e.g., schema) at [https://github.com/bolucunecva/table\\_generation](https://github.com/bolucunecva/table_generation)

| Domain           | Dataset         | # Documents | Avg. Pages | # Records | Records / Document | Schema Size | Reference                   |
|------------------|-----------------|-------------|------------|-----------|--------------------|-------------|-----------------------------|
| Machine Learning | TDMS            | 332         | 10.52      | 904       | 2.72               | 4           | Hou et al. (2019)           |
|                  | SciREX          | 372         | 11.83      | 1,897     | 5.34               | 5           | Jain et al. (2020)          |
| Material Science | MPEA            | 264         | 8.70       | 1,544     | 5.85               | 17          | Borg et al. (2020)          |
|                  | Diffusion       | 93          | 14.02      | 3,533     | 37.99              | 18          | Zhang et al. (2010)         |
|                  | YSHEAY          | 219         | 15.65      | 837       | 3.82               | 3           | Polak et al. (2024)         |
|                  | CCRMG           | 24          | 11.96      | 43        | 1.79               | 3           | Polak et al. (2024)         |
|                  | Doping          | 66          | 6.66       | 544       | 8.24               | 3           | Dunn et al. (2022)          |
|                  | MMD             | 3           | 6.00       | 140       | 46.67              | 16          | Xie et al. (2023)           |
|                  | MRL             | 100         | 5.81       | 993       | 9.93               | 19          | Xie et al. (2023)           |
|                  | PNCEXtract      | 155         | 8.92       | 838       | 5.41               | 6           | Khalighinejad et al. (2024) |
|                  | PolyIE          | 76          | 8.58       | 2,337     | 30.75              | 3           | Cheung et al. (2024)        |
| Chemistry        | BRENDA-enzyme   | 155         | 12.45      | 4,210     | 27.16              | 13          | Jiang et al. (2025)         |
|                  | BRENDA-ribozyme | 163         | 11.05      | 1,756     | 10.77              | 17          | Jiang et al. (2025)         |
|                  | OPE             | 104         | 7.55       | 255       | 2.45               | 7           | Cai et al. (2025)           |
|                  | PPE             | 109         | 6.97       | 265       | 2.43               | 6           | Cai et al. (2025)           |
|                  | SE              | 96          | 8.01       | 2,363     | 24.61              | 3           | Cai et al. (2025)           |
| Medicine         | AE              | 40          | 4.30       | 406       | 10.15              | 3           | Cai et al. (2025)           |
| Physics          | SuperMat        | 142         | 8.55       | 1,301     | 9.16               | 4           | Foppiano et al. (2021)      |

Table 7: Statistics of datasets across five domains.

### Algorithm 1: Mimic HAT-DC Workflow

**Input:** Dataset  $D$ ; pool size  $k$ ; samples  $m$ ; schema  $H$ ; LLM  $LLM$   
**Output:** Table  $T$   
 $T \leftarrow []$ ;  
 $S \leftarrow$  samples of  $D$ ;  $N \leftarrow |S|$ ;  
**for**  $t \leftarrow 1$  **to**  $N$  **do**  
     $test \leftarrow S[t]$ ;  
     $trainCandidates \leftarrow S \setminus \{test\}$ ;  
     $pool \leftarrow RandomSample(trainCandidates, k)$ ;  
     $ranked \leftarrow BM25Rank(test, pool)$ ;  
     $sample \leftarrow ranked[1:m]$ ;  
     $records \leftarrow Prediction(LLM, test, sample, H)$ ;  
    append  $\langle records \rangle$  to  $T$ ;  
**return**  $T$

### C.3 Evaluation metrics

We adapt table-generation metrics for record-level evaluation (Ghosh et al., 2024; Khalighinejad et al., 2024; Cheung et al., 2024; Feng et al., 2024; Jiang et al., 2025). Precision, Recall, and  $F_1$  are computed with, where a cell counts as correct only if it exactly matches the aligned reference. ChrF is also reported in the record-aligned setting, measuring character n-gram overlap to capture partial matches and minor differences.

### C.4 Models

The details of the models used in the record generation module of SCILIRE (Section 4.2) are given in Table 8.

### C.5 All Experimental Results

The overall experimental results across all datasets and models are summarised in Table 9. Table 10 the full per-dataset results using GPT-5 as LLM.

| Model                 | # of Par. | Context Length | Open-Source | Family    |
|-----------------------|-----------|----------------|-------------|-----------|
| GPT-OSS:20b           | 20B       | 128K           | ✓           | OpenAI    |
| GPT-OSS:120b          | 120B      | 128K           | ✓           | OpenAI    |
| Gemma3:1B             | 1B        | 32K            | ✓           | Google    |
| Gemma3:4B             | 4B        | 128K           | ✓           | Google    |
| Gemma3:12B            | 12B       | 128K           | ✓           | Google    |
| Gemma3:27B            | 27B       | 128K           | ✓           | Google    |
| Qwen3:0.6B            | 0.6B      | 40K            | ✓           | Qwen      |
| Qwen3:4B              | 4B        | 256K           | ✓           | Qwen      |
| Qwen3:14B             | 14B       | 40K            | ✓           | Qwen      |
| Qwen3:32B             | 32B       | 40K            | ✓           | Qwen      |
| Phi-4                 | 14B       | 16K            | ✓           | Microsoft |
| DeepSeek-R1-Llama:8B  | 8B        | 128K           | ✓           | DeepSeek  |
| DeepSeek-R1-Llama:70B | 70B       | 128K           | ✓           | DeepSeek  |
| DeepSeek-R1-Qwen3:14B | 14B       | 128K           | ✓           | DeepSeek  |
| DeepSeek-R1-Qwen3:32B | 32B       | 128K           | ✓           | DeepSeek  |
| GPT-5                 | ?         | 400K           | ✗           | OpenAI    |

Table 8: Details of models used in record generation module (Section 4.2) of SCILIRE.

Detailed tool-level comparisons are reported in Table 11, and Table 12 presents a comprehensive comparison between SCILIRE and SciSpace across all datasets using the complete document collections.

## D User Data

In Figure 4, we see that most of the acceptances and rejections of AI extracted records are performed via a verification step of either checking the tool’s built-in verification support features or checking the reference. This supports the premise that scientists wish to verify and curate data, consistent with our HAT-DC approach. That pure automation is not necessarily what scientists are looking for.

## E Demo Walkthrough

In the demonstration software accompanying this paper, we consider the scenario where a user aims to curate a dataset for a well-known scenario in

| Dataset         | Gemma3:1B | Gemma3:4B | Gemma3:12B | Gemma3:27B | Qwen3:0.6B | Qwen3:4B | Qwen3:14B | Qwen3:32B    | Phi-4 | GPT-OSS:20B  | GPT-OSS:120B | Deepseek-R1-Llama:8B | Deepseek-R1-Llama:70B | Deepseek-R1-Qwen3:14B | Deepseek-R1-Qwen3:32B | GPT-5        |
|-----------------|-----------|-----------|------------|------------|------------|----------|-----------|--------------|-------|--------------|--------------|----------------------|-----------------------|-----------------------|-----------------------|--------------|
| TDMS            | 12.20     | 12.13     | 18.80      | 21.90      | 17.13      | 25.09    | 17.98     | 21.91        | 16.10 | 22.10        | 23.02        | 21.19                | 22.53                 | <b>26.04</b>          | 23.13                 | 25.02        |
| SciREX          | 4.10      | 3.97      | 14.45      | 15.08      | 7.82       | 15.69    | 17.06     | 16.53        | 10.56 | 18.12        | 17.65        | 6.26                 | 08.45                 | 15.99                 | 14.56                 | <b>18.27</b> |
| MPEA            | 9.29      | 10.45     | 13.45      | 18.36      | 8.43       | 15.23    | 18.43     | 17.56        | 13.54 | 19.45        | 16.34        | 7.21                 | 14.56                 | 12.42                 | 13.21                 | <b>30.64</b> |
| Diffusion       | 0.45      | 4.20      | 5.45       | 8.36       | 0.68       | 6.34     | 9.12      | 12.25        | 3.12  | 14.21        | 16.34        | 7.10                 | 9.21                  | 12.10                 | 12.46                 | <b>17.20</b> |
| YSHEAY          | 9.23      | 10.13     | 12.21      | 13.10      | 17.35      | 14.25    | 15.99     | <b>16.48</b> | 11.24 | 16.23        | 14.56        | 12.48                | 13.21                 | 12.02                 | 10.45                 | 7.93         |
| CCRMG           | 5.85      | 5.23      | 5.71       | 3.37       | 4.96       | 5.01     | 5.89      | <b>9.03</b>  | 8.91  | 6.10         | 7.03         | 2.24                 | 3.87                  | 5.02                  | 4.65                  | 2.89         |
| Doping          | 6.44      | 11.99     | 13.70      | 15.82      | 3.68       | 11.83    | 14.78     | 15.20        | 13.15 | 15.33        | <b>16.45</b> | 9.08                 | 10.21                 | 14.04                 | 13.67                 | 12.95        |
| MMD             | 1.22      | 3.15      | 4.7        | 3.96       | 1.26       | 1.30     | 11.05     | 5.92         | 4.23  | <b>11.08</b> | 10.98        | 1.39                 | 8.45                  | 2.48                  | 6.10                  | 8.54         |
| MRL             | 0.23      | 0.45      | 0.59       | 1.23       | 0.67       | 1.34     | 1.97      | <b>1.95</b>  | 0.45  | 1.78         | 1.82         | 0.56                 | 1.10                  | 1.13                  | 1.12                  | 1.82         |
| PNCEextract     | 12.66     | 29.29     | 34.30      | 35.78      | 20.82      | 21.07    | 42.49     | 27.86        | 19.92 | 26.12        | 29.48        | 20.73                | 32.87                 | 30.35                 | 28.23                 | <b>43.59</b> |
| PolyIE          | 3.54      | 17.04     | 21.38      | 23.69      | 6.56       | 16.47    | 20.60     | 19.05        | 17.25 | <b>21.10</b> | 20.03        | 15.12                | 14.32                 | 17.51                 | 16.18                 | 18.34        |
| BRENDA_enzyme   | 3.08      | 17.04     | 16.85      | 22.97      | 3.11       | 5.44     | 24.07     | 30.83        | 15.12 | 31.27        | 32.33        | 7.31                 | 13.23                 | 14.45                 | 13.45                 | <b>36.59</b> |
| BRENDA_ribozyme | 1.16      | 3.80      | 6.05       | 7.74       | 2.05       | 2.81     | 6.47      | 14.98        | 4.36  | 15.98        | 14.65        | 5.82                 | 6.23                  | 5.81                  | 8.10                  | <b>18.48</b> |
| OPE             | 5.34      | 9.23      | 12.27      | 15.20      | 13.21      | 16.05    | 17.54     | 16.89        | 9.24  | 20.13        | 21.05        | 15.32                | 14.23                 | 13.45                 | 17.89                 | <b>22.12</b> |
| PPE             | 16.38     | 41.16     | 45.65      | 60.48      | 50.21      | 54.62    | 63.43     | <b>66.80</b> | 55.52 | 66.10        | 65.86        | 34.63                | 54.19                 | 56.86                 | 55.10                 | 64.60        |
| SE              | 4.35      | 24.86     | 29.13      | 31.06      | 9.36       | 23.23    | 43.87     | 43.42        | 32.73 | 45.10        | 44.61        | 17.50                | 22.10                 | 30.75                 | 31.20                 | <b>46.78</b> |
| AE              | 1.92      | 8.0       | 19.19      | 21.75      | 2.97       | 21.28    | 23.20     | <b>25.89</b> | 12.54 | 24.45        | 23.76        | 19.22                | 20.10                 | 19.63                 | 18.65                 | 11.53        |
| SuperMat        | 12.74     | 22.95     | 28.40      | 23.51      | 10.47      | 11.55    | 13.48     | 13.80        | 19.53 | <b>21.0</b>  | 20.14        | 8.75                 | 9.56                  | 13.41                 | 14.32                 | 17.14        |
| AVG.            | 6.12      | 13.06     | 16.79      | 19.07      | 10.04      | 14.92    | 20.41     | 20.91        | 14.86 | 21.98        | 22.01        | 11.77                | 15.47                 | 16.86                 | 16.80                 | <b>22.47</b> |

Table 9:  $F_1$  results across datasets for multiple LLMs and their ability to benefit from ICL Dynamic sampling ( $n=all$ ).  $F_1$  reported with 0–100 scale; best score is **boldfaced**. For detailed results, see <https://github.com/bolucunecva/scilire>.

machine learning and computer science: generating a leaderboard in the computer science domain (for an overview, see Timmer et al. (2025)).

Notionally, the user would first create a project by uploading the schema and a collection of documents (Figure 5). For demonstration purposes, papers have been uploaded and the project created in advance.<sup>21</sup>

In Figure 6, we see how the user navigates to their project. By clicking on *Projects* in the left-hand navigational menu, the user can click the option “Manage” for the registered project, here called “TDMS” (For the leaderboard dataset of the same name, which stands for “Task, Dataset, Metric, Score” (Hou et al., 2019)).

In Figure 7, we see how the user can revisit the outcomes of the Human-AI Team approach and the iterative data curation workflow. The interface contains a table labelled “Previous Extractions”, where Samples 1-3 represent iterations through the pilot phase.

Figure 8 shows a table in Sample 1, which has been edited and curated by the user. Sample 1 is the first batch, in which SCILIRE generates records

for the selected PDFs whereby the LLM generates results under the zero-shot setting. The yellow cells in Sample 1 show which data the user has reviewed and corrected.

Subsequent samples (e.g., Samples 2-3) proceed iteratively, with SCILIRE leveraging the user-corrected and verified records via dynamic sampling to generate records (Figure 9).

Once the user is satisfied, they trigger the batch phase, where all remaining documents are processed to generate records, producing a complete, curated dataset ready. This is represented by the output in Sample 4, which benefits from dynamic sampling drawn from the pool of data in Samples 1-3.

<sup>21</sup>Due to legal constraints, we are unable to provide a non-licensed user account (i.e., a demo account) that demonstrates the uploading of the given dataset due to copyright legislation in Australia.

| Models          | Zero-shot |       |                |       | ICL -10 |       |                |       | ICL -50 |       |                |       | ICL -100 |       |                |       | ICL -all |       |                |       |
|-----------------|-----------|-------|----------------|-------|---------|-------|----------------|-------|---------|-------|----------------|-------|----------|-------|----------------|-------|----------|-------|----------------|-------|
|                 | P         | R     | F <sub>1</sub> | ChrF  | P       | R     | F <sub>1</sub> | ChrF  | P       | R     | F <sub>1</sub> | ChrF  | P        | R     | F <sub>1</sub> | ChrF  | P        | R     | F <sub>1</sub> | ChrF  |
| TDMS            | 6.95      | 18.75 | 10.14          | 11.94 | 13.88   | 30.14 | 19.01          | 13.88 | 16.47   | 38.2  | 23.01          | 14.70 | 17.68    | 40.13 | 24.54          | 15.11 | 18.12    | 40.38 | <b>25.02</b>   | 15.05 |
| SciREX          | 2.8       | 5.29  | 3.66           | 7.16  | 10.67   | 18.39 | 13.51          | 10.16 | 12.61   | 20.07 | 15.49          | 10.78 | 12.53    | 21.06 | 15.71          | 11.09 | 14.52    | 24.61 | <b>18.27</b>   | 11.63 |
| MPEA            | 31.12     | 27.56 | 29.23          | 0.74  | 34.59   | 30.68 | 32.52          | 0.91  | 32.9    | 28.24 | 30.39          | 0.77  | 32.45    | 29.32 | <b>30.81</b>   | 0.74  | 31.92    | 29.45 | 30.64          | 0.74  |
| Diffusion       | 33.51     | 11.86 | 17.52          | 1.54  | 32.51   | 12.43 | <b>17.99</b>   | 1.58  | 33.13   | 11.98 | 17.59          | 1.6   | -        | -     | -              | -     | 30.35    | 12.0  | 17.20          | 1.52  |
| YSHEAY          | 2.94      | 29.51 | 5.34           | 12.32 | 4.45    | 34.41 | 7.87           | 12.96 | 4.69    | 35.96 | <b>8.30</b>    | 13.14 | 4.54     | 34.85 | 8.03           | 13.11 | 4.47     | 35.36 | 7.93           | 13.28 |
| CCRMG           | 0.96      | 18.7  | 1.82           | 11.64 | 1.31    | 23.58 | 2.48           | 13.59 | -       | -     | -              | -     | -        | -     | -              | -     | 1.53     | 25.20 | <b>2.89</b>    | 12.48 |
| Doping          | 20.66     | 4.62  | 7.55           | 10.35 | 23.65   | 9.18  | 13.23          | 14.85 | 26.19   | 10.17 | <b>14.65</b>   | 14.10 | -        | -     | -              | -     | 22.32    | 9.12  | 12.95          | 14.57 |
| MMD             | 2.98      | 0.45  | 0.78           | 1.08  | -       | -     | -              | -     | -       | -     | -              | -     | -        | -     | -              | -     | 37.5     | 4.82  | <b>8.54</b>    | 2.12  |
| MRL             | 3.54      | 1.20  | 1.80           | 1.12  | 3.95    | 1.33  | 1.99           | 1.14  | 4.22    | 1.34  | <b>2.04</b>    | 1.19  | -        | -     | -              | -     | 4.25     | 1.43  | 1.82           | 2.01  |
| PNCEExtract     | 30.61     | 31.7  | 31.14          | 9.66  | 40.4    | 40.02 | 40.21          | 10.73 | 42.71   | 42.2  | 42.46          | 11.24 | 45.82    | 44.29 | <b>45.04</b>   | 11.49 | 44.83    | 42.42 | 43.59          | 11.50 |
| PolyIE          | 11.43     | 19.04 | 14.29          | 15.67 | 15.78   | 22.61 | 18.58          | 15.79 | 16.33   | 21.69 | <b>18.64</b>   | 15.75 | -        | -     | -              | -     | 15.88    | 21.69 | 18.34          | 15.75 |
| BRENDA_enzyme   | 36.2      | 21.0  | 26.58          | 4.14  | 52.59   | 28.31 | <b>36.81</b>   | 4.94  | 49.72   | 26.21 | 34.33          | 4.77  | 51.14    | 27.94 | 36.14          | 4.91  | 53.08    | 27.92 | 36.59          | 5.06  |
| BRENDA_ribozyme | 14.08     | 10.06 | 11.74          | 2.09  | 21.17   | 15.42 | 17.84          | 2.36  | 22.95   | 16.39 | <b>19.12</b>   | 2.51  | 22.36    | 16.05 | 18.69          | 2.49  | 22.47    | 15.69 | 18.48          | 2.51  |
| OPE             | 16.33     | 35.34 | 22.33          | 5.78  | 21.15   | 45.64 | <b>28.91</b>   | 7.19  | 20.69   | 44.09 | <b>28.17</b>   | 7.16  | 17.38    | 37.76 | 23.80          | 5.97  | 16.03    | 35.65 | 22.12          | 5.88  |
| PPE             | 41.65     | 68.05 | 51.67          | 11.62 | 58.68   | 80.38 | <b>67.83</b>   | 14.07 | 55.78   | 78.93 | 65.36          | 13.72 | 53.78    | 81.38 | 64.76          | 13.83 | 53.91    | 80.57 | 64.60          | 13.81 |
| SE              | 40.50     | 33.38 | 36.6           | 11.42 | 49.16   | 37.76 | 42.71          | 12.5  | 51.61   | 43.15 | <b>47.00</b>   | 12.41 | -        | -     | -              | -     | 51.83    | 42.62 | 46.78          | 12.8  |
| AE              | 17.25     | 13.53 | 15.17          | 5.12  | 20.25   | 16.04 | <b>17.90</b>   | 10.03 | 17.77   | 14.79 | 16.14          | 8.87  | -        | -     | -              | -     | 13.56    | 10.03 | 11.53          | 9.34  |
| SuperMat        | 11.48     | 3.76  | 5.66           | 9.30  | 40.05   | 12.72 | <b>19.31</b>   | 11.06 | 33.08   | 11.06 | 16.58          | 11.23 | 34.24    | 11.16 | 16.83          | 11.32 | 34.63    | 11.39 | 17.14          | 11.06 |

Table 10: Evaluation results across datasets. Cells marked with ‘-’ indicate that the dataset does not have enough train data for evaluation. All scores are reported on a 0–100 scale, with the best F<sub>1</sub> score highlighted in **boldfaced**. LLM: GPT-5.

| Dataset         | SciSpace  |       |                |      |              |       |                |      |           |      |                |      |              |       |                |      | Elicit    |       |                |       |               |       |                |       |           |     |                |      |               |      |      |      | SciLIRE best (ICL large) |      |      |      |       |       |              |       |       |       |             |       |       |       |              |       |      |      |      |      |
|-----------------|-----------|-------|----------------|------|--------------|-------|----------------|------|-----------|------|----------------|------|--------------|-------|----------------|------|-----------|-------|----------------|-------|---------------|-------|----------------|-------|-----------|-----|----------------|------|---------------|------|------|------|--------------------------|------|------|------|-------|-------|--------------|-------|-------|-------|-------------|-------|-------|-------|--------------|-------|------|------|------|------|
|                 | Zero-shot |       |                |      | ICL (Static) |       |                |      | Zero-shot |      |                |      | ICL (Static) |       |                |      | Zero-shot |       |                |       | ICL (Dynamic) |       |                |       | Zero-shot |     |                |      | ICL (Dynamic) |      |      |      |                          |      |      |      |       |       |              |       |       |       |             |       |       |       |              |       |      |      |      |      |
|                 | P         | R     | F <sub>1</sub> | ChrF | P            | R     | F <sub>1</sub> | ChrF | P         | R    | F <sub>1</sub> | ChrF | P            | R     | F <sub>1</sub> | ChrF | P         | R     | F <sub>1</sub> | ChrF  | P             | R     | F <sub>1</sub> | ChrF  | P         | R   | F <sub>1</sub> | ChrF |               |      |      |      |                          |      |      |      |       |       |              |       |       |       |             |       |       |       |              |       |      |      |      |      |
| TDMS            | 0.0       | 0.0   | 0.0            | 2.97 | 0.0          | 0.0   | 0.0            | 4.36 | 0.0       | 0.0  | 0.0            | 5.43 | 5.0          | 2.27  | 3.13           | 7.70 | 2.4       | 11.36 | 3.97           | 9.42  | 7.94          | 22.73 | <b>11.76</b>   | 11.08 | 0.0       | 0.0 | 0.0            | 2.97 | 0.0           | 0.0  | 0.0  | 2.97 | 2.0                      | 0.74 | 1.08 | 3.73 | 12.0  | 4.44  | 6.49         | 6.80  | 2.09  | 5.19  | 2.98        | 8.09  | 12.43 | 34.07 | <b>18.22</b> | 14.21 |      |      |      |      |
| SciREX          | 0.0       | 0.0   | 0.0            | 2.97 | 0.0          | 0.0   | 0.0            | 4.36 | 0.0       | 0.0  | 0.0            | 5.43 | 5.0          | 2.27  | 3.13           | 7.70 | 2.4       | 11.36 | 3.97           | 9.42  | 7.94          | 22.73 | <b>11.76</b>   | 11.08 | 0.0       | 0.0 | 0.0            | 2.97 | 0.0           | 0.0  | 0.0  | 2.97 | 2.0                      | 0.74 | 1.08 | 3.73 | 12.0  | 4.44  | 6.49         | 6.80  | 2.09  | 5.19  | 2.98        | 8.09  | 12.43 | 34.07 | <b>18.22</b> | 14.21 |      |      |      |      |
| MPEA            | 47.06     | 7.71  | 13.26          | 0.24 | 47.06        | 7.71  | 13.26          | 0.24 | 0.0       | 0.0  | 0.0            | 0.22 | 0.0          | 0.0   | 0.0            | 0.29 | 39.13     | 42.33 | 40.67          | 0.73  | 39.59         | 45.35 | <b>42.27</b>   | 0.74  | 0.0       | 0.0 | 0.0            | 0.55 | 1.85          | 0.03 | 0.06 | 0.55 | 17.28                    | 0.27 | 0.53 | 0.96 | 23.29 | 3.98  | 6.80         | 1.41  | 27.51 | 5.17  | <b>8.71</b> | 1.59  |       |       |              |       |      |      |      |      |
| Diffusion       | 20.99     | 0.33  | 0.65           | 0.35 | 20.99        | 0.33  | 0.65           | 0.35 | 1.85      | 0.03 | 0.06           | 0.55 | 17.28        | 0.27  | 0.53           | 0.96 | 23.29     | 3.98  | 6.80           | 1.41  | 27.51         | 5.17  | <b>8.71</b>    | 1.59  | 0.0       | 0.0 | 0.0            | 5.12 | 0.0           | 0.0  | 0.0  | 5.12 | 3.33                     | 1.67 | 2.22 | 5.35 | 20.0  | 10.0  | <b>13.33</b> | 11.28 | 1.73  | 33.33 | 3.29        | 11.72 | 2.97  | 41.67 | 5.54         | 12.22 |      |      |      |      |
| YSHEAY          | 0.0       | 0.0   | 0.0            | 4.47 | 0.0          | 0.0   | 0.0            | 4.04 | 13.33     | 2.08 | 3.60           | 5.05 | 33.33        | 5.21  | 9.01           | 7.61 | 20.0      | 3.12  | 5.41           | 8.32  | 22.22         | 27.92 | <b>12.12</b>   | 9.67  | 0.0       | 0.0 | 0.0            | 3.83 | 0.0           | 0.0  | 0.0  | 3.92 | 0.0                      | 0.0  | 0.0  | 4.12 | 33.33 | 16.67 | <b>22.22</b> | 8.36  | 0.96  | 15.0  | 1.80        | 11.82 | 1.42  | 21.67 | 2.67         | 11.73 |      |      |      |      |
| CCRMG           | 0.0       | 0.0   | 0.0            | 4.47 | 0.0          | 0.0   | 0.0            | 4.04 | 13.33     | 2.08 | 3.60           | 5.05 | 33.33        | 5.21  | 9.01           | 7.61 | 20.0      | 3.12  | 5.41           | 8.32  | 22.22         | 27.92 | <b>12.12</b>   | 9.67  | 0.0       | 0.0 | 0.0            | 3.83 | 0.0           | 0.0  | 0.0  | 3.92 | 0.0                      | 0.0  | 0.0  | 4.12 | 33.33 | 16.67 | <b>22.22</b> | 8.36  | 0.96  | 15.0  | 1.80        | 11.82 | 1.42  | 21.67 | 2.67         | 11.73 |      |      |      |      |
| Doping          | 0.0       | 0.0   | 0.0            | 4.47 | 0.0          | 0.0   | 0.0            | 4.04 | 13.33     | 2.08 | 3.60           | 5.05 | 33.33        | 5.21  | 9.01           | 7.61 | 20.0      | 3.12  | 5.41           | 8.32  | 22.22         | 27.92 | <b>12.12</b>   | 9.67  | 0.0       | 0.0 | 0.0            | 3.83 | 0.0           | 0.0  | 0.0  | 3.92 | 0.0                      | 0.0  | 0.0  | 4.12 | 33.33 | 16.67 | <b>22.22</b> | 8.36  | 0.96  | 15.0  | 1.80        | 11.82 | 1.42  | 21.67 | 2.67         | 11.73 |      |      |      |      |
| MMD             | 0.0       | 0.0   | 0.0            | 0.29 | 0.0          | 0.0   | 0.0            | 0.25 | 0.0       | 0.0  | 0.0            | 0.30 | 6.25         | 0.13  | 0.26           | 0.46 | 2.98      | 0.45  | 0.78           | 1.08  | 36.18         | 4.91  | <b>8.65</b>    | 1.96  | 0.0       | 0.0 | 0.0            | 0.29 | 1.05          | 0.07 | 0.13 | 0.24 | 1.05                     | 0.07 | 0.13 | 0.23 | 0.0   | 0.0   | 0.0          | 0.40  | 4.74  | 0.31  | 0.58        | 0.64  | 4.66  | 1.07  | <b>1.75</b>  | 1.08  | 4.90 | 0.93 | 1.57 | 1.10 |
| MRL             | 1.05      | 0.07  | 0.13           | 0.24 | 1.05         | 0.07  | 0.13           | 0.23 | 0.0       | 0.0  | 0.0            | 0.40 | 4.74         | 0.31  | 0.58           | 0.64 | 4.66      | 1.07  | <b>1.75</b>    | 1.08  | 4.90          | 0.93  | 1.57           | 1.10  | 0.0       | 0.0 | 0.0            | 0.29 | 1.05          | 0.07 | 0.13 | 0.24 | 1.05                     | 0.07 | 0.13 | 0.23 | 0.0   | 0.0   | 0.0          | 0.40  | 4.74  | 0.31  | 0.58        | 0.64  | 4.66  | 1.07  | <b>1.75</b>  | 1.08  | 4.90 | 0.93 | 1.57 | 1.10 |
| PNCEExtract     | 11.67     | 1.44  | 2.56           | 1.56 | 11.67        | 1.44  | 2.56           | 1.56 | 23.33     | 2.88 | 5.13           | 4.43 | 26.67        | 3.29  | 5.86           | 5.05 | 40.43     | 23.46 | 29.69          | 9.64  | 51.19         | 26.54 | <b>34.96</b>   | 10.83 | 0.0       | 0.0 | 0.0            | 4.71 | 0.0           | 0.0  | 0.0  | 4.89 | 0.0                      | 0.0  | 0.0  | 4.0  | 13.33 | 0.38  | 0.74         | 7.33  | 10.56 | 14.04 | 12.05       | 13.60 | 17.24 | 20.15 | <b>18.58</b> | 15.11 |      |      |      |      |
| PolyIE          | 0.0       | 0.0   | 0.0            | 4.71 | 0.0          | 0.0   | 0.0            | 4.89 | 0.0       | 0.0  | 0.0            | 4.0  | 13.33        | 0.38  | 0.74           | 7.33 | 10.56     | 14.04 | 12.05          | 13.60 | 17.24         | 20.15 | <b>18.58</b>   | 15.11 | 0.0       | 0.0 | 0.0            | 0.55 | 5.38          | 0.17 | 0.33 | 0.64 | 6.92                     | 0.22 | 0.42 | 0.59 | 22.31 | 0.70  | 1.35         | 1.51  | 47.03 | 27.04 | 34.34       | 4.37  | 65.2  | 37.28 | <b>47.44</b> | 5.56  |      |      |      |      |
| BRENDA_enzyme   | 0.77      | 0.02  | 0.05           | 0.61 | 5.38         | 0.17  | 0.33           | 0.64 | 6.92      | 0.22 | 0.42           | 0.59 | 22.31        | 0.70  | 1.35           | 1.51 | 47.03     | 27.04 | 34.34          | 4.37  | 65.2          | 37.28 | <b>47.44</b>   | 5.56  | 0.0       | 0.0 | 0.0            | 0.55 | 5.38          | 0.17 | 0.33 | 0.64 | 6.92                     | 0.22 | 0.42 | 0.59 | 22.31 | 0.70  | 1.35         | 1.51  | 47.03 | 27.04 | 34.34       | 4.37  | 65.2  | 37.28 | <b>47.44</b> | 5.56  |      |      |      |      |
| BRENDA_ribozyme | 0.0       | 0.0   | 0.0            | 0.55 | 0.0          | 0.0   | 0.0            | 0.51 | 8.24      | 1.11 | 1.96           | 0.48 | 18.24        | 2.46  | 4.34           | 0.86 | 26.39     | 25.68 | 26.03          | 2.24  | 30.75         | 31.16 | <b>30.95</b>   | 2.67  | 0.0       | 0.0 | 0.0            | 0.55 | 5.38          | 0.17 | 0.33 | 0.64 | 6.92                     | 0.22 | 0.42 | 0.59 | 22.31 | 0.70  | 1.35         | 1.51  | 47.03 | 27.04 | 34.34       | 4.37  | 65.2  | 37.28 | <b>47.44</b> | 5.56  |      |      |      |      |
| OPE             | 24.44     | 12.87 | 16.86          | 1.02 | 30.0         | 15.79 | 20.69          | 0.99 | 15.56     | 8.19 | 10.73          | 1.85 | 24.44        | 12.87 | 16.86          | 2.92 | 13.22     | 36.18 | <b>19.37</b>   | 4.84  | 10.66         | 34.21 | 16.25          | 4.40  | 0.0       | 0.0 | 0.0            | 0.55 | 5.38          | 0.17 | 0.33 | 0.64 | 6.92                     | 0.22 | 0.42 | 0.59 | 22.31 | 0.70  | 1.35         | 1.51  | 47.03 | 27.04 | 34.34       | 4.37  | 65.2  | 37.28 | <b>47.44</b> | 5.56  |      |      |      |      |
| PPE             | 10.0      | 3.7   | 5.41           | 1.75 | 0.0          | 0.0   | 0.0            | 1.21 | 3.33      | 1.23 | 1.80           | 2.11 | 23.33        | 8.64  | 12.61          | 3.55 | 39.39     | 64.20 | 48.83          | 11.61 | 52.5          | 77.78 | <b>62.69</b>   | 14.09 | 0.0       | 0.0 | 0.0            | 0.55 | 5.38          | 0.17 | 0.33 | 0.64 | 6.92                     | 0.22 | 0.42 | 0.59 | 22.31 | 0.70  | 1.35         | 1.51  | 47.03 | 27.04 | 34.34       |       |       |       |              |       |      |      |      |      |

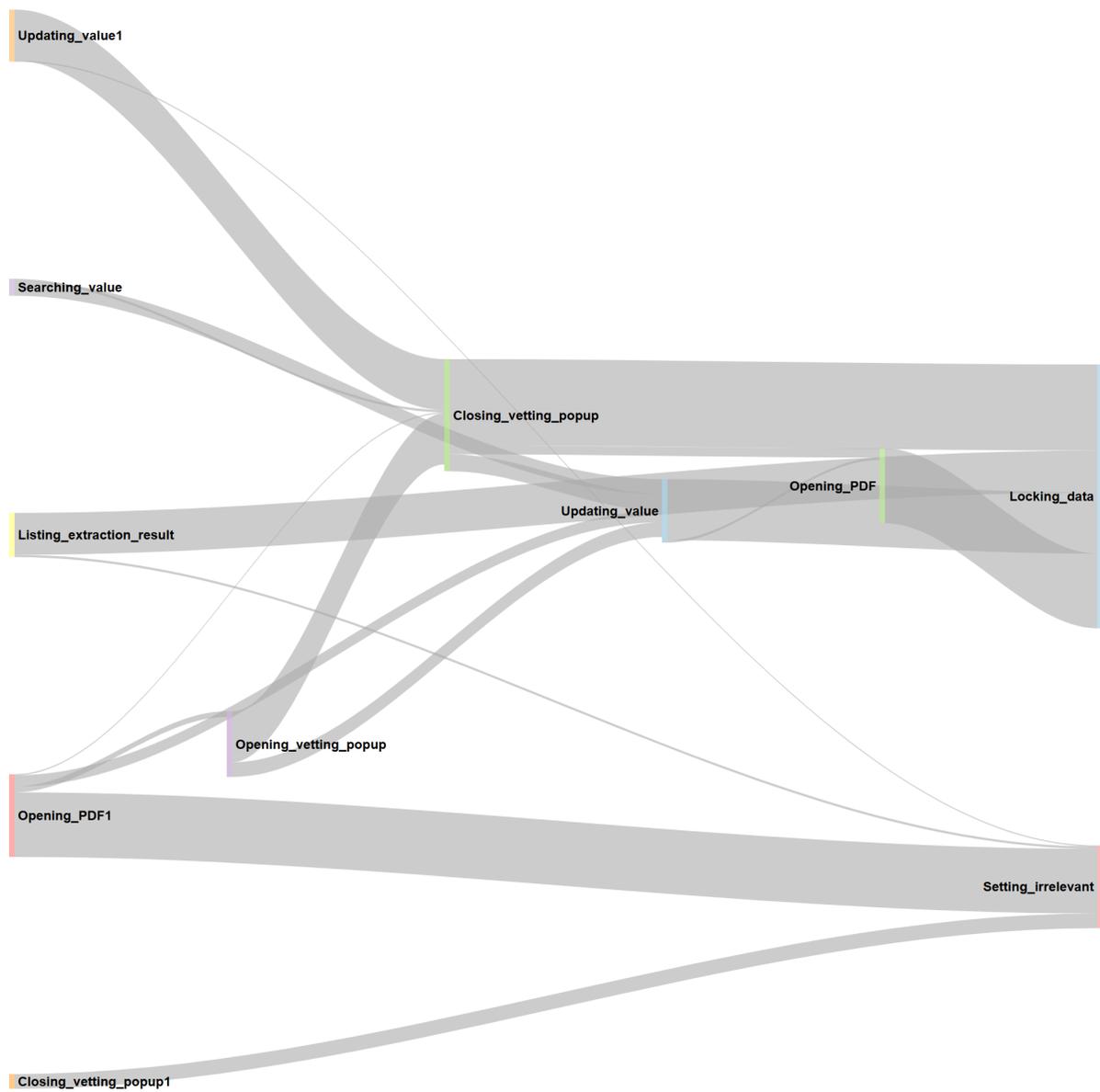


Figure 4: Interaction flows within SCILIRE for the early adopter trials. Results show that most data acceptance (*locking\_data*) or rejections (*setting\_irrelevant*) occur via a data verification step (either checking the provenance data or the original source PDF). “Vetting popup” here refers to the verification support tools. “Updating\_value” refers to human editing and manual data curation activities. Actions with “1” at the end are used to eliminate cycles for the purposes of visualisation with a Sankey diagram.

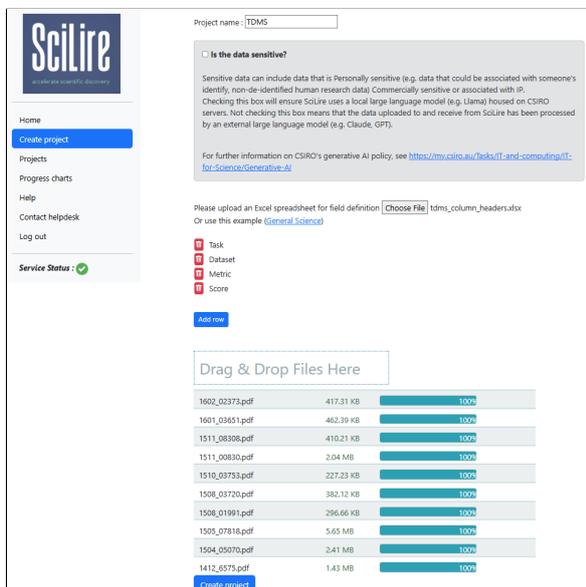


Figure 5: A screenshot of SciLIRE for project creation.



Home  
Create project  
**Projects**  
Progress charts  
Help  
Contact helpdesk  
Log out

Service Status : ✔

| # | Project name | Date created         | Users       | Status        | Actions                                                                   |
|---|--------------|----------------------|-------------|---------------|---------------------------------------------------------------------------|
| 1 | TDMS         | 12/1/2025 7:29:22 PM | demo_access | Ready (26/26) | <a href="#">Manage</a> <a href="#">Upload PDFs</a> <a href="#">Delete</a> |

Figure 6: A listing of a user's projects.



Home  
Create project  
**Projects**  
Progress charts  
Help  
Contact helpdesk  
Log out

Service Status : ✔

**Project Name : TDMS**

Project Overview | [Update table structure](#) | [Extract data](#) | [Project members](#)

**Metadata**      **Data**

Project name: TDMS  
Date created: 12/1/2025 7:29:22 PM

Previous extractions

| Sample#                   | Date                  | Doc count | Files                                                                                                                                                                                                                                                                                                                                                                                      | Training status                 | Status |
|---------------------------|-----------------------|-----------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------|--------|
| <a href="#">Sample #4</a> | 12/2/2025 10:04:38 AM | 19        | f2205c56378e715d8d12c521d045c0756a76.pdf<br>P18_4013.pdf<br>view.pdf<br>C18_1139.pdf<br>1603_01354.pdf<br>P18_2038.pdf<br>N18_5012.pdf<br>1603_01360.pdf<br>1810_13097.pdf<br>1703_06345.pdf<br>1802_05365.pdf<br>1504_05070.pdf<br>1709_04109.pdf<br>1508_03720.pdf<br>27496a2ee337db705e7c611dea1fd8e6f41437c2.pdf<br>1602_02373.pdf<br>N18_1127.pdf<br>1809_08370.pdf<br>1806_03489.pdf | Checked pdfs use                | Done   |
| <a href="#">Sample #3</a> | 12/1/2025 7:50:14 PM  | 2         | 1412_6575.pdf<br>1508_01991.pdf                                                                                                                                                                                                                                                                                                                                                            | Checked pdfs use                | Done   |
| <a href="#">Sample #2</a> | 12/1/2025 7:46:23 PM  | 3         | 1601_03651.pdf<br>1505_07818.pdf<br>1511_08308.pdf                                                                                                                                                                                                                                                                                                                                         | Checked pdfs use                | Done   |
| <a href="#">Sample #1</a> | 12/1/2025 7:37:34 PM  | 2         | 1511_00830.pdf<br>1510_03753.pdf                                                                                                                                                                                                                                                                                                                                                           | Baseline (no checked pdfs used) | Done   |

[Show all results](#)

Users: demo\_user, necva.bolucu@data61.csiro.au, stephen.wan@csiro.au, demo\_access [Add user](#)

AI settings: Data sensitivity : no

Fields: Task, Dataset, Metric, Score [Update field definition](#)

Static: active

Figure 7: The “samples” representing the table output from the iterative workflow.

| Checked PDF                         | Relevant                            | File Name      | Title                                                      | Author                                                                | Publication Date | Task                    | Dataset                        | Metric                           | Score                                                 |
|-------------------------------------|-------------------------------------|----------------|------------------------------------------------------------|-----------------------------------------------------------------------|------------------|-------------------------|--------------------------------|----------------------------------|-------------------------------------------------------|
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1510_03753.pdf | Improved Deep Learning Baselines for Ubuntu Corpus Dialogs | Rudolf Kadlec, Martin Schmid, Jan Kleindienst                         | 2015-11-03       | retrieval-based_chatbot | Ubuntu Corpus                  | R_2@1                            | 89.5                                                  |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1510_03753.pdf | Improved Deep Learning Baselines for Ubuntu Corpus Dialogs | Rudolf Kadlec, Martin Schmid, Jan Kleindienst                         | 2015-11-03       | retrieval-based_chatbot | Ubuntu Corpus                  | R_10@1                           | 63                                                    |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1510_03753.pdf | Improved Deep Learning Baselines for Ubuntu Corpus Dialogs | Rudolf Kadlec, Martin Schmid, Jan Kleindienst                         | 2015-11-03       | retrieval-based_chatbot | Ubuntu Corpus                  | R_2@1                            | 89.5                                                  |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1510_03753.pdf | Improved Deep Learning Baselines for Ubuntu Corpus Dialogs | Rudolf Kadlec, Martin Schmid, Jan Kleindienst                         | 2015-11-03       | retrieval-based_chatbot | Ubuntu Corpus                  | R_10@1                           | 63                                                    |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | sentiment_analysis      | Multi-Domain Sentiment Dataset | Accuracy on DVD                  | 76.57                                                 |
| <input checked="" type="checkbox"/> | <input type="checkbox"/>            | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | sentiment_analysis      | Multi-Domain Sentiment Dataset | Accuracy on Books                | 73.4                                                  |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | sentiment_analysis      | Multi-Domain Sentiment Dataset | Accuracy on Electronics          | 80.53                                                 |
| <input checked="" type="checkbox"/> | <input type="checkbox"/>            | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | sentiment_analysis      | Multi-Domain Sentiment Dataset | Accuracy on Kitchen              | 82.93                                                 |
| <input checked="" type="checkbox"/> | <input type="checkbox"/>            | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | sentiment_analysis      | Multi-Domain Sentiment Dataset | Accuracy on Average              | 78.36                                                 |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | sentiment_analysis      | Amazon reviews dataset         | Accuracy on label y              | Higher accuracy on 9 out of 12 tasks compared to DANN |
| <input checked="" type="checkbox"/> | <input type="checkbox"/>            | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | Domain adaptation       | Amazon reviews dataset         | Accuracy on domain information s | Towards random chance (0.5)                           |
| <input checked="" type="checkbox"/> | <input type="checkbox"/>            | 1511_00830.pdf | THE VARIATIONAL FAIR AUTOENCODER                           | Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, Richard Zemel | 2017-08-10       | Learning invariant      | Extended Yale B                | Accuracy on label y              | Improved from 78%                                     |

Figure 8: The table in Sample 1, which has been edited and curated by the user.

| Checked PDF              | Relevant                            | File Name      | Title                                              | Author                        | Publication Date | Task                     | Dataset                           | Metric   | Score                                                             |
|--------------------------|-------------------------------------|----------------|----------------------------------------------------|-------------------------------|------------------|--------------------------|-----------------------------------|----------|-------------------------------------------------------------------|
| <input type="checkbox"/> | <input checked="" type="checkbox"/> | 1508_01991.pdf | Bidirectional LSTM-CRF Models for Sequence Tagging | Zhiheng Huang, Wei Xu, Kai Yu | 2015-08-09       | part_of_speech_tagging   | VLSP 2013 POS tagging shared task | Accuracy | 95.06                                                             |
| <input type="checkbox"/> | <input type="checkbox"/>            | 1508_01991.pdf | Bidirectional LSTM-CRF Models for Sequence Tagging | Zhiheng Huang, Wei Xu, Kai Yu | 2015-08-09       | chunking                 | CoNLL2000 dataset                 | F1 score | 94.13 (Bi-LSTM-CRF without Senna), 94.46 (Bi-LSTM-CRF with Senna) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> | 1508_01991.pdf | Bidirectional LSTM-CRF Models for Sequence Tagging | Zhiheng Huang, Wei Xu, Kai Yu | 2015-08-09       | named_entity_recognition | VLSP 2016 NER shared task         | F1 score | 86.48                                                             |

Figure 9: A screenshot from the pilot phase (Sample 3) of the AI-augmented curation workflow in the demo.

# xLM: A Python Package for Non-Autoregressive Language Models

Dhruvesh Patel Durga Prasad Maram\* Sai Sreenivas Chintha\* Benjamin Rozonoyer  
Andrew McCallum

University of Massachusetts Amherst

{dhruveshpate, dmaram, saisreenivas, brozonoyer, mccallum}@umass.edu

 xlm-core, xlm-models  Docs  Video  Code

## Abstract

In recent years, there has been a resurgence of interest in non-autoregressive text generation in the context of general language modeling. Unlike the well-established autoregressive language modeling paradigm, which has a plethora of standard training and inference libraries, implementations of non-autoregressive language modeling have largely been bespoke making it difficult to perform systematic comparisons of different methods. Moreover, each non-autoregressive language model typically requires its own data collation, loss, and prediction logic, making it challenging to reuse common components. In this work, we present the xLM python package, which is designed to make implementing small non-autoregressive language models faster. With a secondary goal of providing a suite of small pre-trained models (through a companion xlm-models package) that can be used by the research community.

## 1 Introduction

Autoregressive language models (ARLMs), which generate text sequentially from left to right by adding one token at a time, are well established with a plethora of standard training and inference libraries (Wolf et al., 2020; OLMo et al., 2024). However, recently, there has been a resurgence of research interest in non-autoregressive text generation due to its potential for faster inference speeds and better generation quality for certain tasks. Unlike left-to-right generation, non-autoregressive text generation can be done in many ways, for example, using masked diffusion language models (Sahoo et al., 2024), Gaussian diffusion language models (Gulrajani and Hashimoto, 2023), insertion language models (Patel et al., 2025), edit-based language models (Havasi et al., 2025), etc. Moreover, each method typically requires its own data

collation, loss, and prediction logic, making it challenging to reuse common components across different methods. The rapidly expanding landscape of these methods has led to many bespoke implementations, making it extremely difficult to compare them systematically. In this work, we present the xLM python package, which aims to provide a unified framework for developing and comparing small non-autoregressive language models. xLM uses Pytorch (Paszke et al., 2019) as the deep learning framework, Pytorch Lightning (Falcon and The PyTorch Lightning team, 2019) for training utilities, and Hydra (Yadan, 2019) for configuration management. xLM is designed to make implementing small non-autoregressive language models faster without sacrificing flexibility.

The rest of the paper is organized as follows. Section 2 discusses the core design principles of xLM. In section 3, we discuss how xLM serves a unique purpose in the landscape of LLM libraries. Section 4 presents a high-level overview of the three core components of xLM, followed by section 5 that provides a step-by-step demonstration of how one would implement a new language model using xLM. Finally, Section 6 presents a set of benchmarking results where we implement three models in xLM to reproduce known results.

## 2 Design Principles

xLM follows the principle of maximal independence. The core library provides access to a small number of shared components, which are designed to be model independent, and can be used by any kind of language model. Each model implementation lives in its own folder/package and is completely independent of other models. This allows researchers to keep their model code clean, self-contained, and easy to share. It also allows them to use their model outside the xLM framework without refactoring. Maximal independence is achieved by

\*Equal contribution second authors.

following design choices.

**Composition over inheritance.** The maximal independence is achieved by using composition over inheritance (Gamma et al., 1995), wherein the core components delegate model specific logic to the specific model instance. For example, as shown in Figure 1, the `DataModule` carries a collection of `DatasetManager` instances, one for each dataset, and the creation of dataloaders is delegated to the respective `DatasetManager` instance depending on the stage (train, val, test, or predict). Similarly, the `Harness` carries instances of `Model`, `LossFunction`, `Predictor`, to which it delegates the model specific logic for forward pass, loss computation, and generation respectively. Moreover, one can swap out one or more of the four components with a different but compatible implementation without having to change all four.

**Copy over branching.** This principle goes against the common wisdom of not copying code. However, in case of research codebases, copying reduces code complexity and helps increase the speed of development (Wolf et al., 2020). It naturally creates independence. Moreover, it also allows templating the process of creating a brand new model which helps rapid prototyping by humans as well as LLMs (Li et al., 2025).

**Arbitrary code injection.** Python is a highly dynamic language, which allows one to inject arbitrary code at runtime. In production and public facing codebases, this creates a security risk, but this flexibility is a boon for research codebases, as it allows rapid prototyping and experimentation. As we will discuss in section 4.3, Hydra (Yadan, 2019) provides a powerful mechanism to inject arbitrary code at runtime by allowing one to fill a specific *slot* with an instance of any class.

### 3 Related Work

Due to the rapid development of LLMs, there are many libraries for training (Wolf et al., 2020; von Werra et al., 2020; Lightning-AI, 2023; OLMo et al., 2024) and inference using auto-regressive LLMs (Kwon et al., 2023). On the other hand, there are only a handful of python libraries that support non-autoregressive sequence modeling like FairSeq (Ott et al., 2019), and AllenNLP (Gardner et al., 2017). Moreover, even these libraries are no longer actively maintained and do not support non-autoregressive language modeling. To the best

of our knowledge, xLM is the only library that supports fast prototyping of small non-autoregressive language models, and is geared towards providing a suite of reference implementations for up and coming non-autoregressive language modeling methods.

## 4 Core Components of xLM

In this section, we will discuss the the `Harness`, the `DataModule`, and the configuration management, which together handle the execution flow of all the supported workflows (section 5.7) like training, evaluation, prediction and debugging. In most use cases, these components need not be touched by the user, removing the need for most of the boilerplate code.<sup>1</sup>

### 4.1 DataModule

The base `TextDataModule`, which builds on top of `Lightning DataModule`, provides a generic, model and task agnostic interface for managing arbitrary number of text datasets.<sup>2</sup> This is achieved by using one `DatasetManager` per dataset as shown in fig. 1. Each `DatasetManager` instance is responsible for managing the complete lifecycle of a single dataset, including downloading, preprocessing, caching and managing the data collator and data-loader options. It has slots for the following components that allow injecting custom logic:

- `Dataset` which could point to a HuggingFace dataset or a custom dataset.
- `Collator` and `Preprocessor`, both of which can depend on the model type as well as the task.

Complete flexibility and independence is achieved by allowing a many-to-many mapping between the `DatasetManager` instances and the workflow stages (train, val, test, predict), wherein a single `DatasetManager` instance can be mapped to multiple workflow stages, and vice versa. This ensures that only a single copy of the dataset is loaded into memory but if needed it can be used at multiple places in the workflow.

The core implementation of `DatasetManager` supports all common training strategies for small models: single-node single-GPU, single-node multi-GPU and multi-node multi-GPU, with map style and iterable style dataset support for each.

<sup>1</sup>xLM also has some useful additional features described in appendix F.

<sup>2</sup>See <https://lightning.ai/docs/pytorch/stable/data/datamodule.html>. The interface is not specific to text datasets, and can be used for any kind of sequence datasets.

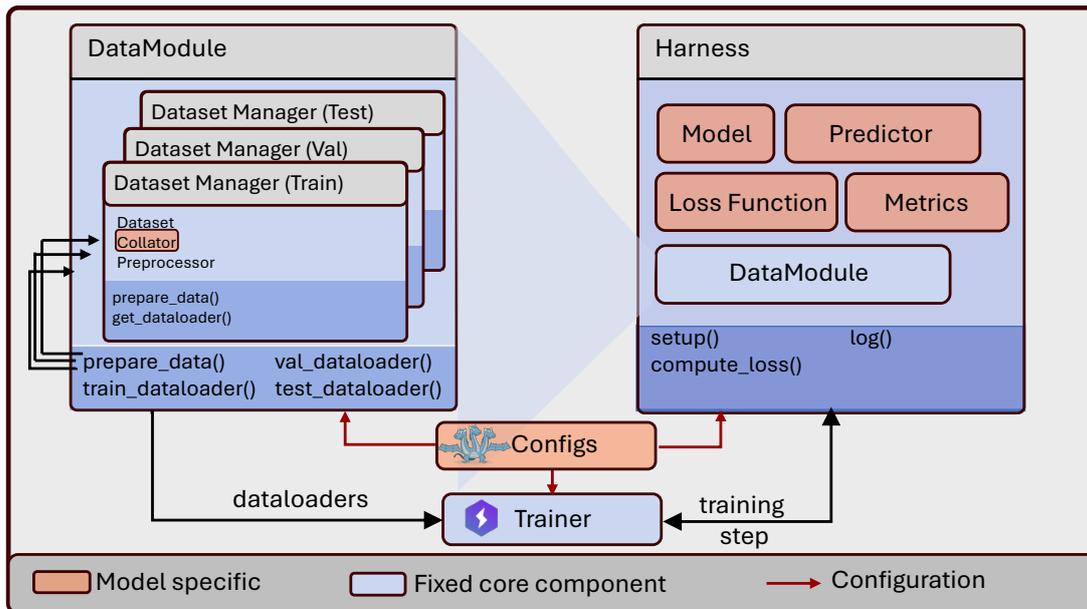


Figure 1: Overview of xLM design. It consists of two classes of components: the core components (Harness and DataModule) and the model-specific components, whose implementations depend on the model logic. These components are defined in the configuration files managed by Hydra (see fig. 3), enabling arbitrary component swapping. The Harness component is responsible for instantiating all components (model, loss, predictor, etc.) and delegating their respective functionalities. The DataModule component manages multiple datasets across workflow stages using DatasetManager objects, each handling a dataset and an appropriate Collator.

The user simply needs to provide the respective arguments in the config file.

## 4.2 Harness

The Harness is the main class that inherits from the PyTorch lightning’s LightningModule<sup>3</sup>, and is responsible for instantiating all the components like the model, loss function, predictor, etc., based on the configuration files. As shown in fig. 1, the Harness has slots (attributes) for all the core components, and it delegates the model specific logic to the respective components’ methods. Inheriting from the LightningModule allows us to use all the features of PyTorch lightning, such as logging, checkpointing, saving, etc., and also allows us to use the LightningTrainer. See appendix C for more details.

## 4.3 Configuration Management

In xLM, the configuration files have two roles. First, like any other configuration system (e.g. Python’s ArgParse), it allows the user to specify various parameter values that dictate the behavior of the system. Second, through the use of `hydra.utils.instantiate`, it allows swapping out entire components directly from the configuration

file, without changing a single line of python code. This enables arbitrary code injection at runtime. Hydra configs themselves can be arbitrarily nested, and one config file can be referred in another config file, the composition of which is automatically taken care of by Hydra. This enables a modularization of the configuration files themselves.<sup>4</sup>

## 5 Demonstration

In this section, we will walk through, step by step, the process of implementing a new language model using the xLM library. In order to demonstrate the flexibility of the library, we pick a non-standard language modeling paradigm for this demonstration. Specifically, we will implement the Insertion Language Model (ILM) of Patel et al. (2025), which generates text by iteratively inserting tokens in the existing sequence by selecting the position and the vocabulary item to insert. In order to keep the demonstration simple, we will use the synthetic seq2seq task of generating a path on a star shaped graph as the task (Patel et al., 2025) on the StarEasy dataset<sup>3</sup>. We also provide a demonstration of training ILM on the LM1B corpus (Chelba et al., 2014) in Appendix A.

<sup>3</sup>[https://lightning.ai/docs/pytorch/stable/common/lightning\\_module.html](https://lightning.ai/docs/pytorch/stable/common/lightning_module.html)

<sup>4</sup>Please refer to the Hydra documentation for more details <https://hydra.cc/docs/intro/>.

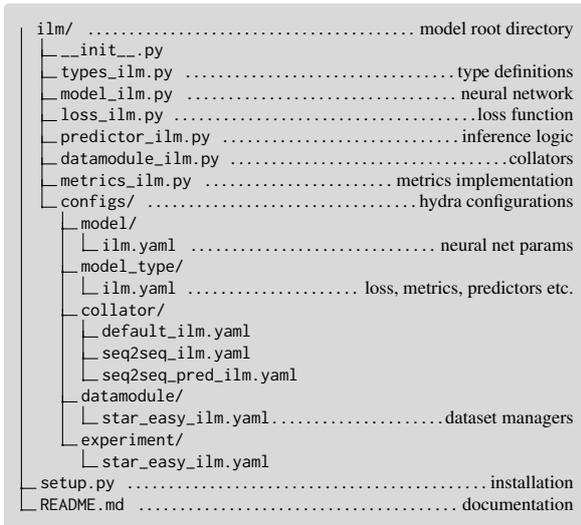


Figure 2: Directory structure generated by the scaffolding script.

To construct a new model, we need to create a fresh directory for the model, which will contain the implementations of the model specific components: Model, LossFunction, Predictor, and Collators. It will also contain their respective configuration files. In order to make the process of creating a new model easier, we provide a scaffolding script that automatically generates the necessary files. It can be invoked by executing `xlm-scaffold ilm` which will generate the directory structure as shown in fig. 2.

The scaffold files already contain placeholder empty class and function declarations. Next we will walk through the process of adding implementation for each of the python files shown in fig. 2. Finally, we will show how to reference these components configuration files.

## 5.1 Model

In `model_ilm.py`, we define the neural network backbone for the model which inherits from `torch.nn.Module` and implements the `forward()` method. We have the complete freedom to decide the arguments that the `forward()` method takes.

```

from xlm.model import Model
from xlm.modules.rotary_transformer import RotaryTransformerLayer

class ILMModel(torch.nn.Module, Model):
    ...
    def forward(self, input_ids, attention_mask,...):
        self.encoder_layer = RotaryTransformerLayer(
            d_model,
            nhead,
            dim_feedforward,
            dropout,
            activation,
            layer_norm_eps,
        )

```

```

...
return vocab_logits, stopping_logits

```

Many of the xLM modules detailed in section F.1 can be used for building the architecture like the use of `RotaryTransformerLayer` here for the encoder. After defining the model class, the constructor arguments needed for default instantiation along with the fully qualified class path are stored in the config `configs/model/ilm.yaml` file:

```

# @package _global_
model:
  _target_: ilm.model_ilm.ILMModel
  ...

```

## 5.2 LossFunction

The LossFunction is a callable that takes in a batch of inputs and returns a dictionary of values, with a mandatory "loss" key and any other optional values that we want to log. In ILM, the loss function computes two components: (1) a cross-entropy loss over only the positions where tokens were dropped, and (2) a binary classification loss that predicts whether input sequence is complete.

```

from xlm.harness import LossFunction
from types_ilm import ILMBatch, ILMLossDict

class ILMLoss(LossFunction[ILMBatch, ILMLossDict]):
    def loss_fn(self, batch: ILMBatch, ...) -> ILMLossDict:
        vocab_logits, stopping_logits = self.model(**batch)

        vocab_logit_loss = masked_cross_entropy(vocab_logits,
        batch["target_ids"])

        stopping_loss = binary_cross_entropy(stopping_logits)
        return {"loss": vocab_logit_loss + stopping_loss, ...}

```

The default loss parameters are stored in the config `configs/model_type/ilm.yaml` file under the `loss` key as shown in fig. 5.

## 5.3 Data Pipeline

To setup the data pipeline, we just need to configure the dataset and implement model specific collators for each dataset, and add the configuration files for the DatasetManagers and Collators. The xLM comes with some synthetic seq2seq, and language modeling datasets preconfigured. This includes the synthetic StarEasy dataset.<sup>5</sup> For this demonstration, we will use the preconfigured StarEasy dataset, wherein each example consists of a prompt and a target path. The prompt contains an edge list (in random order) of a star graph and the start and end nodes. The target contains the gold path from the start node to the end node (Patel et al., 2025). An example prompt and target is shown in fig. 4.

<sup>5</sup>See appendix D.1 for learning how to configure a new dataset, and appendix G for the list of preconfigured tasks.

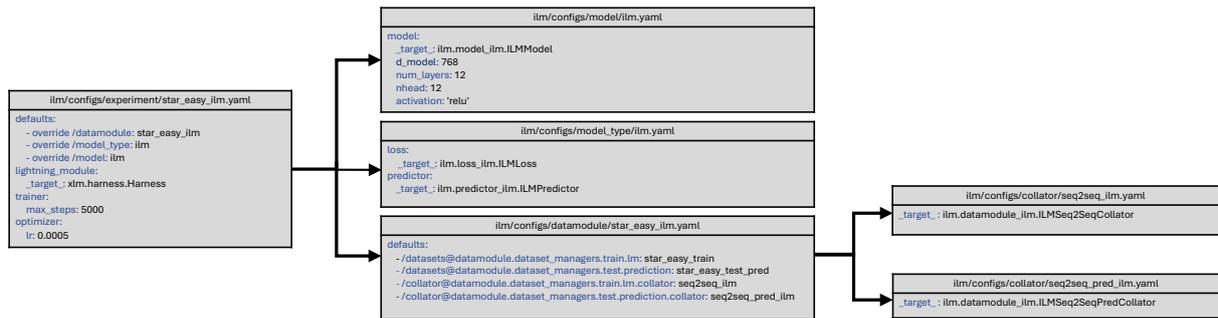


Figure 3: Configuration tree for a typical experiment (e.g. for ILM for a seq2seq planning task on the StarEasy dataset). The `experiment` config is at the root of the nesting structure, contains global parameters, and composes component configs (`model`, `model_type`, and `datamodule`). The `model/ilm.yaml` file stores the parameters for the model class. The `model_type/ilm.yaml` file contains the information needed to instantiate the loss function, predictor, and metric components. The `datamodule/star_easy_ilm.yaml` composes the configs of the DatasetManagers and Collators (here, StarEasy and seq2seq collators). Note: Only partial entries are shown in the figure for brevity.

```
edge_list = [[1, 5], [1, 7], [7, 9]]
source = 1
goal = 9
path = [1, 7, 9]
sequence = "CLS 1 7 1 5 7 9 1 9 BOS 1 7 9 PAD PAD"
```

Figure 4: An example of a prompt and target for the StarEasy dataset.

## 5.4 Collators

A Collator is a callable that takes in a list of raw examples from the dataset and returns a batch to be fed to the LossFunction or the Predictor. Typically, a Collator’s implementation depends on the model type and the type of task, e.g. the collator for a seq2seq task will be different from that of unconditional language modeling task. Moreover, a collator for loss computation will be different from that used for prediction. For the synthetic seq2seq task of star graphs, we need to implement two collators: one for training and one for prediction. The training collator randomly drops tokens from the gold path and places the dropped tokens under the `"target_ids"` key in the batch. The prediction collator only keep the prompt under the `"input_ids"` key in the batch leaving all the tokens in the path to be predicted by the model.

```
from xlm.datamodule import Collator
from types_ilm import ILMBatch

class ILMSeq2SeqCollator(Collator):
    def __call__(self, examples: List) -> ILMBatch:
        prefix = prepare_prefix_ids([e["prompt_ids"] for e in
examples])
        suffix = drop_tokens([e["input_ids"] for e in
examples])
        return {
            "input_ids": torch.cat([prefix["input_ids"],
suffix["input_ids"]], dim=1),
            "target_ids": suffix["target_ids"],
            "n_drops": suffix["n_drops"]
        }
```

```
class ILMSeq2SeqPredCollator(ILMSeq2SeqCollator):
    def __call__(self, examples: List) -> ILMBatch:
        prefix = prepare_prefix_ids([e["prompt_ids"] for e in
examples])
        target_ids =
prepare_target_ids_for_test([e["input_ids"] for e in
examples])
        return {
            "input_ids": prefix["input_ids"],
            "target_ids": target_ids["target_ids"],
            "n_drops": None
        }
```

The collators are configured by adding `configs/collator/seq2seq_ilm.yaml`:

```
_target_: ilm.datamodule_ilm.ILMSeq2SeqCollator
```

and `configs/collator/seq2seq_pred_ilm.yaml`:

```
_target_: ilm.datamodule_ilm.ILMSeq2SeqPredCollator
```

Finally, the entire data pipeline is configured by adding `configs/datamodule/star_easy_ilm.yaml` config file as shown in fig. 6.

## 5.5 Predictor

The Predictor is a callable that implements the inference loop for the model. In ILM, the Predictor implements iterative infilling by repeatedly sampling the location of insertion and the vocabulary item to insert till the stopping classification head indicates that generation should stop.

```
from xlm.harness import Predictor
from types_ilm import ILMBatch, ILMPredictionDict

class ILMPredictor(Predictor[ILMBatch, ILMPredictionDict]):
    def predict(self, batch: ILMBatch, ...) ->
ILMPredictionDict:
        ...
        return {"text": decoded_texts, "history":
generation_history}
```

xLM automatically invokes the predictor during evaluation and handles decoding and metric computation. The predictor class path is stored in the

config `configs/model_type/ilm.yaml` file under the `predictor` key as shown in fig. 5

## 5.6 Experiment Configuration

Once the individual components are implemented and configured as elaborated previously, the hierarchical configuration tree is formed through the *experiment config*. Figure 3 depicts the configuration tree for the ILM model on the StarEasy dataset. The arrows in the figure indicate the nesting structure: the main experiment config `experiments/star_easy_ilm.yaml` refers to the model `model/ilm.yaml`, `model_type` `model_type/ilm.yaml`, and the datamodule `datamodule/star_easy_ilm.yaml`, which in turn refers to the the collators.

## 5.7 Workflows

**Model Discovery** After creating the experiment configuration, we are ready to use the model. However, before we can do that we need to make the model *discoverable*. There are two ways to do that:

1. **Python Package:** The model can be installed as a standalone python package, followed by setting the environment variable `XML_MODELS_PACKAGES` as “:” separated list of installed model packages, e.g. `XML_MODELS_PACKAGES=ilm:mlm`.
2. **Directory:** Alternatively, one can simply set `XML_MODELS_PATH` to point to the parent directory that contains the model folder, e.g. if the model is located at `/path/to/ilm`, then `XML_MODELS_PATH=/path/to`.

Once the model is discoverable, xLM provides three main workflows: training, evaluation, and generation. They can be run using the following command

```
$ xlm job_type=[JOB_TYPE] job_name=[JOB_NAME]
experiment=[CONFIG_PATH]
```

The `job_type` argument can be one of `train`, `eval` and `generate`. The name of the experiment config file (without the `.yaml` extension) containing all necessary overrides is given under the `experiment` option. xLM also provides a group of debug settings that can be used to perform a quick debug run that tries to overfit on a single batch of data. This can be used by appending the `debug=overfit` option to the command. Finally, the model code can be packaged as a standalone python package.

Table 1: Benchmark performance on planning seq2seq task on star graphs. The columns represent token and sequence accuracies for each model.

| Model | Easy  |       | Medium |       | Hard |       |
|-------|-------|-------|--------|-------|------|-------|
|       | Seq.  | Token | Seq.   | Token | Seq. | Token |
| ARLM  | 33.1  | 81.7  | 77.2   | 82.1  | 25.2 | 43.7  |
| ILM   | 100.0 | 100.0 | 100.0  | 100.0 | 97.5 | 98.2  |
| MLM   | 100.0 | 100.0 | 83.1   | 98.0  | 25.3 | 79.6  |
| MDLM  | 100.0 | 100.0 | 36.5   | 90.6  | 21.0 | 54.9  |

```
# Packaging
python setup.py sdist bdist_wheel
# Installation
pip install ilm
```

Post training, the model weights can be extracted and uploaded to the Hugging Face model hub (see appendix F for the details).

## 6 Benchmarks

xLM is a framework aimed at making research prototyping of small non-autoregressive language models easier. The purpose of this section is to show that the library can be used to reproduce known results for small non-autoregressive language models. As representative tasks, we pick one seq2seq task and one unconditional language modeling task. Specifically, we reproduce the results on the synthetic seq2seq of path finding on star graphs and unconditional language modeling tasks on LM1B. **Synthetic Planning Tasks** We use the hyperparameters reported in for training auto-regressive model (ARLM), masked diffusion model (MDLM) (Sahoo et al., 2024), masked language model (MLM) and insertion language model (ILM) on the three variants of the synthetic planning tasks (Patel et al., 2025). The results (table 1) are within 2% of the reported results in the original papers.

**Language Modeling** We train a 12 layer transformer as ARLM, MDLM, and ILM, respectively, on the LM1B corpus (Chelba et al., 2014). In order to make the results comparable, we use negative loglikelihood under Llama 3.2 8B model as the metric. The results are reported in table 2, which are close to the results reported in Patel et al. (2025) for the same settings.

## 7 Conclusion and Future Work

We presented a modular python package that makes prototyping small non-autoregressive language models easier. We plan on adding refer-

Table 2: Benchmark performance on unconditional language modeling on LM1B. The rows represent negative loglikelihood (under Llama 3.2) and entropy of the generated sequences for each model.

|         | corpus | ARLM | MDLM | ILM  |
|---------|--------|------|------|------|
| NLL     | 3.71   | 3.94 | 4.81 | 4.72 |
| Entropy | 3.08   | 3.12 | 3.70 | 2.81 |

ence implementations and benchmark more models (Havasi et al., 2025) as well as add support for non-text sequence generation tasks like molecule generation (see appendix H for details).

## References

- Gregor Bachmann and Vaishnavh Nagarajan. 2024. [The pitfalls of next-token prediction](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 2296–2318. PMLR.
- Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2014. [One billion word benchmark for measuring progress in statistical language modeling](#). *Preprint*, arXiv:1312.3005.
- Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. 2024. [Flex attention: A programming model for generating optimized attention kernels](#). *Preprint*, arXiv:2412.05496.
- William Falcon and The PyTorch Lightning team. 2019. [PyTorch Lightning](#).
- Erich Gamma, Richard Helm, Ralph Johnson, and John Vlissides. 1995. *Design patterns: elements of reusable object-oriented software*. Addison-Wesley Longman Publishing Co., Inc., USA.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2017. [Allennlp: A deep semantic natural language processing platform](#).
- Ishaan Gulrajani and Tatsunori Hashimoto. 2023. [Likelihood-based diffusion language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Marton Havasi, Brian Karrer, Itai Gat, and Ricky T. Q. Chen. 2025. [Edit Flows: Flow Matching with Edit Operations](#). *Preprint*, arXiv:2506.09018.
- John J Irwin and Brian K Shoichet. 2005. Zinc- a free database of commercially available compounds for virtual screening. *Journal of chemical information and modeling*, 45(1):177–182.
- Jaeyeon Kim, Lee Cheuk-Kit, Carles Domingo-Enrich, Yilun Du, Sham Kakade, Timothy Ngotiaoco, Sitant Chen, and Michael Albergo. 2025. [Any-Order Flexible Length Masked Diffusion](#). *Preprint*, arXiv:2509.01025.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Pengfei Li, Qichang Zheng, and Ziyi Jiang. 2025. [An empirical study on the accuracy of large language models in api documentation understanding: A cross-programming language analysis](#). *Journal of Computing Innovations and Applications*, 3(2):1–14.
- Lightning-AI. 2023. Litgpt. <https://github.com/Lightning-AI/litgpt>.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2024. [2 olmo 2 furious](#). *Preprint*, arXiv:2501.00656.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). *Preprint*, arXiv:1912.01703.
- Dhruv Patel, Aishwarya Sahoo, Avinash Amballa, Tahira Naseem, Tim GJ Rudner, and Andrew McCallum. 2025. Insertion language models: Sequence generation with arbitrary-position insertions. *arXiv preprint arXiv:2505.05755*.
- William Peebles and Saining Xie. 2023. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205.
- Lars Ruddigkeit, Ruud Van Deursen, Lorenz C Blum, and Jean-Louis Reymond. 2012. Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of chemical information and modeling*, 52(11):2864–2875.

Subham Sekhar Sahoo, Marianne Arriola, Aaron Gokaslan, Edgar Mariano Marroquin, Alexander M. Rush, Yair Schiff, Justin T. Chiu, and Volodymyr Kuleshov. 2024. [Simple and Effective Masked Diffusion Language Models](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Galouédec. 2020. Trl: Transformer reinforcement learning. <https://github.com/huggingface/trl>.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. [Hugging-face’s transformers: State-of-the-art natural language processing](#). *Preprint*, arXiv:1910.03771.

Omry Yadan. 2019. [Hydra - a framework for elegantly configuring complex applications](#). Github.

## A Demonstration on LM1B

In section 5, we demonstrated ILM on a seq2seq task where the model generates a path given a graph description as a prompt. Here, we show how to implement ILM for unconditional language modeling on a real-world text corpus (LM1B), where there is no prompt and the model generates text from scratch. The steps translate to any other type of language model.

### A.1 Collator for Unconditional LM

The first difference from the seq2seq task is the collator. In the case of ILM, instead of `ILMSeq2SeqCollator`, which protects the prompt tokens from being dropped, we need to implement a new collator, which we call `DefaultILMCollator`, which drops tokens uniformly from the entire sequence:

```
from xlm.datamodule import Collator

class DefaultILMCollator(Collator):
    """Used for pre-training."""

    def __call__(self, examples: List) -> ILMBatch:
        batch = ilm_single_segment_collate_target_fn(
            [e["input_ids"] for e in examples],
            ...
            sample_n_drops_fn=_n_drop_uniformly, # uniform
            drop_count
            drop_indices_fn=_drop_uniformly, # uniform
            drop positions
```

```
)
return batch
```

## A.2 Datamodule Configuration

The datamodule configuration for unconditional language modeling differs from the seq2seq task at two places: (1) we need to swap out the collator configuration to use our new `DefaultILMCollator`, and (2) we need to add a new dataset manager for unconditional generation, which manages dataset of blank sequences.

```
# configs/datamodule/lm1b_ilm.yaml
#@package _global_
defaults:
- /datasets@datamodule.dataset_managers.train.lm:
  lm1b_train
- /datasets@datamodule.dataset_managers.val.lm: lm1b_test
- /datasets@datamodule.dataset_managers.val.unconditional
  _prediction:
    text_unconditional_prediction
- /collator@datamodule.dataset_managers.train.lm.collator:
  default_ilm
- /collator@datamodule.dataset_managers.val.lm.collator:
  default_ilm
- /collator@datamodule.dataset_managers.val.unconditional
  _prediction.collator:
  default_ilm
```

For unconditional generation, the model must generate text starting from an empty sequence. This is handled by `ILMEmptyDataset`, which generates empty examples that the model fills entirely:

```
class ILMEmptyDataset(IterableDataset):
    def __init__(self, tokenizer: Tokenizer, num_examples:
        int):
        self.tokenizer = tokenizer
        self.num_examples = num_examples

    def __iter__(self):
        for _ in range(self.num_examples):
            yield self.tokenizer("", add_special_tokens=False)
```

This dataset is referenced in the experiment configuration:

```
# configs/experiment/lm1b_ilm.yaml
datamodule:
  dataset_managers:
    val:
      unconditional_prediction:
        dataset_constructor_str:
          ilm.datamodule_ilm.ILMEmptyDataset
```

## B Hydra Configs

Figure 5 and fig. 6 and show the model type and datamodule configs, respectively.

## C Harness

This section describes the functionality implemented in Harness.

```

# @package _global_
loss:
  _target_: ilm.loss_ilm.ILMLoss
  ... # other arguments

predictor:
  _target_: ilm.predictor_ilm.ILMPredictor
  ... # other arguments

reported_metrics:
  train: # reported during training loop
    lm: # dataloader name
      accumulated_loss: # metric name
        prefix: train/lm # str prefix for logging
        update_fn: ilm.metrics_ilm.mean_metric_update_fn # callable
        ... # metrics for additional dataloaders
    val:
      ...
  test:
    ...

```

Figure 5: The `configs/model_type/ilm.yaml` config file for the ILM model. It contains sections for LossFunction, Predictor and Metrics.

### C.1 Loggers

Logger components can be registered through the `loggers` key. xLM provides pre-configured tensorBoard and WandB logger configurations, with tensorboard being the default. However, all PyTorch Lightning-supported loggers can also be used.

```

defaults:
  - override /loggers:
    - wandb

```

Figure 7: The entries for loggers in the `experiment` config file.

### C.2 Logging Metrics

The library provides various preconfigured metrics for different stages, such as **accumulated loss** (mean loss value), **exact match**, and **token accuracy**. Each of these metric components inherits from the `torchmetrics.Metric` class and is wrapped by default using the `xlm.metrics.MetricWrapper` module, which manages the computation of its value. Another key method to define is the `update_fn` function, which takes raw input batch sequences and loss function outputs, and transforms them into a dictionary of entries used by the Metric class to compute the final value. This allows for customization, enabling custom metric logic depending on the model and task. Different metrics can be configured for various workflow stages as depicted in fig. 8.

```

# @package _global_
defaults:
  - /datasets@datamodule.dataset_managers.train.lm:
    star_easy_train
  - /datasets@datamodule.dataset_managers.val.lm:
    star_easy_val
  - /datasets@datamodule.dataset_managers.val.prediction:
    star_easy_val_pred
  - /datasets@datamodule.dataset_managers.test.lm:
    star_easy_test
  - /datasets@datamodule.dataset_managers.test.prediction:
    star_easy_test_pred
  -
  - /datasets@datamodule.dataset_managers.predict.prediction:
    star_easy_test_pred
  - /collator@datamodule.dataset_managers.train.lm.collator:
    seq2seq_ilm
  - /collator@datamodule.dataset_managers.val.lm.collator:
    seq2seq_ilm
  - /collator@datamodule.dataset_managers.val.prediction.collator:
    seq2seq_pred_ilm
  - /collator@datamodule.dataset_managers.test.lm.collator:
    seq2seq_ilm
  - /collator@datamodule.dataset_managers.test.prediction.collator:
    seq2seq_pred_ilm
  - /collator@datamodule.dataset_managers.predict.prediction.collator:
    seq2seq_pred_ilm
  ...

```

Figure 6: The `configs/datamodule/star_easy_ilm.yaml` config file for the ILM model. It contains sections for datasetmanagers and collators.

### C.3 Logging Predictions

The model predictions for validation and test sets are logged under the `logs/runs` directory. The configuration for this is specified using the `log_predictions` key, and xLM’s `xlm.log_predictions.LogPredictions` component should be used. The logger file contains a **text** field, which contains the prefix and generated sequence, and a **truth** field, which contains the ground-truth sequence. Predictions can be logged to a local file, trainer loggers, or the console by using the values ‘file’, ‘logger’, or ‘console’.

```

log_predictions:
  _target_: xlm.log_predictions.LogPredictions
  fields_to_keep_in_output:
    - text
    - truth
  inject_target: target_ids
  writers:
    - file
    - logger

```

Figure 9: The entries for logging predictions in the `experiment` config file.

```

reported_metrics:
  train:
    lm:
      accumulated_loss:
        _target_: xlm.metrics.MetricWrapper
        name: accumulated_loss
        metric:
          _target_: torchmetrics.MeanMetric
          prefix: train/lm
          update_fn:
            xlm.lm.ilm.metrics_ilm.mean_metric_update_fn
      val:
        ...
      test:
        lm:
          accumulated_loss:
            _target_: xlm.metrics.MetricWrapper
            name: accumulated_loss
            metric:
              _target_: torchmetrics.MeanMetric
              prefix: test/lm
              update_fn:
                xlm.lm.ilm.metrics_ilm.mean_metric_update_fn
          prediction:
            exact_match:
              _target_: xlm.metrics.MetricWrapper
              name: exact_match
              metric:
                _target_: xlm.metrics.ExactMatch
                prefix: test/prediction
                update_fn: xlm.metrics.seq2seq_exact_match_update_fn
            token_accuracy:
              _target_: xlm.metrics.MetricWrapper
              name: token_accuracy
              metric:
                _target_: xlm.metrics.TokenAccuracy
                prefix: test/prediction
                update_fn:
                  xlm.metrics.seq2seq_token_accuracy_update_fn

```

Figure 8: Metric related entries in `configs/model_type/ilm.yaml` config file for the ILM Model.

## D FAQs

### D.1 How to add a new task/dataset?

To add a new task, one must prepare the corresponding datasets in a Hugging Face (HF) dataset compatible format. Depending on the use case, separate dataset configuration (`.yaml`) files can be created for each stage (train/val/test/pred), providing the flexibility to process the same dataset differently or to use entirely different datasets across stages. The HF-downloadable dataset can be specified using the `full_name` key, and the `xlm.datamodule.DatasetManager` from xLM can be reused for dataset manager instantiation. This manager automatically handles downloading, caching, preprocessing, and data-loading operations according to the dataset configuration entries. Alternatively, datasets can be used locally without being uploaded to the HF Hub by employing `xlm.datamodule.LocalDatasetManager`. Once these configuration files are prepared, they must be registered in the `configs/datamodule/MODEL.yaml` (Fig-

ure 3).

### D.1.1 Implementing custom DatasetManager

In addition to xLM’s datasetmanager component, one can implement a custom datasetmanager for greater flexibility by inheriting from it. The necessary methods can be overridden—for example, the `_download` method shown below, where custom logic can be added to read and process the required type of file format.

```

from xlm.datamodule import DatasetManager
import datasets

class CustomDatasetManager(DatasetManager):

    def _download(self) -> datasets.Dataset:
        ...
        return dataset

```

The component can be configured along with mentioning the necessary arguments by adding `configs/datasets/custom_dataset_train.yaml`

```

_target_: CustomDatasetManager
... # other arguments

```

This config can then be registered in the `configs/datamodule/MODEL.yaml` file

```

defaults:
  -/datasets@datamodule.dataset_managers.train.lm:
    custom_dataset_train
    ... # other dataset and collator entries

```

## E Troubleshooting

### E.1 Hydra Errors

**Unable to find a package ... error by Hydra:** See the name of the package in the error message. For example, if you encounter `Unable to find or instantiate abc.xyz.MyClass`, first try to import it manually in the Python interpreter: `python -c "from abc.xyz import MyClass"`.

**Hydra Composition Errors:** First check the Hydra documentation <https://hydra.cc/docs/intro/>. If the error persists, write a single experiment config without using defaults list for components.

## F Additional Features

### F.1 Modules

The library, under the `models` component, provides several architectural implementations that can be used to easily build diverse model backbones for prevalent non-autoregressive workflows in the literature. We provide modules for standard decoder

only transformer, Diffusion Transformers (Peebles and Xie, 2023), rotary embedder (Su et al., 2024), time embedder, adaptive layer normalization layers (Peebles and Xie, 2023), and some standard noise schedulers. They are available under `xlm.modules`.

## F.2 Push to hub

The library provides a `push_to_hub` command that uploads trained models to the Hugging Face Hub by reconstructing the complete training environment (datamodule, tokenizer and model architecture) from Hydra configurations.

```
$ xlm-push-to-hub experiment=[CONFIG_PATH]
+hub_checkpoint_path=[CKPT_PATH] +hub.repo_id=[HUB_PATH]
```

The environment variable `HF_HUB_KEY` must be assigned a valid Hugging Face access token.

## F.3 Callbacks

The library extends the PyTorch Lightning callback infrastructure, enabling modular components to integrate with the training cycle (e.g., per-batch updates, validation hooks) in a decoupled manner. It provides a set of extensible callbacks, such as an Exponential Moving Average callback `EMACallback` for maintaining smoothed evaluation weights, a Checkpoint Monitoring callback `ModelCheckpoint` and a Performance Monitoring callback `SpeedMonitorCallback` for tracking training speed and identifying bottlenecks. The callback config file names (xLM's or custom callbacks) can be mentioned in the following way to override the default callbacks:

```
defaults:
- override /callbacks:
- ema
- speed_monitor
- checkpoint_monitor
```

## F.4 Checkpointing

The library provides a checkpointing system that saves training state (model weights, optimizer state, and training progress) to enable recovery from failures and long-running training jobs. It integrates with PyTorch Lightning's built-in `ModelCheckpoint` to save the best-performing model. It also supports frequent lightweight checkpoints using a `ThinningCheckpoint` callback that retains only milestone intervals to save storage and an `OnExceptionCheckpoint` callback that preserves state during crashes.

## G Preconfigured Tasks and Models

**Star Graphs** The library provides three synthetic datasets that involve generating the path from a designated start node to a target node on star-shaped graphs (Bachmann and Nagarajan, 2024; Patel et al., 2025). The three variants follow the naming convention and construction of Patel et al. (2025). StarEasy contains symmetric graphs with the start node fixed at the junction, while StarMedium & StarHard introduce asymmetric structures with variable arm lengths and start nodes that may lie off the junction.

**Language Modeling** For text generation, we provide training and testing config setup for two datasets - LM1B, a large-scale corpus from the news domain, consisting of short text sequences (typically 2–3 sentences), and OpenWebText, which are widely used to benchmark the performance of small language models.

**Models** We benchmark and release three preconfigured models: ARLM, MDLM and ILM.<sup>6</sup>

## H Planned Features

**Non-text datasets** Non-autoregressive sequence generation is useful for non-text tasks like molecule generation (Irwin and Shoichet, 2005; Ruddigkeit et al., 2012), path planning, etc. We plan on adding support for external non-text datasets in the future

**New Models** We plan to add support for newer models (Havasi et al., 2025; Kim et al., 2025).

**FlexAttention** Pytorch 2.5 introduces FlexAttention (Dong et al., 2024), dynamically compiled attention layer which allows fast attention with arbitrary masks. This can be very useful for non-autoregressive sequence generation as it can allow sequence packing eliminating the need for padding even for non-autoregressive models.

## I Resources

Table 3 lists the resources provided through this work.

<sup>6</sup>We are working on benchmarking newer models, which will be released soon.

Table 3: Resources

| Name       | Type           | Description                                                     | Link                                                                                                                                                                | License     | Source               |
|------------|----------------|-----------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------|----------------------|
| StarEasy   | Dataset        | Synthetic star graph dataset                                    | <a href="https://github.com/facebookresearch/xlm/tree/main/xlm/datasets/star_easy">https://github.com/facebookresearch/xlm/tree/main/xlm/datasets/star_easy</a>     | CC-BY-NC-SA | (Patel et al., 2025) |
| StarMedium | Dataset        | Synthetic star graph dataset                                    | <a href="https://github.com/facebookresearch/xlm/tree/main/xlm/datasets/star_medium">https://github.com/facebookresearch/xlm/tree/main/xlm/datasets/star_medium</a> | CC-BY-NC-SA | (Patel et al., 2025) |
| StarHard   | Dataset        | Synthetic star graph dataset                                    | <a href="https://github.com/facebookresearch/xlm/tree/main/xlm/datasets/star_hard">https://github.com/facebookresearch/xlm/tree/main/xlm/datasets/star_hard</a>     | CC-BY-NC-SA | (Patel et al., 2025) |
| xlm-core   | Python Package | xlm-core package for non-autoregressive language modeling       | <a href="https://pypi.org/project/xlm-core/">https://pypi.org/project/xlm-core/</a>                                                                                 | MIT         |                      |
| xlm-models | Python Package | Companion package for xlm-core containing model implementations | <a href="https://pypi.org/project/xlm-models/">https://pypi.org/project/xlm-models/</a>                                                                             | MIT         |                      |



# AITutor-EvalKit: Exploring the Capabilities of AI Tutors

Numaan Naeem\*, Kaushal Kumar Maurya\*,  
Kseniia Petukhova and Ekaterina Kochmar

Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE

{Numaan.Naeem, kaushal.maurya, kseniia.petukhova, ekaterina.kochmar}@mbzuai.ac.ae

## Abstract

We present AITutor-EvalKit, an application that uses language technology to evaluate the pedagogical quality of AI tutors, provides software for demonstration and evaluation, as well as model inspection and data visualization. This tool is aimed at education stakeholders as well as \*ACL community at large, as it supports learning and can also be used to collect user feedback and annotation.

## 1 Introduction

Personalized one-on-one tutoring has long been recognized as a highly effective educational approach (Bloom, 1984). Yet, its widespread adoption is constrained by the limited availability of qualified tutors (Wang et al., 2024b) and the high costs associated with tutor training (Kelly et al., 2020), among other impediments (Yoon et al., 2007; Boyd et al., 2008). An alternative to human tutoring is provided by AI tutoring systems, especially those relying on recent advances in large language models (LLMs), such as Khanmigo (Khan, 2024) and Tutorly.<sup>1</sup> Despite the remarkable success of LLMs in various tasks (Minaee et al., 2024), their adoption in education is hindered by lack of a clear understanding of what these models are capable of (Tack et al., 2023; Jurenka et al., 2024) and how pedagogically useful they are (Macina et al., 2023b), which results in lack of trust on the part of key educational stakeholders. With the fast development of LLMs and their easy integration into learning tools, questions about the evaluation of AI-driven tutor performance become increasingly relevant (Kosmyna et al., 2025). The goal of our tool is two-fold: (1) via the open-source and open-access code, we provide a practical, customizable and versatile evaluation tool that

\* Equal contribution.

<sup>1</sup><https://tutorly.io>

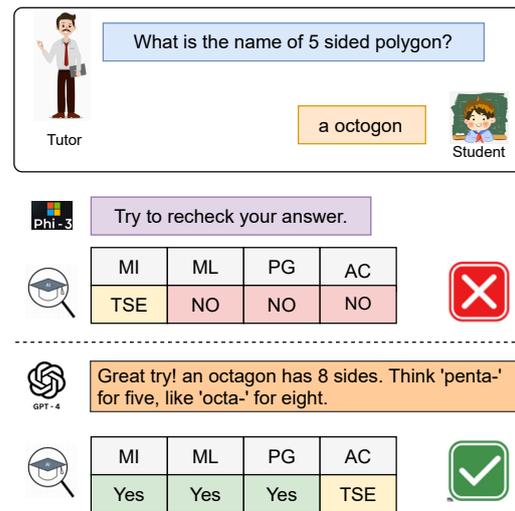


Figure 1: This example shows a sample dialogue and its pedagogical-ability evaluation by the LoMtl model using the AITutor-EvalKit. The evaluation follows the four dimensions from Kochmar et al. (2025): **MI** (Mistake Identification), **ML** (Mistake Location), **PG** (Providing Guidance), and **AC** (Actionability). **TSE**: To some extent.

can be applied to a variety of educational scenarios; and (2) via informative demonstrations, we aim to raise awareness of the stakeholders (e.g., students and teachers) as well as practitioners interested in the state of the art in AI in Education (i.e., \*ACL community at large) of the current capabilities of LLMs-as-tutors. Although our tool is an *early research prototype*, through its extendable nature, we aim to facilitate research in this exciting and emerging area of applied NLP research.

The current version of our tool focuses on the pedagogical quality evaluation of tutor responses in the context of Student Mistake Remediation (SMR) (Boaler, 2013; Handa et al., 2023) in the mathematical domain, where tutors address errors or misconceptions that hinder students' progress (Wang et al., 2024a,b; Macina et al., 2023a). As the foundation for evaluation, we use the established taxonomy from Maurya et al. (2025), which is grounded in the learning

sciences principles and which allows us to assess the quality of tutor responses along four key SMR dimensions (Kochmar et al., 2025): (1) *mistake identification*, concerned with whether tutor’s response notifies the student of the committed mistake; (2) *mistake location*, focusing on whether the tutor clearly points to the erroneous part in the student’s solution; (3) *providing guidance*, evaluating the quality of the pedagogical guidance; and (4) *actionability*, assessing whether tutor’s response makes it clear what the student should do next. These dimensions are further outlined in Table 3, with an illustrative example provided in Figure 1.

The structure and implementation of the front-end and backend of our tool are described in Section 3. Using MRBench dataset (Maurya et al., 2025) and taking inspiration from the BEA 2025 shared task on AI tutor response evaluation (Kochmar et al., 2025), we introduce a novel, efficient, and lightweight multi-task learning model that addresses the four dimensions of pedagogical quality evaluation (§3.1). The outputs of the model, as well as those of an open-source (Prometheus (Kim et al., 2024)) and a commercial (GPT-5) LLMs used as judges, are displayed using an interactive browser-based UI with helpful visualizations (§3.3). The evaluation results (§4) suggest that while our model achieves competitive results when evaluated against gold standard annotation, its outputs are also perceived by users to be at least as accurate as those of the commercial LLM-as-judge models, and the UI is considered informative and easy to use. We present and publicly release:

- The first of its kind, open-access and open-source model aimed at evaluation of the pedagogical quality of AI tutor responses available at MIT-licensed python repository: <https://github.com/kaushal0494/AITutor-EvalKit>. We believe it to be useful for AI-in-Education practitioners and developers, as it is highly customizable, allowing researchers to apply it to further educational contexts and dialogues, as well as to extend it to other scenarios and domains.
- An interactive UI available at <https://demo-ai-tutor.vercel.app>, which communicates the results and showcases the capabilities of AI tutors in an interpretable way, which we consider to be of interest to education stakeholders and the \*ACL community at large. The demo tool can also be run

locally with the user’s own data and models following the instructions provided in the GitHub repository.

- A short video demonstrating the tool available at <https://www.youtube.com/watch?v=9qgDfrhz0vg>.

## 2 Related Work

Over the years, the NLP community has seen significant advances in the development of publicly available toolkits for modeling and evaluation, including Hugging Face Transformers (Wolf, 2019), NLTK (Bird, 2006), and Scikit-learn (Pedregosa et al., 2011). These toolkits have enhanced code reusability, enabling researchers to focus on developing more sophisticated models and metrics.

However, in the educational domain, especially in conversational dialogues, there remains a lack of robust tools to push research boundaries. Some progress has been made with frameworks like ConvoKit (Chang et al., 2020), which facilitates the manipulation and analysis of general conversational data, and the social interactions embedded within. More recently, Edu-ConvoKit (Wang and Demszky, 2024) was developed for preprocessing, annotation, and analysis specifically for educational dialogues. While these toolkits contribute significantly to data handling and analysis, they fall short of addressing the evaluation of pedagogical quality in AI-driven educational systems. To the best of our knowledge, there is no toolkit that supports *on-the-fly* assessment of the pedagogical abilities of AI tutors. With AITutor-EvalKit, we aim to fill this critical gap by providing an open-source toolkit for the systematic evaluation of AI tutors. The toolkit integrates a popular taxonomy proposed by Maurya et al. (2025) into an automated evaluation framework. It is also designed to be extensible to additional evaluation aspects such as the ones proposed by Macina et al. (2025), who augment pedagogical assessment with measures of student expertise and understanding. By supporting and unifying such evaluation methods, AITutor-EvalKit aims to facilitate progress in this underexplored yet important research area.

## 3 System and Demonstration Description

AITutor-EvalKit consists of two major modules: **backend** and **frontend**, with the pipeline illustrated in Figure 2. The backend module includes several models for tutor response evalu-

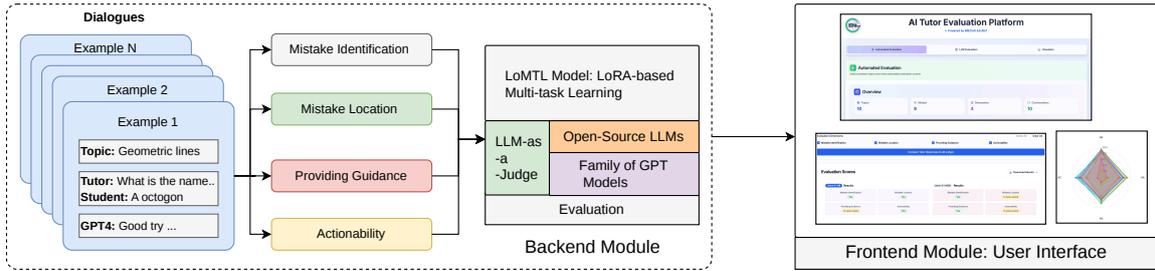


Figure 2: AITutor-EvalKit pipeline: The backend module includes several model options to assess the pedagogical soundness of tutors’ responses, and the frontend presents evaluation outputs in an interactive user interface.

ation, and frontend seamlessly integrates these evaluation outputs with an interactive, customizable, and flexible UI. Different audience groups can benefit from different functionalities of the toolkit: e.g., researchers and developers can use both modules, train their own automated models, use their own datasets, and launch the demo UI locally, while teachers, policymakers and other educational stakeholders can use the frontend module to understand the capabilities of LLMs-as-tutors and make decisions accordingly.

### 3.1 Backend Module: Evaluation Models

The backend module has two components: *a specialized automated evaluation model* and *an evaluation pipeline for LLMs-as-judges*. In this section, we provide details on the task, the training data (for the automated evaluation model), the test data, and the automated evaluation model itself, as well as the usage of an open-source LLM and a closed-source GPT model as judges.

### 3.2 Student Mistake Remediation Task

This task considers mathematical educational dialogues between a student and a tutor, where interactions are driven by student’s mistakes or confusions, and the AI tutor aims to remediate them through pedagogically appropriate responses. Formally, let the conversation history be  $H = \{(T_1, S_1), (T_2, S_2), \dots, (T_t, S_t)\}$ , where  $T_i$  and  $S_i$  denote the tutor’s and student’s  $i$ -th utterances, respectively. Let  $S_k$ ,  $k \in [1, \dots, t]$ , denote the most recent student’s utterance containing a mistake or confusion; the tutor then produces  $T_{t+1}$  to address it. The proposed toolkit evaluates the pedagogical quality of  $T_{t+1}$  along eight dimensions defined by Maurya et al. (2025).

#### 3.2.1 Dataset

As discussed in Section 1, we used MRBench dataset by Maurya et al. (2025), which has 491 dialogues (300 in the development set and 191 in the test set), each paired with seven LLM-generated tutor responses and one or two human tutor responses. Each response is annotated by human annotators for each of the four dimensions using the categories “Yes,” “To some extent,” and “No.” We also created a randomly selected *demonstration set* consisting of 10 dialogues, which is a subset of the test set. More details on the dimensions, annotation, and data statistics can be found in Appendix A.

#### 3.2.2 Specialized Automated Model: LoMTL

This component implements a ternary classifier to evaluate tutor responses across the four considered dimensions. The models developed by participants in the BEA shared task provide a good starting point, as summarized by Kochmar et al. (2025). However, upon closer inspection, we found that several teams relied on closed-source models, fully fine-tuned and LoRA-tuned models, and some even used multiple models or ensembles for each dimension, which makes these approaches difficult to scale across all dimensions.

Considering this, we propose a novel LoRA-based multi-task model (called LoMTL) that fine-tunes a single small LLM in a LoRA setting across all four dimensions, each treated as a task in a multi-task learning setting. This modeling is naturally suited to the four tasks, all of which are ternary, closely related classification problems, allowing them to benefit from shared learning during training. Additionally, we use sampling and balanced batching methods to improve performance. We observe that a small google/gemma-2-2b-it (Team et al., 2024) model using LoMTL achieves performance compet-

itive with the top-performing teams in the shared task while remaining highly efficient and scalable. Appendix B presents more details on model development, prompts, configuration, and comparative results with the BEA shared task’s top-performing teams.

### 3.2.3 LLM-as-a-Judge Evaluation

This component provides functionality for evaluation using LLMs-as-judges. Following Maurya et al. (2025), we have selected Prometheus2 (Kim et al., 2024) as the primary open-source LLM. However, the implementation of this module is flexible enough to support the family of Llama and other open-source causal LLMs for evaluation. Additionally, we have included an option to evaluate using the closed-source OpenAI GPT-5 model, although the implementation supports any model from the GPT family. Users must provide their OpenAI API key to run evaluations. We selected the above two models considering their open- and closed-source nature, as well as their human-like performance on public benchmarks (Kim et al., 2024). However, the codebase is flexible enough to require only minimal adaptation to support other LLMs.

### 3.2.4 Flexible Design

Each of the three evaluation setups (automated evaluation, evaluation with open-source LLMs, and evaluation with closed-source GPT models) provides high degree of flexibility in the choice of the base LLM, prompting strategy, and hyperparameter tuning. For example, each evaluation setup has its own prompting file and configuration, allowing users to customize and use the components of the tool as per their needs. Evaluation can be run using the following short commands:

```
# Training LoMTL model
cd src/autoeval && ./lora_finetune_runner.sh

# Evaluation with automated (i.e., LoMTL) model
cd src/autoeval && ./lora_evaluation_runner.sh

# Evaluation with open-source LLM
cd src/llmeval &&
python run_open_llm_as_judge_evaluation.py

# Evaluation with GPT5
cd src/llmeval && ./gpt5_eval_runner.sh
```

## 3.3 Frontend Module: Demo App

We present an interactive prototype built on top of MRBench, designed to evaluate the pedagogi-

cal abilities of AI tutors in educational dialogues. Users can explore the demo in two modes: (1) a **Static Mode**, where educators can access a deployed version containing 10 conversations from the MRBench test set evaluated by our fine-tuned model and share feedback on the helpfulness of tutor responses; and (2) a **Customized Mode**, where developers can run the interface locally to analyze their own datasets using either our LoMTL model or their own one, supported by the provided code (the exact step-by-step details are supplied in the official GitHub README file<sup>2</sup>). The UI consists of three key modules **Automated Evaluation**, **LLM Evaluation**, and **Visualizer**, that together enable exploration, analysis, and visualization of the key aspects of educational dialogues.

- **Automated Evaluation** provides automated evaluation results using our LoMTL model on 4 evaluation dimensions.
- **LLM Evaluation** leverages LLMs-as-judges for human-aligned evaluations.
- **Visualizer** enables rich visual analytics for interpreting evaluation scores across 4 pedagogical dimensions on the MRBench development set.

The UI supports the selection of evaluation methods and provides instant visualization of performance metrics through plots, bar charts, and spider graphs.

### 3.3.1 Automated Evaluation UI

The Automated Evaluation module presents results generated by our LoMTL model, assessing AI tutors across four pedagogical dimensions. As shown in Figure 6, users are first presented with an overview panel summarizing key statistics such as the number of topics, models, dimensions, and conversations. Users can evaluate a single tutor or compare two tutors side by side.

**Single Tutor Evaluation:** In this mode, users select a problem topic to view the complete student-tutor dialogue in the “Context” block (Figure 7). A drop-down menu allows selection of a tutor model, and the corresponding response is displayed in the “Tutor Response” block. Users can rate the usefulness of the response (*Helpful*, *To Some Extent*, or *Not Helpful*) and optionally view the ground truth solution for better context. Upon clicking “Get Auto-Evaluation Results,” the

<sup>2</sup><https://github.com/kaushal0494/AITutor-EvalKit/blob/main/README.md>

system generates performance evaluations across the chosen dimensions, with results downloadable in a PNG, JPG, or JSON format. The “Best Performance by Dimension” panel highlights the top tutor(s) for each dimension, helping users quickly identify their pedagogical strengths.

**Tutor Comparison Mode:** This mode allows users to directly compare two tutors by enabling the “Tutor Comparison Mode” option. The selected tutors’ responses are displayed side by side (Figure 8), and users can provide quick feedback indicating which tutor performed better or mark both as equally good or bad. After choosing the evaluation dimensions and clicking “Compare Tutor Responses,” the interface presents a detailed two-column comparison of scores across all selected dimensions. To facilitate interpretation, the “Comparison Visualization” panel (Figure 9) provides four interactive views: the **Summary** view highlights the leading tutor for each evaluation dimension and identifies the overall winner; the **Spider Chart** offers a radar-style visualization comparing performance patterns across dimensions; the **Bar Chart** displays side-by-side scores for each dimension; and the **Differences** view illustrates the magnitude of score gaps between tutors. All visualizations can be exported in PNG or JPG formats for reporting or analysis. Finally, the “Best Performance by Dimension” panel summarizes the comparative strengths of the tutor pair, providing a concise overview of pedagogical differences. This mode supports structured, interpretable, and visually grounded benchmarking of AI tutor performance.

### 3.3.2 LLM Evaluation UI

The LLM Evaluation module enables advanced pedagogical assessment of AI tutor responses using LLMs as judges. It extends the functionality of the automated evaluation pipeline by leveraging LLMs to judge and compare tutor responses across four pedagogical dimensions. The UI supports three evaluation modes: *single tutor evaluation*, *tutor model comparison*, and *LLM judge comparison*. As shown in Figure 10, the overview panel summarizes available topics, conversations, tutor models, evaluation dimensions, and judge LLMs, providing a quick snapshot of the evaluation setup. Currently, two LLM judges are supported: GPT-5 and Prometheus-7B-v2.0 (Kim et al., 2024).

**Single Tutor Evaluation:** This mode allows users to analyze how a selected tutor performs on

a specific problem using an LLM as a judge. After selecting the problem, tutor, and LLM judge, users can generate dimension-wise evaluation results that reflect the LLM’s assessment of the tutor’s pedagogical performance. A “Best Performance by Dimension” panel highlights the top-performing tutor(s) for each dimension.

**Tutor Comparison Mode:** This mode enables side-by-side comparison of two tutor models on the same problem, judged by a selected LLM. Tutor responses are displayed together (Figure 11), and upon comparison, the system presents dimension-wise evaluations and visualizations such as Summary, Spider Chart, Bar Chart, and Differences to clearly show relative strengths across pedagogical aspects.

**Judge Comparison Mode:** This mode compares how different LLM judges evaluate the same tutor response. The tutor’s response is displayed once, and evaluations from both judges are presented in parallel with corresponding visualizations. This feature helps assess consistency between LLM judges and identify possible biases in their evaluations.

Together, these modes enable fine-grained, interpretable analysis of tutor behavior, offering insights into both model performance and evaluation reliability across different judging LLMs.

### 3.3.3 Visualizer UI

The Visualizer module provides a high-level overview of the MRBench development set, using gold-standard annotations across four evaluation dimensions through intuitive visual analytics. Users are first presented with a “Dataset Overview” panel summarizing key statistics, including the number of conversations, tutor models, and evaluation dimensions. This module includes three main visualization panels: *Tutor Performance Summary*, *Visualization Controls*, and *Dataset Visualization*.

The **Tutor Performance Summary** panel presents average scores for each tutor model across all four dimensions, where “Yes,” “To some extent,” and “No” correspond to 1.0, 0.5, and 0.0, respectively (Figure 12). The **Visualization Controls** panel allows users to select specific tutors and dimensions to generate detailed visualizations (Figure 13). The **Dataset Visualization** panel then displays the results through spider and bar charts, where spider plots (the *most informative* visualizations from our perspective) summarize tutors’

strengths and weaknesses across dimensions (Figure 14), while bar charts show detailed score distributions for selected dimensions (Figure 15), along with averages for each response label.

This visual analytics module enables users to explore and interpret pedagogical quality effectively, supporting comparative analysis and data-driven insights. The same functionality is also available in the customized mode, allowing developers to visualize and analyze their own datasets locally in a similar way.

## 4 Evaluation

To assess the toolkit and evaluation models, we first measure models’ performance using the metrics from the BEA 2025 shared task (Kochmar et al., 2025) – accuracy and macro-F1, and then conduct a human evaluation study, in which participants assess both the prediction quality of the LoMTL evaluation model and the usability of our demo tool and UI.

### 4.1 Intrinsic Evaluation: Quantitative Analysis

For intrinsic evaluation, we compare our LoMTL model’s predictions on the test set from Kochmar et al. (2025) with the gold human annotations and also run GPT-5 and Prometheus2 on the same data. The average results across evaluation dimensions, presented in Table 1, show that Prometheus2 performs substantially worse than both our model and GPT-5. Our model achieves the highest averaged accuracy and macro-F1 on the full test set, and it also performs competitively on the demonstration subset. While GPT-5 achieves a slightly higher averaged macro-F1 on the demonstration subset, the overall trend indicates that our model provides more reliable evaluations than Prometheus2 and performs on par with, or better than, GPT-5. A close inspection of the confusion matrix between human annotations and the LoMTL model shows strong overall agreement, with the majority of instances concentrated on the diagonal, particularly for clear *Yes* (3,106) and *No* (1,057) cases. However, the model exhibits systematic confusion on borderline instances and a mild tendency to over-predict *Yes*, especially when humans assign *To some extent* or *No* labels. Extended results and a related discussion, including average precision and recall scores, are provided in Table 7 in Appendix C.

| Model        | Full Test Set |             | Demonstration Set |             |
|--------------|---------------|-------------|-------------------|-------------|
|              | Accuracy      | Macro-F1    | Accuracy          | Macro-F1    |
| LoMTL (ours) | <b>0.72</b>   | <b>0.60</b> | <b>0.68</b>       | 0.55        |
| Prometheus2  | 0.47          | 0.41        | 0.41              | 0.34        |
| GPT-5        | 0.66          | 0.58        | 0.66              | <b>0.59</b> |

Table 1: Accuracy and macro-F1 scores (averaged across Mistake Identification (MI), Mistake Location (ML), Providing Guidance (PG), and Actionability (AC)) for LoMTL, Prometheus2, and GPT-5 on the full test set from Kochmar et al. (2025) and on the demonstration set. Best results are shown in **bold**.

**Performance across dimensions:** The per-dimension results, presented in Table 2, show that Prometheus2 performs poorly, while our model notably outperforms GPT-5 on the Mistake Identification and Actionability dimensions. However, it underperforms GPT-5 by 3 and 6 percentage points in terms of macro-F1 for Mistake Location and Providing Guidance, respectively. A similar trend is observed in the ten dialogues used for demonstration.

### 4.2 Extrinsic Evaluation: Human Study

Participants were given access to the demo website, detailed guidelines, and an evaluation form. The form instructed them to assess three components: the Automated Evaluation tab, the LLM Evaluation tab, and the Visualizer tab, as detailed in Section 3.3. For the first two tabs, participants were asked to explore at least five different dialogues, review at least two tutor responses per dialogue, and use the feedback field to rate each response as *helpful*, *helpful to some extent*, or *not helpful*. They also examined the model’s evaluation for each response across at least one assessment dimension. Additionally, for each dialogue, participants compared at least one pair of tutor responses in comparison mode and indicated which response they considered a better one. After using Automated Evaluation tab, they reported how frequently they agreed with the model’s judgments on a 1-5 scale, both in single-response and comparison modes. In the LLM Evaluation tab, participants evaluated how often they agreed with GPT-5 and Prometheus2, again in both modes. They were also asked which evaluation model they perceived as more accurate: GPT-5 or the model from the first tab (which corresponds to our LoMTL model), and Prometheus2 or the model from the first tab. Finally, in the Visualizer tab, participants explored visualizations for at least two tutors and

| Model        | Full Test Set |             |             |             |             |             |             |             | Demonstration Set |             |             |             |             |             |             |             |
|--------------|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|              | Accuracy      |             |             |             | Macro-F1    |             |             |             | Accuracy          |             |             |             | Macro-F1    |             |             |             |
|              | MI            | ML          | PG          | AC          | MI          | ML          | PG          | AC          | MI                | ML          | PG          | AC          | MI          | ML          | PG          | AC          |
| LoMTL (ours) | <b>0.86</b>   | 0.67        | 0.63        | <b>0.70</b> | <b>0.66</b> | 0.55        | 0.54        | <b>0.65</b> | <b>0.76</b>       | <b>0.69</b> | 0.60        | <b>0.68</b> | <b>0.57</b> | 0.50        | 0.47        | <b>0.65</b> |
| Prometheus   | 0.58          | 0.53        | 0.31        | 0.46        | 0.48        | 0.42        | 0.32        | 0.43        | 0.48              | 0.41        | 0.42        | 0.32        | 0.34        | 0.30        | 0.38        | 0.30        |
| GPT-5        | 0.67          | <b>0.68</b> | <b>0.70</b> | 0.58        | 0.53        | <b>0.58</b> | <b>0.61</b> | 0.55        | 0.67              | 0.64        | <b>0.67</b> | 0.66        | 0.53        | <b>0.56</b> | <b>0.59</b> | 0.64        |

Table 2: Accuracy and macro-F1 scores of our model, Prometheus, and GPT-5 across Mistake Identification (MI), Mistake Location (ML), Providing Guidance (PG), and Actionability (AC) on the full test set from Kochmar et al. (2025) and on the demonstration set. Best results are shown in **bold**.

rated how informative they found them on a 1-5 scale. For every tab, participants also rated the ease of use on a 1-5 scale. The full questionnaire is provided in Appendix C.

A total of 14 participants took part in the study. Their educational background included individuals pursuing a Master’s degree, those holding a Master’s degree, and those holding a PhD. Eleven participants had teaching experience, and eight had prior experience using an AI tutor.

Most participants perceived the LoMTL evaluation model as more accurate in its judgments than both GPT-5 and Prometheus2. As shown in Figure 3, the majority of participants agreed with the LoMTL model’s assessments more than half of the time in both single-response and comparison modes. A similar trend was observed for GPT-5: most participants agreed with its judgments more than half of the time, although a larger proportion reported agreement only about half of the time. In contrast, participants tended to agree with Prometheus2 less than half of the time in the single-response mode, but more than half of the time in the comparison mode. Overall, most participants rated the ease of use of all tabs as *very easy* and found the visualizations *very informative*.

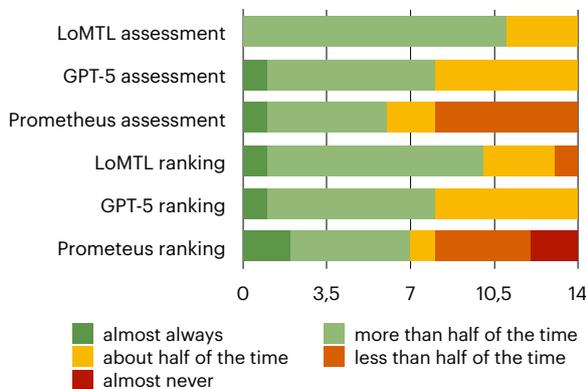


Figure 3: Participants’ responses indicating how often they agreed with the models’ judgments in single-response and comparison modes.

In total, we collected 95 annotations for single-tutor responses and 115 for pairwise comparisons. The analysis of these annotations is provided in Appendix C.2. We do not present it here because evaluating tutor performance is not the focus of this work. Instead, our goal is to demonstrate how our interactive UI can be used for annotation purposes by collecting data that can later be used to train evaluation models or align language models.

## 5 Conclusions and Future Work

This paper introduces the first open-access, open-source model for pedagogical quality evaluation of AI tutor responses, released under an MIT license. It is grounded on pedagogical principles and presents multiple evaluation choices including our proposed light-weight multi-tasking LoMTL model. The toolkit is highly customizable, allowing researchers and practitioners to extend it to diverse educational contexts and dialogues across domains, and it can be easily set up and run locally. In addition, we provide an interactive web-based UI that offers interpretable evaluations of the pedagogical capabilities of state-of-the-art LLMs acting as AI tutors for education stakeholders, policy-makers, and non-technical audience.

Future work will extend the toolkit by (1) integrating new user dialogues and LLMs in the front-end module to generate real-time responses and evaluations; (2) expanding the evaluation dimensions; (3) enabling new data upload option in UI; and (4) broadening coverage to additional subjects, grade levels, and languages.

## Limitations

We acknowledge that our work has several limitations. Below, we summarize the major ones among them.

**Domain and grade level:** In this work, we focus on the mathematical tasks at the middle-school

level. This decision is motivated by the availability of the data at this level and in this domain, but we plan to extend our work to other domains and levels as we elaborate in Section 5. Moreover, since users can run our tool on their own data, new domains and levels can already be integrated on the user’s side.

**Language:** Similarly, our current work focuses on English only. In the future, we hope to extend it to other languages, as we specify in Section 5, while users of our tool can also apply it to data in other languages with their own evaluation models and LLMs-as-judges on their side.

**Educational scenario:** Building on previous work in this domain (Maurya et al., 2025; Kochmar et al., 2025), we only address student mistake remediation as an educational scenario. While this is one of the most salient and challenging scenarios in educational dialogues, we recognize this as one of the limitations of our work and plan to address it in the future.

**Context length:** Similarly, following up on the previous work, our current evaluation approach is limited to a single turn in the dialogue. We acknowledge this as a limitation, and believe that future work should extend single-turn approaches to multiple-turn or full-dialogue ones.

**Taxonomy:** We have built our prototype tool around a well-established evaluation taxonomy of Maurya et al. (2025). While using a single taxonomy is a limitation, our code is open-source and can be extended to incorporate other data, dialogues, taxonomies, and evaluation models on the user’s side.

**Models:** Finally, via our demo tool, we only showcase a few LLMs-as-tutors and deploy only two LLMs-as-judges. Since our code is open-source and extendable, more tutor models and LLMs-as-judges can be integrated on the user’s side.

## Ethical Considerations

As this work is exploratory, we do not anticipate any significant ethical risks associated with it. Moreover, one of our main goals in this research is to raise awareness of the education stakeholders as well as practitioners interested in the state of the art in AI in Education and the current capabilities of LLMs-as-tutors. Through transparent eval-

uation of such models, we aim to improve their interpretability, which we hope will help avoid potential future risks associated with a wider adoption of these models in education.

This work uses the MRBench dataset (Maurya et al., 2025), which integrates the MathDial (Macina et al., 2023a) and Bridge (Wang et al., 2024a) datasets. In these datasets, the identities of the tutors are not revealed, while the student profiles are either synthetically created or anonymized. As a result, we do not anticipate any direct ethical risks associated with the datasets used.

## Acknowledgments

We are grateful to the Google Academic Research Award (GARA) 2024 for supporting this research.

## References

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, and 1 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Jiyuan An, Xiang Fu, Bo Liu, Xuquan Zong, Cunliang Kong, Shuliang Liu, Shuo Wang, Zhenghao Liu, Liner Yang, Hanghang Fan, and Erhong Yang. 2025. BLCU-ICALL at BEA 2025 shared task: Multi-strategy evaluation of AI tutors. In *Proceedings of the 20th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2025)*, pages 1084–1097, Vienna, Austria. Association for Computational Linguistics.
- Anthropic. 2024. The Claude 3 Model Family: Opus, Sonnet, Haiku. In <https://api.semanticscholar.org/CorpusID:268232499>.
- Steven Bird. 2006. NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions*, pages 69–72.
- Benjamin Samuel Bloom. 1984. *The 2 Sigma Problem: The Search for Methods of Group Instruction as Effective as One-to-One Tutoring*. *Educational Researcher*, 13:16 – 4.
- Jo Boaler. 2013. Ability and mathematics: The mindset revolution that is reshaping education. Forum.

- Donald J. Boyd, Pam L. Grossman, Hamilton Lankford, Susanna Loeb, and James Humphrey Wyckoff. 2008. [Teacher Preparation and Student Achievement](#). *Educational Evaluation and Policy Analysis*, 31:416–440.
- Jonathan P Chang, Caleb Chiam, Liye Fu, Andrew Z Wang, Justine Zhang, and Cristian Danescu-Niculescu-Mizil. 2020. Convokit: A toolkit for the analysis of conversations. *arXiv preprint arXiv:2005.04246*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The LLaMA 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yuming Fan, Chuangchuang Tan, and Wenyu Song. 2025. [BJTU at BEA 2025 shared task: Task-aware prompt tuning and data augmentation for evaluating AI math tutors](#). In *Proceedings of the 20th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2025)*, pages 1073–1077, Vienna, Austria. Association for Computational Linguistics.
- Kunal Handa, Margaret Clapper, Jessica Boyle, Rose Wang, Diyi Yang, David Yeager, and Dorottya Demszky. 2023. [“Mistakes Help Us Grow”: Facilitating and Evaluating Growth Mindset Supportive Language in Classrooms](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8877–8897, Singapore. Association for Computational Linguistics.
- Baraa Hikal, Mohamed Basem, Islam Oshallah, and Ali Hamdi. 2025. [MSA at BEA 2025 shared task: Disagreement-aware instruction tuning for multi-dimensional evaluation of LLMs as math tutors](#). In *Proceedings of the 20th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2025)*, pages 1194–1202, Vienna, Austria. 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, and 1 others. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825*.
- Irina Jurenka, Markus Kunesch, Kevin R McKee, Daniel Gillick, Shaojian Zhu, Sara Wiltberger, Shubham Milind Phal, Katherine Hermann, Daniel Kasenberg, Avishkar Bhoopchand, and 1 others. 2024. Towards responsible development of generative AI for education: An evaluation-driven approach. *arXiv preprint arXiv:2407.12687*.
- Sean Kelly, Robert Bringe, Esteban Aucejo, and Jane Cooley Fruehwirth. 2020. Using global observation protocols to inform research on teaching effectiveness and school improvement: Strengths and emerging limitations. *Education Policy Analysis Archives*, 28:62–62.
- Sal Khan. 2024. Khanmigo.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. [Prometheus 2: An Open Source Language Model Specialized in Evaluating Other Language Models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4334–4353, Miami, Florida, USA. Association for Computational Linguistics.
- Ekaterina Kochmar, Kaushal Maurya, Kseniia Petukhova, KV Aditya Srivatsa, Anaïs Tack, and Justin Vasselli. 2025. [Findings of the BEA 2025 shared task on pedagogical ability assessment of AI-powered tutors](#). In *Proceedings of the 20th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2025)*, pages 1011–1033, Vienna, Austria. Association for Computational Linguistics.
- Nataliya Kosmyrna, Eugene Hauptmann, Ye Tong Yuan, Jessica Situ, Xian-Hao Liao, Ashly Vivian Beresnitzky, Iris Braunstein, and Pattie Maes. 2025. Your brain on chatgpt: Accumulation of cognitive debt when using an ai assistant for essay writing task. *arXiv preprint arXiv:2506.08872*.
- 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 the Association for Computational Linguistics: EMNLP 2023*, pages 5602–5621, Singapore. Association for Computational Linguistics.
- Jakub Macina, Nico Daheim, Ido Hakimi, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2025. [MathTutorBench: A benchmark for measuring open-ended pedagogical capabilities of LLM tutors](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 204–221, Suzhou, China. Association for Computational Linguistics.
- 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, Dubrovnik, Croatia. Association for Computational Linguistics.
- Yuren Mao, Yuhang Ge, Yijiang Fan, Wenyi Xu, Yu Mi, Zhonghao Hu, and Yunjun Gao. 2025. A survey on lora of large language models. *Frontiers of Computer Science*, 19(7):197605.
- Kaushal Kumar Maurya, Kv Aditya Srivatsa, Kseniia Petukhova, and Ekaterina Kochmar. 2025. [Unifying AI tutor evaluation: An evaluation taxonomy for pedagogical ability assessment of LLM-powered AI](#)

- tutors. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1234–1251, Albuquerque, New Mexico. Association for Computational Linguistics.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196*.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Mathieu Brucher, Matthieu Perrot, and Édouard Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Jihyeon Roh and Jinhyun Bang. 2025. [bea-jh at BEA 2025 shared task: Evaluating AI-powered tutors through pedagogically-informed reasoning](#). In *Proceedings of the 20th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2025)*, pages 1049–1059, Vienna, Austria. Association for Computational Linguistics.
- Anaïs Tack, Ekaterina Kochmar, Zheng Yuan, Serge Bibauw, and Chris Piech. 2023. [The BEA 2023 shared task on generating AI teacher responses in educational dialogues](#). In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 785–795, Toronto, Canada. Association for Computational Linguistics.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Rose Wang and Dorottya Demszky. 2024. [EduConvoKit: An Open-Source Library for Education Conversation Data](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: System Demonstrations)*, pages 61–69, Mexico City, Mexico. Association for Computational Linguistics.
- Rose Wang, Qingyang Zhang, Carly Robinson, Susanna Loeb, and Dorottya Demszky. 2024a. [Bridging the Novice-Expert Gap via Models of Decision-Making: A Case Study on Remediating Math Mistakes](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2174–2199, Mexico City, Mexico. Association for Computational Linguistics.
- Rose E Wang, Ana T Ribeiro, Carly D Robinson, Susanna Loeb, and Dorottya Demszky. 2024b. Tutor CoPilot: A Human-AI Approach for Scaling Real-Time Expertise. EdWorkingPaper No. 24-1054. *Annenberg Institute for School Reform at Brown University*.
- T Wolf. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Kwang Suk Yoon, Teresa Duncan, Silvia Wen-Yu Lee, Beth Scarloss, and Kathy L Shapley. 2007. Reviewing the evidence on how teacher professional development affects student achievement. *issues & answers. rel 2007-no. 033*. ERIC.

## A Pedagogical Dimensions and Data

### A.1 Evaluation Dimensions

Maurya et al. (2025) proposed an evaluation taxonomy with eight dimensions to assess the pedagogical soundness of  $T_{t+1}$  tutor response in the context of SMR. These dimensions are grounded in learning science research and prior work on tutor evaluation, defining the assessment via eight concrete criteria. Furthermore, the authors validated these dimensions as necessary and sufficient through a human pilot study. They also released the associated MRBench dataset with human annotations (see §A.2). Following this, Kochmar et al. (2025) focused on four key dimensions from this taxonomy for the BEA shared task, where participating teams were challenged to develop a ternary classification model for each of the four dimensions. These classifiers aimed to assess the pedagogical quality of the  $T_{t+1}$  response *on-the-fly*, enabling evaluation scalability with new data and tutors. Details on each of the four dimensions, along with their definitions, annotated labels, and desiderata, are provided in Table 3. Building on this, our AITutor-EvalKit toolkit focuses on the four key dimensions used in the BEA shared task.

### A.2 Dataset Details and Statistics

We used the MRBench dataset to develop the toolkit and the automated evaluation model (i.e., LoMTL). The initial version of MRBench, released by Maurya et al. (2025), contains 191 dialogues. This was extended by Kochmar et al. (2025) which results in 300 dialogues in the development set and 191 dialogues in the test set. In this work, we use the extended version of the MRBench dataset.

The dataset is built on top of two public datasets – MathDial (Macina et al., 2023a) and Bridge (Wang et al., 2024a) – which provide partial conversational histories from secondary and primary school-level mathematics, respectively, along with human tutor responses as  $T_{t+1}$ . MathDial includes only one expert tutor response, whereas Bridge includes two responses from expert and novice human tutors. Additionally, each dialogue includes seven  $T_{t+1}$  responses generated by seven *state-of-the-art* LLMs-as-tutors, including GPT-4 (Achiam et al., 2023), Gemini (Reid et al., 2024), Sonnet (Anthropic, 2024), Mistral (Jiang et al., 2023), Llama-3.1-8B and Llama-3.1-405B (Dubey et al., 2024), and Phi3 (Abdin et al., 2024).

All  $T_{t+1}$  responses (whether from human tu-

tors or LLMs) are annotated by human annotators across four selected dimensions with labels "Yes," "To some extent," and "No," as detailed in Table 3. The proposed LoMTL model is trained on the development set and evaluated on the test set; all results presented in this work are reported on the test set. We further split the development set into a 9:1 ratio for training and validation when developing the LoMTL model. All model checkpoints were selected based on validation performance. Finally, we randomly selected a subset of 10 dialogues from the test set as a demonstration set, which is used in the demo app. These details are summarized in Table 4.

## B LoMTL Evaluation Model

### B.1 Motivation

Building LoMTL has a two-fold motivation: (1) The current human annotation-based pedagogical ability assessments presented by Maurya et al. (2025) are static in nature. They are not scalable to new LLMs and tutoring systems, which are being developed very frequently nowadays. We need a reliable *automated* evaluation model that can provide pedagogical assessment *on-the-fly* for new tutors or responses and help track the progress of AI tutor abilities. (2) The BEA shared task (Kochmar et al., 2025) is a good starting point for developing an automated evaluation model. More than 50 international teams participated in the challenge and proposed several novel modeling approaches, including diverse prompting strategies, full instruction tuning, LoRA-based finetuning, supervised finetuning, data augmentation, label balancing, ensembling and so on. However, most teams that participated in all four tracks did not develop a unified approach (with the exception of the MSA team (Hikal et al., 2025)) and instead used models with a large number of parameters. This hinders the adaptability of these approaches, as their deployment becomes challenging and costly.

These limitations motivated us to develop LoMTL, a lightweight model with only 2 billion parameters, created by training the google/gemma-2-2b-it model using LoRA in a multi-task learning setting. It achieves comparable performance while being significantly more efficient (see Table 5 for comparison). For instance, the top-performing BJTU team (Fan et al., 2025) achieved a macro-F1 score of 0.645 using 288 billion parameters (across all

| Dimension              | Definition                                                                                                          | Labels                                  | Desiderata |
|------------------------|---------------------------------------------------------------------------------------------------------------------|-----------------------------------------|------------|
| Mistake identification | Has the tutor identified a mistake in a student's response?                                                         | (1) Yes<br>(2) To some extent<br>(3) No | Yes        |
| Mistake location       | Does the tutor's response accurately point to a genuine mistake and its location?                                   | (1) Yes<br>(2) To some extent<br>(3) No | Yes        |
| Providing guidance     | Does the tutor offer correct and relevant guidance, such as an explanation, elaboration, hint, examples, and so on? | (1) Yes<br>(2) To some extent<br>(3) No | Yes        |
| Actionability          | Is it clear from the tutor's feedback what the student should do next?                                              | (1) Yes<br>(2) To some extent<br>(3) No | Yes        |

Table 3: An overview of the evaluation taxonomy, associated definitions, annotation labels, and desired labels from Maurya et al. (2025).

| Parameters                                     | Value/Details                                                       |
|------------------------------------------------|---------------------------------------------------------------------|
| Number of dialogues (Dev / Test / Total)       | 300 / 191 / 491                                                     |
| Number of tutor responses (Dev / Test / Total) | 2,476 / 1,547 / 4,023                                               |
| Number of tutors (Total)                       | 9                                                                   |
| Number of human tutors                         | 2 (1 expert, 1 novice)                                              |
| Number of LLM tutors                           | 7                                                                   |
| LLM tutor models                               | GPT-4, Gemini, Sonnet, Mistral, Llama-3.1-8B, Llama-3.1-405B, Phi-3 |
| Source datasets                                | MathDial, Bridge                                                    |
| MathDial                                       | Expert human tutor only                                             |
| Bridge                                         | Expert and Novice human tutors                                      |
| Demonstration set size                         | 10 dialogues (from test set)                                        |

Table 4: Key details of the extended MRBench dataset.

four dimensions). In contrast, the LoMTL model achieved 0.60 with only 2 billion parameters, approximately 0.7% of the BJTU model's parameter count.

## B.2 Training and Inference

In this section, we provide details on the training and evaluation of the LoMTL model. Inspired by the success of LoRA-based modeling from the BEA shared task (Kochmar et al., 2025) and by the community (Mao et al., 2025), we adapted a LoRA-based fine-tuning approach. Since we have a small amount of training data (approximately 2500 examples for each dimension) and the different tasks are somewhat related, the natural modeling choice is multi-task learning, where each evaluation dimension is formulated as a task. This LoRA-based multi-task fine-tuning approach (called LoMTL) resulted in a compact single model and enabled flexibility in model deployment. Further, we experimented with a small google/gemma-2-2b-it model, which resulted in fast inference. We observed two major issues during training: *task imbalance* and *label imbalance*. To mitigate task imbalance, we implemented a balance batching where each batch has an uniform number of examples from each task. For label imbalance, we

explored several approaches such as focal loss, label sampling, loss weighting, and sampling methods. We obtained the best performance with over-sampling where we randomly sample underrepresented examples in the training dataset. Model training and inference were done with a single 48GB A6000 GPU. The best checkpoints were obtained using validation data (10% of the development dataset).

## B.3 Prompts and Configurations

Prompt structure and training/evaluation configurations for the LoMTL model are shown in Figure 4 and Table 6, respectively.

## C Toolkit Evaluation Details and Results

### C.1 Extended Evaluation Results

In addition to the observations on the accuracy and macro-F1 scores in Section 4.1, a closer inspection of precision and recall (from Table 7) further supports our findings. On the full test set, LoMTL achieves the highest macro-precision (0.63) while maintaining competitive recall (0.59), indicating more accurate and consistent positive predictions compared to both baselines. GPT-5 attains slightly higher recall (0.60) but with lower

| Team/Model                     | Macro-F1 | # LLMs                            | # Parameters   | Parameter Size vs. 2B |
|--------------------------------|----------|-----------------------------------|----------------|-----------------------|
| BJTU (Fan et al., 2025)        | 0.646    | 4 (Qwen/Qwen2.5-72B)              | 4 × 72B = 288B | 144× larger           |
| MSA (Hikal et al., 2025)       | 0.643    | 5 × 4 (mistralai/Mistral-7B-v0.1) | 20 × 7B = 140B | 70× larger            |
| BLCU-ICALL (An et al., 2025)   | 0.632    | 4 (Qwen/Qwen2.5-7B)               | 4 × 7B = 28B   | 14× larger            |
| bea-jh (Roh and Bang, 2025)    | 0.625    | 4 (zai-org/glm-4-9b)              | 4 × 9B = 36B   | 18× larger            |
| Prometheus2 (Kim et al., 2024) | 0.410    | 4 (prometheus-7b-v2.0)            | 4 × 7B = 28B   | 14× larger            |
| GPT-5 (Achiam et al., 2023)    | 0.581    | -                                 | -              | -                     |
| LoMTL (ours)                   | 0.601    | 1 (google/gemma-2-2b-it)          | 2B             | -                     |

Table 5: Comparison of macro-F1 scores and model parameters between our automated evaluation model (i.e., LoMTL) and the top-performing teams in the BEA shared task and LLM-as-a-judge models across four dimensions (aka. tasks). Note that, (1) for a fair comparison, we include only the teams that participated in all four tracks of the shared task, and (2) since GPT-5 is a closed-source model, its parameter details are not publicly available.

| SYSTEM PROMPT                                                                                                                                                                                                                                                                                                       |  |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--|
| You are an expert evaluator of AI tutors.<br>For the given ### Task, ### Task Definition, ### Label Definition, ### Conversation History and ### Tutor Response, assess the pedagogical appropriateness of the Tutor Response.<br>Output exactly one label without additional text: Yes, "No", or "To some extent". |  |

| TASK DEFINITIONS       |                                                                                                                     |
|------------------------|---------------------------------------------------------------------------------------------------------------------|
| Mistake Identification | Has the tutor identified or recognized a mistake in a student's response?                                           |
| Mistake Location       | Does the tutor's response accurately pinpoint the location of a genuine mistake?                                    |
| Providing Guidance     | Does the tutor offer correct and relevant guidance, such as an explanation, elaboration, hint, examples, and so on? |
| Actionability          | Is it clear from the tutor's feedback what the student should do next?                                              |

| LABEL DEFINITIONS      |                                                                                                                                                                                                                                              |
|------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Mistake Identification | <b>Yes:</b> The tutor correctly identified the mistake in the student's response.<br><b>To some extent:</b> The tutor partially recognized the mistake but did not fully capture it.<br><b>No:</b> The tutor failed to identify any mistake. |
| Mistake Location       | <b>Yes:</b> The tutor accurately points to the exact mistake and its location.<br><b>To some extent:</b> The tutor points to a mistake but imprecisely or partially.<br><b>No:</b> The tutor fails to indicate the mistake or its location.  |
| Providing Guidance     | <b>Yes:</b> The tutor provides correct and relevant guidance, hints, examples, or explanation.<br><b>To some extent:</b> The guidance is partially correct or not fully helpful.<br><b>No:</b> The tutor fails to provide relevant guidance. |
| Actionability          | <b>Yes:</b> It is clear what the student should do next.<br><b>To some extent:</b> The next steps are somewhat unclear or incomplete.<br><b>No:</b> The feedback does not indicate any actionable steps.                                     |

| FINAL PROMPT STRUCTURE                                                                                                                                                                                                                                                                                                                                    |  |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--|
| <pre> content = (     f'{cls.SYSTEM_PROMPT_WITH_LABEL}\n\n'     f'### Task: {task}\n\n'     f'### Task Definition: {task_def}\n\n'     f'### Label Definition: \n{label_def_str}\n\n'     f'### Conversation History: {conversation.strip()}\n\n'     f'### Tutor Response: {response.strip()}\n\n'     f'Now provide the classification label.' ) </pre> |  |

Figure 4: Overview of the prompt components, their associated definitions and details, and the final prompt structure used in LoMTL.

precision (0.58), suggesting a more recall-oriented behavior. On the demonstration subset, GPT-5 achieves the highest precision (0.60) and recall (0.61), whereas LoMTL remains competitive (0.58 precision, 0.56 recall) and substantially outperforms Prometheus2. Overall, these results show that LoMTL maintains a balanced precision-recall trade-off, reinforcing its robustness across evaluation dimensions.

## C.2 Human Annotation Analysis

We collected 95 annotations for single-tutor responses and 115 for pairwise comparisons. Since annotators were free to choose any dialogues and models, the number of annotations per tutor is not uniform. In pairwise comparisons, Gemini was selected most frequently – 49 out of 115 comparisons included Gemini as one of the tutors. Other tutors appeared in 17-36 comparisons, except for Novice, which was chosen only five times. Over-

| Category                        | Parameter                 | Value                                                                       |
|---------------------------------|---------------------------|-----------------------------------------------------------------------------|
| <b>Common Settings</b>          |                           |                                                                             |
| Model                           | MODEL_NAME                | google/gemma-2-2b-it                                                        |
| Task Dimensions                 | DIMENSIONS                | Mistake_Identification, Mistake_Location, Providing_Guidance, Actionability |
| Input Length                    | MAX_LENGTH                | 1024                                                                        |
| Prompt                          | include_label_definitions | Enabled                                                                     |
| <b>Training-Only Settings</b>   |                           |                                                                             |
| Batching                        | BATCH_SIZE                | 4                                                                           |
| Batching                        | GRAD_ACCUM                | 1                                                                           |
| Training Schedule               | EPOCHS                    | 3                                                                           |
| Training Schedule               | LEARNING_RATE             | 1e-4                                                                        |
| Training Schedule               | WEIGHT_DECAY              | 0.1                                                                         |
| Logging                         | LOGGING_STEPS             | 50                                                                          |
| Saving                          | SAVE_STEPS                | 300                                                                         |
| Evaluation Cycle                | EVAL_STEPS                | 300                                                                         |
| Oversampling                    | OVERSAMPLE_METHOD         | "random"                                                                    |
| Metric for best model           | METRIC_FOR_BEST           | "eval_loss"                                                                 |
| LoRA                            | LORA_R                    | 8                                                                           |
| LoRA                            | LORA_ALPHA                | 16                                                                          |
| LoRA                            | LORA_DROPOUT              | 0.1                                                                         |
| Early Stopping                  | EARLY_PATIENCE            | 5                                                                           |
| Early Stopping                  | EARLY_THRESHOLD           | 0.0                                                                         |
| <b>Evaluation-Only Settings</b> |                           |                                                                             |
| Generation                      | TEMPERATURE               | 1.0                                                                         |

Table 6: Summary of training and evaluation configurations along with their corresponding parameter names.

all, Expert responses were preferred most often, winning in 60% of the cases in which the Expert’s response appeared. Sonnet and Llama-3.1-405B were also selected more than half of the time. Novice’s responses never won.

In the single-tutor mode, 37 annotations marked responses as *helpful*, 31 as *not helpful*, and 27 as *to some extent helpful*. The most helpful responses came from Expert, Sonnet, Gemini, and Mistral, each of which was rated *helpful* in at least half of the corresponding cases. The least helpful responses were from Phi3 and Llama-3.1-8B. Mistral and Llama-3.1-405B were the only tutors without any *not helpful* annotations, as they had the highest proportion of *to some extent helpful* ratings.

| Model        | Full Test Set |             |             |             | Demonstration Set |             |             |             |
|--------------|---------------|-------------|-------------|-------------|-------------------|-------------|-------------|-------------|
|              | Accuracy      | Macro-F1    | Precision   | Recall      | Accuracy          | Macro-F1    | Precision   | Recall      |
| LoMTL (ours) | <b>0.72</b>   | <b>0.60</b> | <b>0.63</b> | 0.59        | <b>0.68</b>       | 0.55        | 0.58        | 0.56        |
| Prometheus2  | 0.47          | 0.41        | 0.44        | 0.45        | 0.41              | 0.34        | 0.38        | 0.36        |
| GPT-5        | 0.66          | 0.58        | 0.58        | <b>0.60</b> | 0.66              | <b>0.59</b> | <b>0.60</b> | <b>0.61</b> |

Table 7: Accuracy, macro-F1, macro-precision and macro-recall scores (averaged across Mistake Identification (MI), Mistake Location (ML), Providing Guidance (PG), and Actionability (AC)) for our LoMTL model, Prometheus2, and GPT-5 on the full test set from Kochmar et al. (2025) and on the demonstration set. Best results are shown in **bold**.

**Informed consent form:** In this study, we will ask you to look through a small set \* of tutorial dialogues taking place between a student and a tutor (in some cases the tutor is an AI-based tutoring model) in the mathematical domain at the middle-school level, where students made mistakes. You will be asked to consider quality judgements by various evaluation models on tutors' responses.

1. Participant selection criteria: There are no specific criteria, as we believe anyone who has done math at school is qualified to participate.
2. Anonymity: The form does not collect or store any personally identifiable information.
3. Use of your responses: We may use your free-form feedback to improve the models or the demo. All quality assessment scores will be used only in their aggregated form and only for research purposes. No individual responses will be publicly shared.
4. Voluntary basis: Participation in this study is completely voluntary. You can withdraw from this study at any point – in that case, the data submitted by you will be deleted and will not be used for any further analysis.

I understand and agree to the conditions of the study

Figure 5: Informed consent form that participants were required to accept before proceeding with their feedback and annotations.

### C.3 Full Questionnaire

Below are the informed consent forms that participants were required to read and accept, as well as the full questionnaire.

#### Background

| Item                                                                                                      | Question / Response Options                                                     |
|-----------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------|
| What is your highest qualification?                                                                       | Response options: <i>Bachelor's degree; Master's degree; PhD degree; Other.</i> |
| Do you have teaching experience (e.g., lecturing, supervising or mentoring students, TA-ing, or similar)? | Response options: <i>Yes; No.</i>                                               |
| Have you ever used an AI tutor before?                                                                    | Response options: <i>Yes; No.</i>                                               |

#### Automated Evaluation

| Item                                                                                    | Question / Response Options                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|-----------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Instructions                                                                            | This tab gives you an opportunity to select among 10 short dialogues on relatively simple (no higher than middle-school level) math problems, check students' misconceptions, and explore various human and AI tutor responses. The quality of these responses is evaluated using a fine-tuned evaluation model.                                                                                                                                                                                                               |
| Steps                                                                                   | <ol style="list-style-type: none"> <li>1) Explore at least 5 different dialogues (the correct answer is also available to help you spot the student's mistake).</li> <li>2) For each dialogue, check at least 2 different tutor responses.</li> <li>3) Rate each response as <i>Helpful</i>, <i>Not Helpful</i>, or <i>To some extent</i>.</li> <li>4) Select at least one quality dimension to view the model's assessment of that response.</li> <li>5) For each dialogue, use the comparison mode at least once.</li> </ol> |
| How many dialogues have you checked?                                                    | Free-form response.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
| How often did you agree with this model's assessment for a single tutor response?       | Scale: 1 = almost never, 2 = less than half of the time, 3 = about half of the time, 4 = more than half of the time, 5 = almost always.                                                                                                                                                                                                                                                                                                                                                                                        |
| Response options: 1, 2, 3, 4, 5.                                                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| How often did you agree with this model's ranking of the tutors in the comparison mode? | Same scale as above.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
| Response options: 1, 2, 3, 4, 5.                                                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| On a scale from 1 to 5, how easy was it to use the "Automated Evaluation" tab?          | Scale: 1 = not easy at all, 5 = very easy.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Response options: 1, 2, 3, 4, 5.                                                        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| Please feel free to give us any further feedback on this tab.                           | Free-form response.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |

#### LLM Evaluation

| Item                                                                         | Question / Response Options                                                                                                                                                                                                                                                                                                                                                                                                       |
|------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Instructions                                                                 | This tab provides access to the same dialogues and tutor responses as in Tab 1, but this time they are evaluated using LLMs as judges. You can choose between GPT-5 and Prometheus.                                                                                                                                                                                                                                               |
| Steps                                                                        | <ol style="list-style-type: none"> <li>1) Explore at least 5 different dialogues (they may be the same ones as before).</li> <li>2) For each dialogue, check at least 2 different tutor responses.</li> <li>3) Use each LLM-as-judge at least twice on different dialogues.</li> <li>4) Select at least one quality dimension for each response.</li> <li>5) For each dialogue, use the comparison mode at least once.</li> </ol> |
| How many dialogues have you checked?                                         | Free-form response.                                                                                                                                                                                                                                                                                                                                                                                                               |
| How often did you agree with GPT-5's assessment for a single tutor response? | Scale: 1 = almost never, 2 = less than half of the time, 3 = about half of the time, 4 = more than half of the time, 5 = almost always.                                                                                                                                                                                                                                                                                           |

| Item                                                                                                                                                                    | Question / Response Options                                                                     |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------|
| Response options: 1, 2, 3, 4, 5.                                                                                                                                        |                                                                                                 |
| How often did you agree with GPT-5's ranking of the tutors in the comparison mode?                                                                                      | Same scale as above.                                                                            |
| Response options: 1, 2, 3, 4, 5.                                                                                                                                        |                                                                                                 |
| How often did you agree with Prometheus' assessment for a single tutor response?                                                                                        | Same scale as above.                                                                            |
| Response options: 1, 2, 3, 4, 5.                                                                                                                                        |                                                                                                 |
| How often did you agree with Prometheus' ranking of the tutors in the comparison mode?                                                                                  | Same scale as above.                                                                            |
| Response options: 1, 2, 3, 4, 5.                                                                                                                                        |                                                                                                 |
| Which evaluation model did you perceive to be more accurate in its judgments of tutor responses GPT-5 or the fine-tuned model from Tab 1 ("Automated Evaluation")?      | Response options: <i>The model from Tab 1; GPT-5; Hard to say: they perform similarly.</i>      |
| Which evaluation model did you perceive to be more accurate in its judgments of tutor responses Prometheus or the fine-tuned model from Tab 1 ("Automated Evaluation")? | Response options: <i>The model from Tab 1; Prometheus; Hard to say: they perform similarly.</i> |
| On a scale from 1 to 5, how easy was it to use the "LLM Evaluation" tab?                                                                                                | Scale: 1 = not easy at all, 5 = very easy.                                                      |
| Response options: 1, 2, 3, 4, 5.                                                                                                                                        |                                                                                                 |
| Any other feedback is welcome.                                                                                                                                          | Free-form response.                                                                             |

## Visualizer

| Item                                                                       | Question / Response Options                                                                                                                                                                                                                                                                                            |
|----------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Instructions                                                               | This tab visualizes statistics from the full dataset of tutorial dialogues and annotated tutor responses. The dataset contains 300 dialogues with responses from 9 tutors (except for the Novice tutor, who has annotations for 76 dialogues). Participants are asked to explore visualizations for at least 2 tutors. |
| On a scale from 1 to 5, how informative did you find these visualizations? | Scale: 1 = not informative at all, 5 = very informative.                                                                                                                                                                                                                                                               |
| Response options: 1, 2, 3, 4, 5.                                           |                                                                                                                                                                                                                                                                                                                        |
| On a scale from 1 to 5, how easy was it to use the "Visualizer" tab?       | Scale: 1 = not easy at all, 5 = very easy.                                                                                                                                                                                                                                                                             |
| Response options: 1, 2, 3, 4, 5.                                           |                                                                                                                                                                                                                                                                                                                        |
| We welcome any further feedback.                                           | Free-form response.                                                                                                                                                                                                                                                                                                    |

## D User Interface (UI) Details

The screenshot displays the 'AI Tutor Evaluation Platform' interface. At the top left is the 'EDU<sub>NLP</sub>' logo. The main title is 'AI Tutor Evaluation Platform', with a sub-note 'Powered by MBZUAI EduNLP'. The navigation bar includes three tabs: 'Automated Evaluation' (active), 'LLM Evaluation', and 'Visualizer'. Below the navigation bar is a green header for the 'Automated Evaluation' section, with the subtext 'Select problem topics and view automated evaluation scores'. The main content area features an 'Overview' section with four metrics: Topics (10), Models (9), Dimensions (4), and Conversations (10). Below this is a toggle for 'Enable Tutor Comparison Mode (Compare Two Tutors)'. There are two dropdown menus: 'Problem Topic' and 'Tutor'. The 'Evaluation Dimensions' section includes four checkboxes: 'Mistake Identification', 'Mistake Location', 'Providing Guidance', and 'Actionability', along with 'Select All' and 'Clear All' buttons. A blue button at the bottom reads 'Get Auto-Evaluation Results'.

Figure 6: Overview of the UI and the Automated Evaluation Tab.

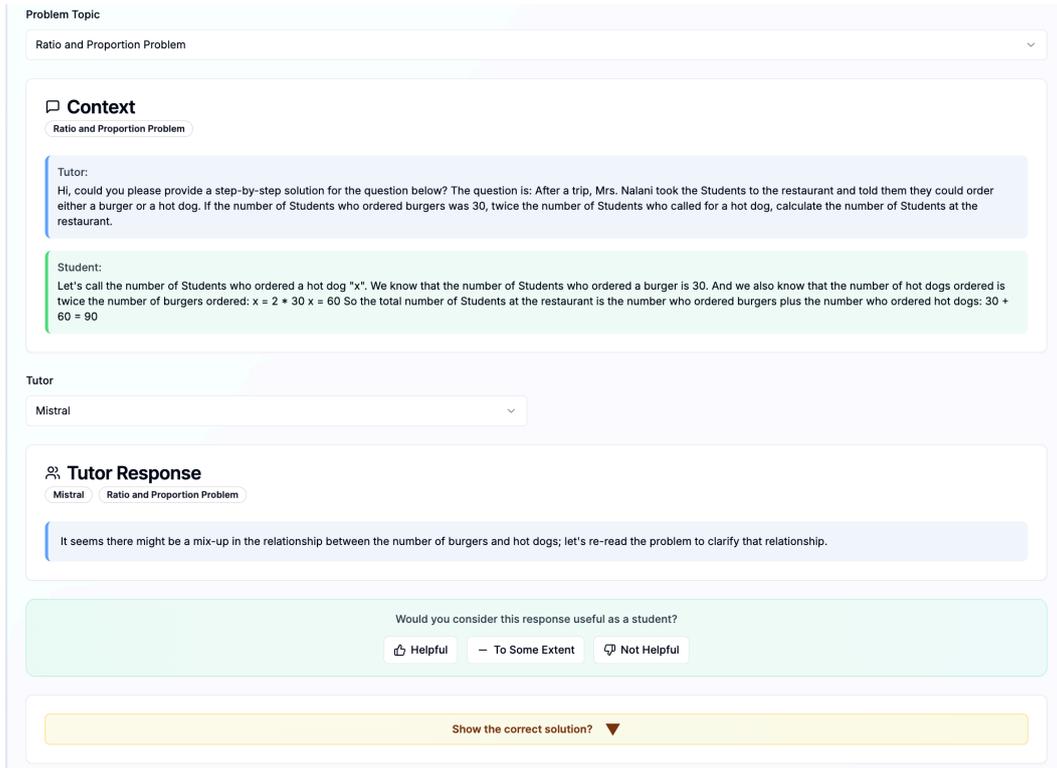


Figure 7: Interface showcasing the selected problem topic with Context, automated Tutor Response, student feedback options, and ground truth verification panel.

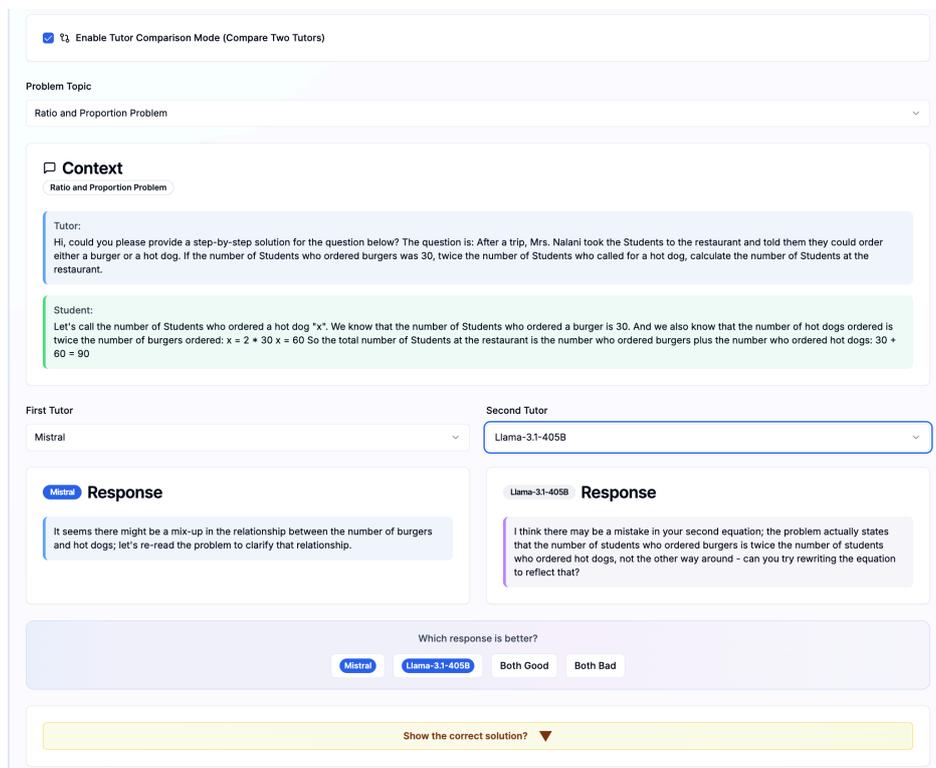


Figure 8: UI displaying the enabled Tutor Comparison Mode, allowing users to compare responses from any two selected tutors for the selected problem topic, along with the Context block, selected tutor responses, feedback options, and ground truth verification option.

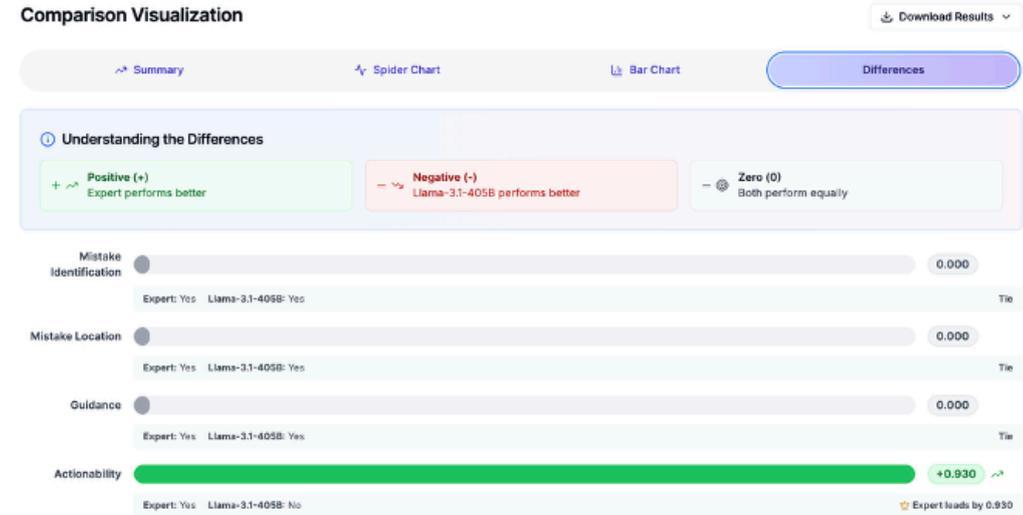
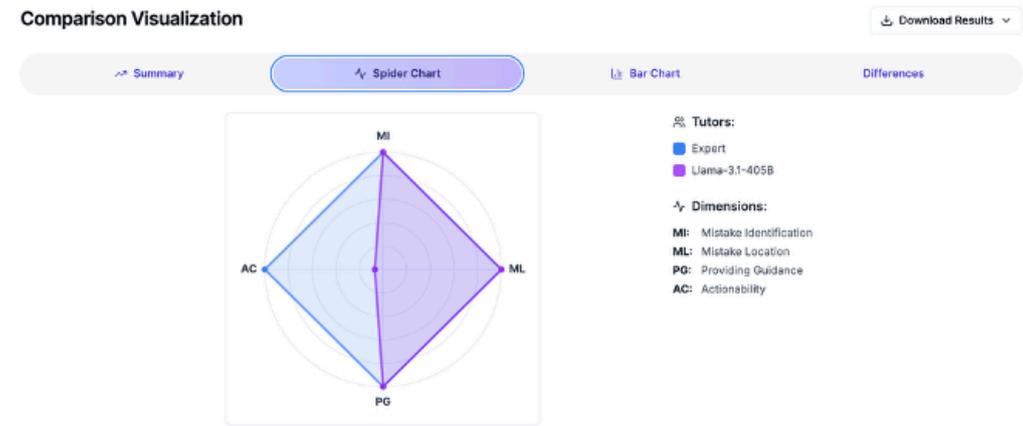
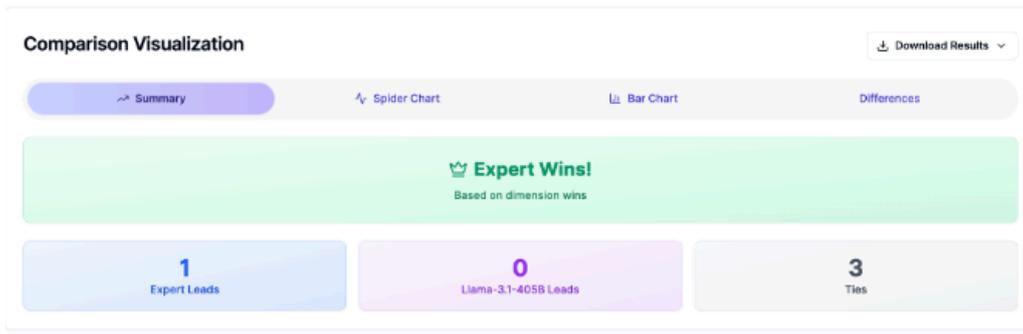


Figure 9: Displayed Tutor Comparison Visualization Panel showcasing Summary metrics, Spider Chart, Bar Chart, and Differences views for the chosen evaluation results.

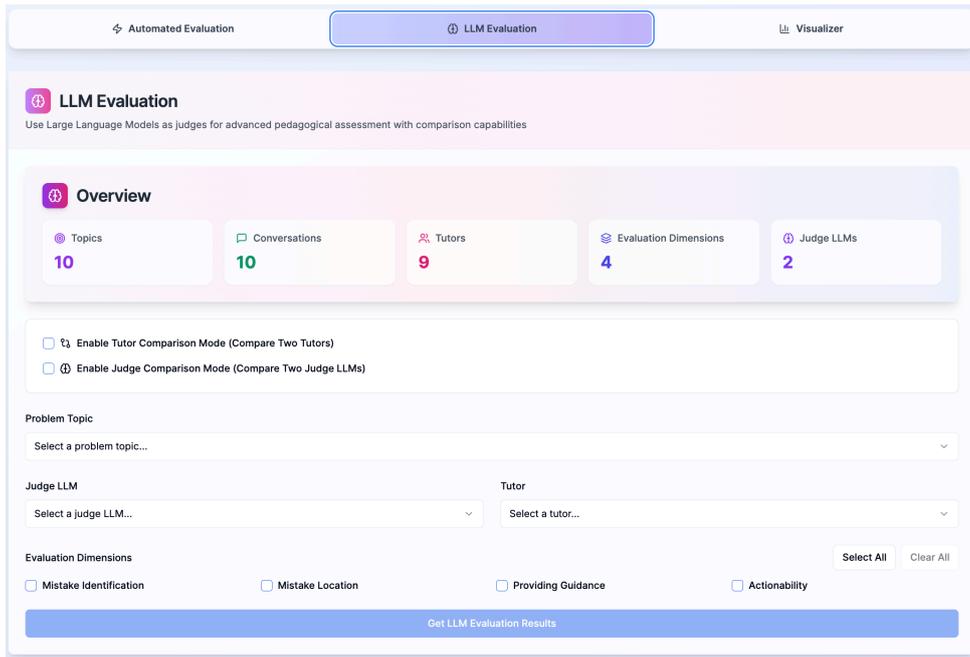


Figure 10: Overview of the LLM Evaluation module showcasing the dashboard panel with statistics on topics, conversations, tutors, and evaluation dimensions. The interface also highlights the provided options for Tutor Comparison Mode and Judge Comparison Mode for advanced pedagogical assessment.

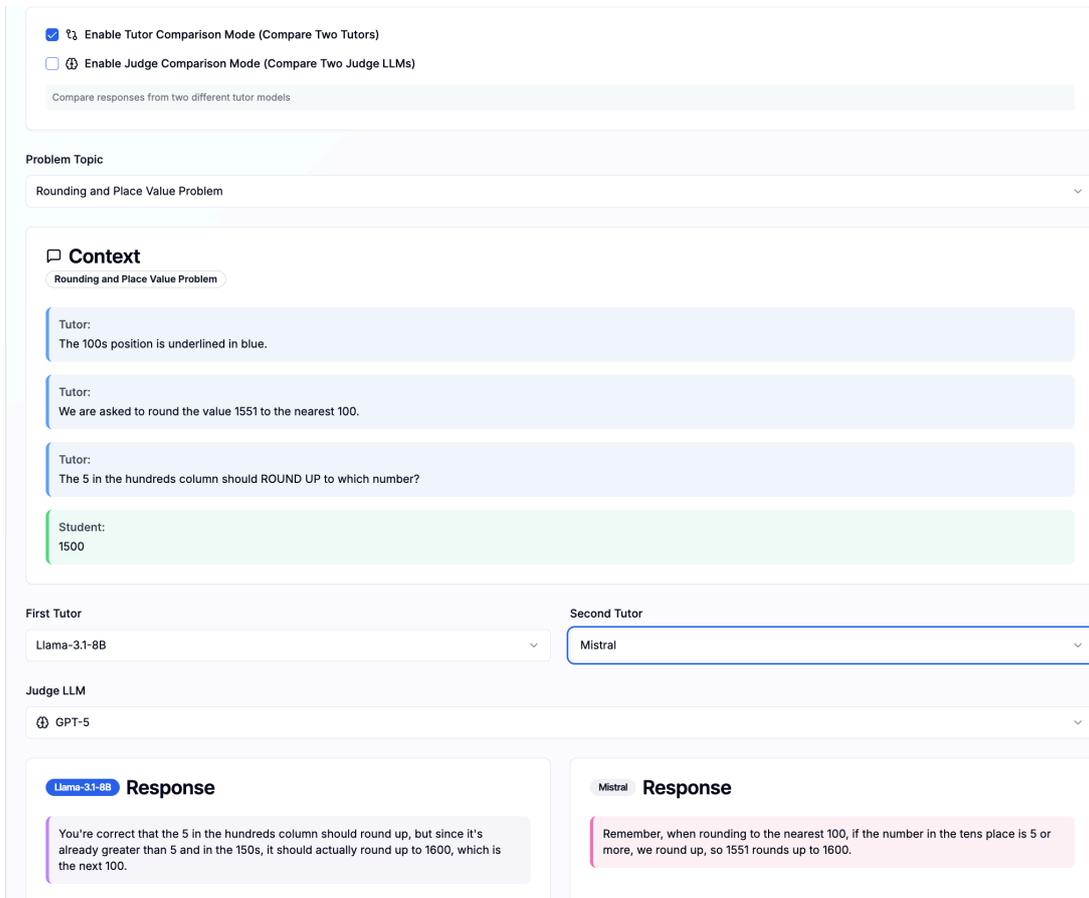


Figure 11: User Interface displaying the enabled Tutor Comparison Mode within the LLM Evaluation module, showing the selected problem, tutor responses from two tutors, and the judge LLM for evaluation.

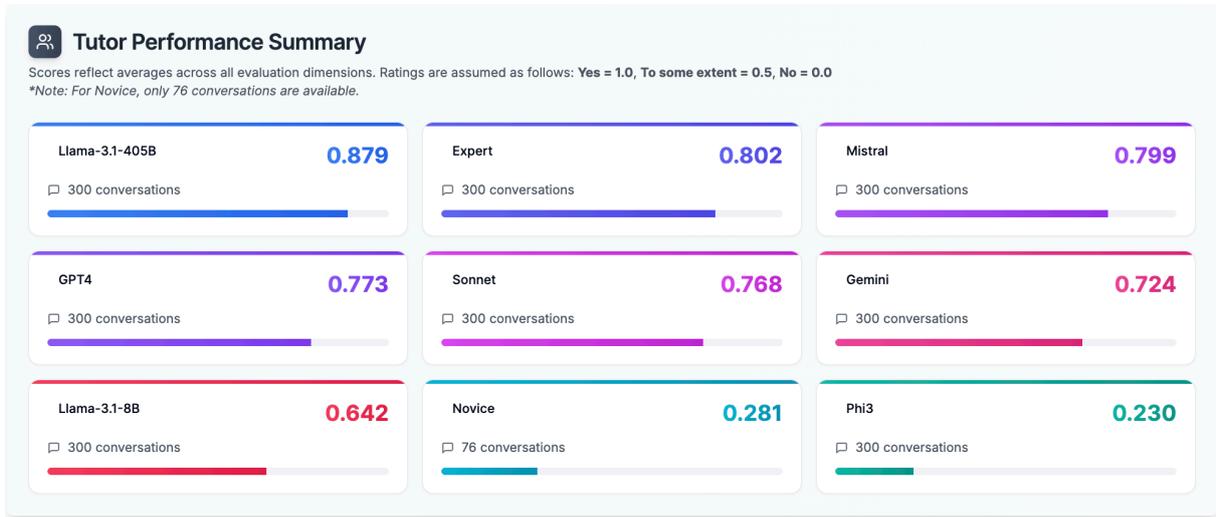


Figure 12: Tutor Performance Summary panel in the Visualizer module, displaying the aggregated evaluation scores for each tutor across all assessment dimensions within the MRBench development dataset.

**Visualization Controls**

**Select Tutors to Compare** Select All | Clear All

Sonnet     
  Llama-3.1-8B     
  Llama-3.1-405B     
  GPT4     
  Mistral  
 Expert     
  Gemini     
  Phi3     
  Novice

Selected: 0 model(s)

**Select Evaluation Dimensions for Spider Plot** Select All | Clear All

Mistake Identification  
 Mistake Location  
 Providing Guidance  
 Actionability

Selected: 0 dimension(s)

**Select Evaluation Dimensions for Bar Plot**

Choose a dimension ▼

Figure 13: Interface of the Visualization Controls Panel, showing configurable options for selecting tutors and evaluation dimensions to generate comparative spider chart and bar plot visualizations.

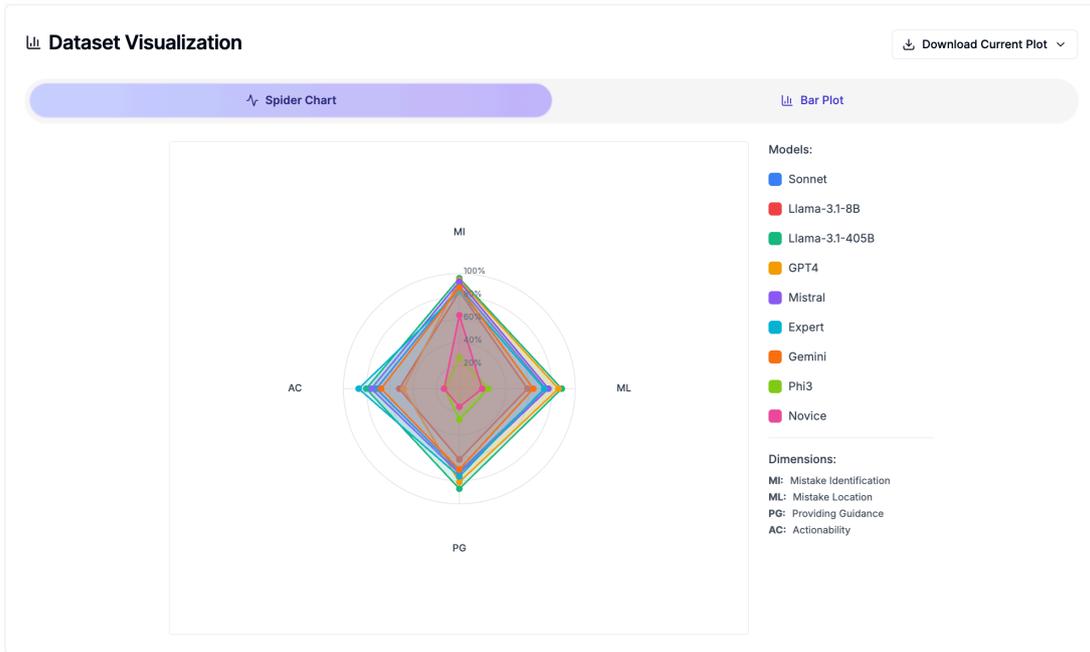


Figure 14: The spider chart representing tutors performance across selected evaluation dimensions, based on configurations chosen in the Visualization Controls Panel.

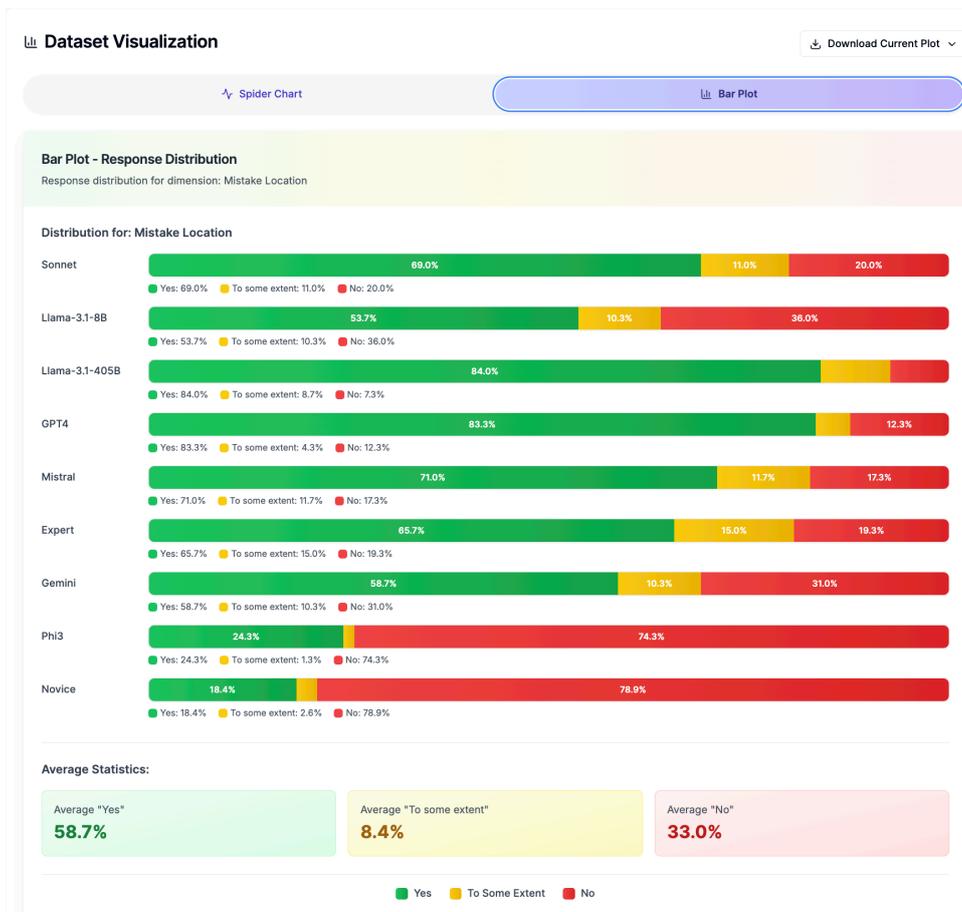


Figure 15: The bar plot representing tutors performance across selected evaluation dimension, based on configurations chosen from the Visualization Controls Panel.

# EvalSense: A Framework for Domain-Specific LLM (Meta-)Evaluation

**Adam Dejl**  
Imperial College London\*  
Department of Computing  
ad5518@ic.ac.uk

**Jonathan Pearson**  
NHS England  
Transformation Directorate  
jonathanpearson@nhs.net

## Abstract

Robust and comprehensive evaluation of large language models (LLMs) is essential for identifying effective LLM system configurations and mitigating risks associated with deploying LLMs in sensitive domains. However, traditional statistical metrics are poorly suited to open-ended generation tasks, leading to growing reliance on LLM-based evaluation methods. These methods, while often more flexible, introduce additional complexity: they depend on carefully chosen models, prompts, parameters, and evaluation strategies, making the evaluation process prone to misconfiguration and bias. In this work, we present EvalSense, a flexible, extensible framework for constructing domain-specific evaluation suites for LLMs. EvalSense provides out-of-the-box support for a broad range of model providers and evaluation strategies, and assists users in selecting and deploying suitable evaluation methods for their specific use-cases. This is achieved through two unique components: (1) an *interactive guide* aiding users in evaluation method selection and (2) *automated meta-evaluation tools* that assess the reliability of different evaluation approaches using perturbed data. We demonstrate the effectiveness of EvalSense in a case study involving the generation of clinical notes from unstructured doctor-patient dialogues, using a popular open dataset. All code, documentation, and assets associated with EvalSense are open-source and publicly available at <https://github.com/nhsengland/evalsense>.

## 1 Introduction

Backed by training on unprecedentedly large quantities of data, large language models (LLMs) have radically advanced the field of machine learning and demonstrated a wide range of impressive capabilities across diverse domains (Bubeck et al., 2023; Van Veen et al., 2024; Luo et al., 2025; McDuff et al., 2025). While these results suggest that

LLMs have the potential to deliver substantial benefits, their use also entails significant risks, including hallucinations (Huang et al., 2025), omissions of crucial information (Busch et al., 2025), unintended disclosure of sensitive personal data (Das et al., 2025), and vulnerability to harmful instructions (Das et al., 2025). Rigorous evaluation of LLMs has been proposed as a key strategy for mitigating these risks and ensuring that LLM-based systems perform reliably on their assigned tasks (WHO, 2023; Ong et al., 2024).

However, reliable evaluation of open-ended texts produced by LLMs remains challenging as a result of the unstructured and complex nature of these texts. Due to the inadequacy of standard statistical metrics, the community has increasingly adopted LLM-as-a-judge approaches (Liu et al., 2023; Fu et al., 2024; Kim et al., 2024a,b), which use LLMs themselves to assess model outputs. These methods tend to be more effective at capturing content-related nuances and generally achieve higher correlations with human judgements (Zheng et al., 2023). Yet, the reliability of LLMs as evaluators may vary depending on the considered task, LLM judge and the used evaluation strategy (Murugadoss et al., 2025; Tan et al., 2025; Han et al., 2025). This motivates the need to carefully choose the evaluation approach suitable for the specific domain and to rigorously *meta-evaluate* its effectiveness (i.e., to evaluate the evaluator), steps that are often neglected in the existing evaluation pipelines.

Several open-source toolkits and frameworks for evaluating LLMs have been introduced, such as lm-evaluation-harness (Gao et al., 2024), OpenCompass (Contributors, 2023), LightEval (Habib et al., 2023), Inspect (UK AI Security Institute, 2024) and Unitxt (Bandel et al., 2024). However, while these tools provide useful infrastructure for running standardised benchmarks or implementing specific evaluation workflows, they typically do not aid users in selecting appropriate

\*Work done while at NHS England.

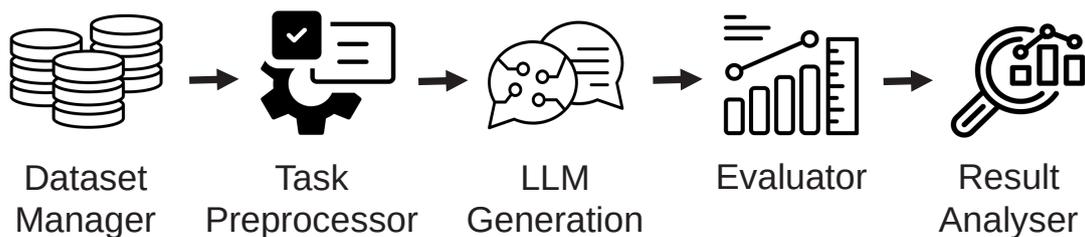


Figure 1: Overview of the LLM evaluation pipeline implemented in EvalSense<sup>1</sup>. After data loading and task-specific preprocessing, model outputs are generated and scored using different evaluators. Result analysers summarise outcomes across experiments, identify higher-level patterns, and support meta-evaluation.

methods or in quantitatively measuring the effectiveness of these methods for a specific domain and task through meta-evaluation.

In response to these gaps, we introduce EvalSense, a highly flexible software framework that enables users to systematically evaluate LLMs on custom datasets. EvalSense offers two key features to help users navigate the spectrum of available evaluation methods:

1. It includes an *interactive evaluation guide*<sup>2</sup>, which prompts users to specify their tasks along with the associated risks and requirements, and then suggests appropriate evaluation strategies. After a subset of methods is selected, the guide generates a coverage report indicating whether the chosen methods comprehensively cover the specified risks and requirements (Figure 2a).
2. EvalSense incorporates *automated meta-evaluation tools* that leverage controlled perturbations to validate evaluator reliability on the user’s own dataset. These tools systematically degrade specific aspects of the output texts, verifying the degree to which these changes are reflected in the scores produced by the different evaluation techniques.

In addition to these features, EvalSense also supports systematic experimentation, a broad range of local and API model providers, configuring evaluations through a graphical user interface (Figure 2b), high-level result analysis and complex generation workflows.

<sup>1</sup>Icons by Noun Project, authors Srinivas Agra, Iconiqu, Gonza Monta, Keyy Creative and suhaiba, CC BY 3.0.

<sup>2</sup>Available on the EvalSense website at <https://nhsengland.github.io/evalsense/>.

To demonstrate and assess the capabilities of EvalSense, we apply it to a realistic evaluation task using ACI-Bench (Yim et al., 2023), which involves generation of structured clinical notes from doctor-patient dialogues. Using EvalSense’s meta-evaluation tools, we demonstrate a non-trivial disparity in the quality of the scores produced by different evaluation methods. This demonstrates the importance of careful method selection and configuration, a process our framework is specifically designed to support.

Overall, we hope that EvalSense contributes to advancing best practices in LLM evaluation by systemizing the process of choosing between different evaluation strategies, both through the interactive guidance provided by the EvalSense guide and the quantitative meta-evaluation supported by the associated open-source library.

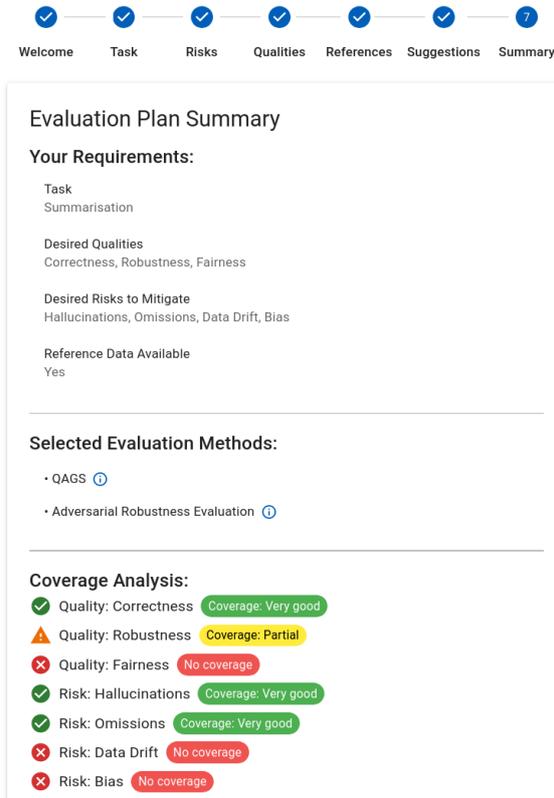
## 2 Background and Related Work

### 2.1 Evaluation Metrics

The growing use of machine learning models for text generation has led to the development of a wide range of evaluation techniques. Broadly, these can be categorised into three groups: traditional statistical metrics, LLM-as-a-judge methods and hybrid approaches.

**Statistical metrics** rely on direct, deterministic comparison of text units extracted from the evaluated text to the ground-truth reference. While still in use, these approaches are often overly simple to reliably assess the quality of open-ended texts. Examples of such metrics include the BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) scores. **LLM-as-a-judge** methods leverage the general capabilities of LLMs to assess generated texts, mitigating many of the drawbacks associated with sta-

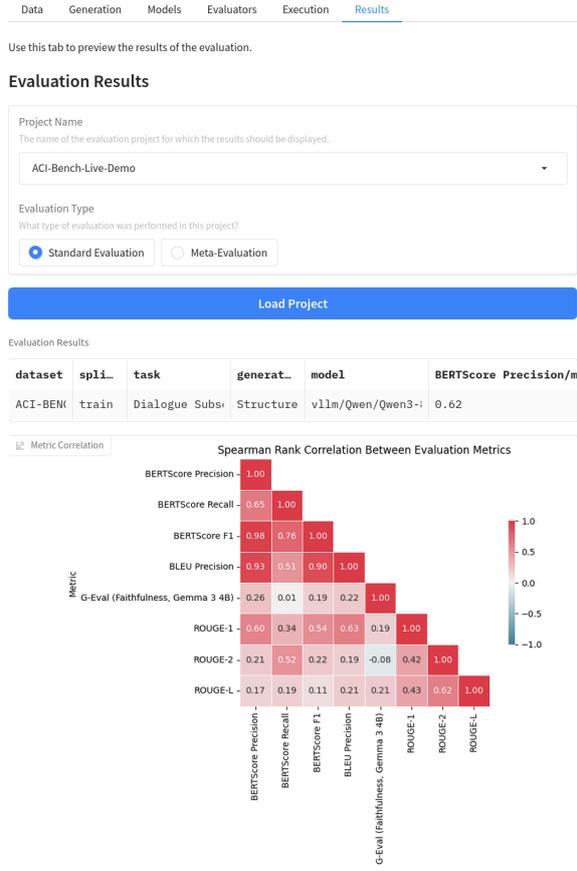
## LLM Evaluation Guide



(a) LLM Evaluation Guide

## EvalSense

To run an evaluation, configure its settings on the individual tabs and start it from the Execution tab. For EvalSense documentation and guidance regarding the available evaluation metrics, please visit the [EvalSense homepage](#).



(b) EvalSense user interface

Figure 2: (a) EvalSense’s LLM Evaluation Guide assists users in selecting suitable evaluation methods based on task-specific risks and requirements. The final evaluation plan summary highlights any risks and requirements not fully covered by the selected methods. The guide is available at <https://nhsengland.github.io/evalsense/guide>. (b) The web-based user interface provided by the EvalSense library can be used to configure and execute evaluations, as well as to view their results. Alternatively, this can be done through code after importing the library.

tistical metrics (Zheng et al., 2023). However, the effectiveness of these methods is highly sensitive to the choice of model, prompt formulation, and general evaluation protocol. Notable examples include G-Eval (Liu et al., 2023) and GPTScore (Fu et al., 2024).

**Hybrid methods** also make use of pre-trained models, but only use these models for targeted subtasks as part of a structured evaluation framework. For instance, BERTScore (Zhang et al., 2020) compares texts based on contextual embeddings, while QAGS (Wang et al., 2020) assesses factual consistency using question answering.

## 2.2 LLM Evaluation Toolkits

A number of open-source toolkits have been introduced to support the evaluation of LLMs. Frame-

works such as such as lm-evaluation-harness (Gao et al., 2024), OpenCompass (Contributors, 2023), and LightEval (Habib et al., 2023) primarily focus on benchmarking models against standardised tasks and datasets. While some of these tools also support evaluation on custom data, they generally lack dedicated mechanisms for guiding evaluation design or assessing the suitability of selected methods in specific domains. The FreeEval framework (Yu et al., 2024) extends beyond benchmarking by incorporating human judgements, bias detection, contamination analysis, and case-by-case inspection. However, it does not support automated meta-evaluation or provide interactive tools for evaluation strategy selection.

Among existing tools, Inspect (UK AI Security Institute, 2024) is especially relevant to our

work. EvalSense uses Inspect as its basis, inheriting its support for multiple model providers, tool use, agentic workflows and detailed logging infrastructure. Nevertheless, EvalSense significantly expands on this foundation through a more versatile and extensible pipeline tailored to custom datasets, improved resource management, support for advanced evaluation methods (including sophisticated LLM-as-a-judge and hybrid approaches), and its unique focus on meta-evaluation and domain-specific guidance. The bespoke components of EvalSense are described in the following section.

### 3 EvalSense Pipeline

EvalSense implements a robust and customisable pipeline that manages the key steps of the evaluation process, from data management and preprocessing, through LLM generation, to final evaluation and result analysis (as illustrated in Figure 1). Uniquely, the generation and result analysis modules provide built-in support for meta-evaluating the reliability of the used metrics in addition to simply returning their scores. The overall design of the pipeline supports reusability and extensibility by making individual components easily replaceable, enabling the use of custom datasets, LLMs, or evaluation methods.

#### 3.1 Dataset Manager

Dataset managers are responsible for loading and generic preprocessing of the data on which the LLM is to be evaluated. For open-source datasets, this may involve downloading the data files from a publicly available repositories, while for internal datasets, the data will typically be loaded from a local file system or secure cloud storage. To simplify data management for custom datasets, EvalSense provides a base `DatasetManager` class that defines a general interface for data managers and implements helper methods for retrieving associated files based on paths specified in a dataset configuration file. These methods can be overridden to support more complex data loading and preprocessing workflows.

#### 3.2 Task Preprocessor

Task preprocessors implement any additional preprocessing steps that may be required to prepare the data for a specific task, as a single dataset may potentially support multiple such tasks. In many simple cases where no additional preprocessing is needed, users can rely on

the `DefaultTaskPreprocessor`, which acts as an identify function. For more complex scenarios, users can define a custom preprocessing function following the `TaskPreprocessingFunction` protocol, which can then be used with the standard `TaskPreprocessor`.

#### 3.3 LLM Generation Steps

After preparing the data for the task, the pipeline generates LLM outputs for evaluation using predefined generation steps. These typically involve prompting the model with specific system and user prompts. Optionally, the generation steps can incorporate more advanced strategies, such as enabling access to external tools (e.g., via the Model Context Protocol<sup>3</sup>), incorporating model self-critiques (Madaan et al., 2023), or using agentic workflows like ReAct (Yao et al., 2023). For the purposes of meta-evaluation, the LLM generation steps can also apply targeted perturbations degrading the quality of the output texts in predictable ways. In EvalSense, generation steps are defined via the `GenerationSteps` class, and the model configuration is specified using the `ModelConfig`.

#### 3.4 Evaluator

Evaluators implement automated methods for scoring model outputs based on predefined quality criteria. EvalSense includes several out-of-the-box score calculators, including:

- BLEU (`BleuPrecisionScoreCalculator`)
- ROUGE (`RougeScoreCalculator`)
- BERTScore (`BertScoreCalculator`)
- G-Eval (`GEvalScoreCalculator`)
- QAGS (`QagsScoreCalculator`)

These calculators can either be used independently or wrapped in an `Evaluator` class to be used as part of an evaluation pipeline. For convenience, EvalSense provides helper functions to easily initialise these evaluators (e.g., `get_bleu_evaluator` for BLEU). All key aspects of the evaluator configurations, such as the used prompts and models for LLM-as-a-judge approaches are fully customisable. Users may also implement new evaluators to be used as part of the pipeline.

<sup>3</sup><https://modelcontextprotocol.io/introduction>

### 3.5 Result Analyser

While evaluators produce fine-grained scores for the individual samples and summary metrics for each configuration evaluated by the pipeline, result analysers can be used to summarise the results from multiple such configurations and surface higher-level trends. EvalSense currently includes three main analysers: `TabularResultAnalyser` (for tabular summaries), `MetricCorrelationAnalyser` (for inter-metric correlation analysis), and `MetaResultAnalyser`. The last-mentioned analyser is crucial for the meta-evaluation capabilities of EvalSense, and can assess the consistency of metric scores either with different levels of automated perturbations (increasingly degrading a specific aspect of the evaluated texts to obtain the ground-truth output rankings) or human annotations. Meanwhile, the correlation analysis provided by the `MetricCorrelationAnalyser` can be particularly helpful for identifying similarities between different evaluation methods.

### 3.6 Pipeline

All components are integrated through the `Pipeline`, which schedules and executes the planned experiments (i.e., different configurations to be evaluated). These experiments can be declared individually or in batches using the `ExperimentConfig` and `ExperimentBatchConfig` data classes, enabling systematic evaluation sweeps. By default, the pipeline attempts to schedule the experiments in an optimal order to minimise the number of necessary model loads for local models, while also enabling users to resume any failed model generation tasks. As outlined above, the pipeline also supports automated meta-evaluation, performing controlled perturbations during the generation stage and assessing the reliability of different evaluation metrics during result analysis.

### 3.7 Project

All outputs, results, and metadata from pipeline execution are tracked through the `Project` class, which maintains a record of all experiments associated with a given project, their status, and links to the relevant logs. This class also provides a high-level interface through which the pipeline and result analysers can access and update these logs.

## 4 Evaluation Case Study

**Task Setup** To demonstrate EvalSense’s effectiveness, we apply it to LLM evaluation on the task of dialogue summarisation using the ACI-Bench dataset (Yim et al., 2023). In this setting, the LLM is tasked with generating a structured clinical note based on a doctor-patient dialogue transcript. Given that correctness and comprehensiveness of the generated notes are the most crucial qualities in this context, our evaluation primarily focuses on these aspects.

We used the 120 samples from the test partitions of the ACI-Bench dataset. Since we are using the original, unchanged dialogues from the dataset, the dataset manager and task preprocessor stages of our pipeline are mostly focused on loading the relevant samples without significant additional preprocessing.

For the LLM generation steps, we used the system and user prompts from (Kanithi et al., 2024). The user prompt specified the intended note structure and section headings, as more general instructions would make the output format ambiguous. We experimented with six different open-weight models: Llama 3.1 8B (Dubey et al., 2024), Phi 4 (Abdin et al., 2024), Qwen3 8B and Qwen3 14B (Yang et al., 2025), and Gemma 3 12B and Gemma 3 27B (Kamath et al., 2025). All models were run locally using vLLM (Kwon et al., 2023) in their default precisions, with the generation temperature set to 0.7, top-p sampling value of 0.95 and a seed of 42.

**Evaluation Setup** The case study involved a total of 13 variants of five major evaluators implemented in EvalSense: BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), BERTScore (Zhang et al., 2020), G-Eval (Liu et al., 2023) and QAGS (Wang et al., 2020). For ROUGE, we considered ROUGE-1, ROUGE-2 and ROUGE-L. The G-Eval metric was used with two different prompt variations: a detailed prompt providing thorough instructions on how to evaluate a note and a brief prompt asking the model to evaluate general faithfulness and accuracy. We also experimented with different G-Eval judge models: Llama 3.1 8B, Qwen3 14B and Gemma 3 27B. For the QAGS metric, we considered two different versions: ternary QAGS that generates questions requiring ternary responses and judge QAGS using more open-ended questions with an LLM judge comparing the responses. Both considered variants of the QAGS score used Llama

Table 1: Results from LLM evaluation on the ACI-Bench case study task using statistical and hybrid evaluation methods. Best results are bolded, second-best results are underlined.

| Model        | BLEU         | ROUGE-1      | ROUGE-2      | ROUGE-L      | BERTScore F1 | Ternary QAGS | Judge QAGS   |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Gemma 3 12B  | <b>0.128</b> | <u>0.514</u> | <u>0.207</u> | <b>0.300</b> | 0.666        | 0.842        | 0.817        |
| Gemma 3 27B  | 0.120        | 0.502        | 0.198        | 0.291        | <u>0.668</u> | <b>0.846</b> | <b>0.821</b> |
| Llama 3.1 8B | <u>0.127</u> | <b>0.534</b> | <b>0.221</b> | <u>0.294</u> | 0.662        | 0.806        | 0.789        |
| Phi 4 14B    | 0.120        | 0.504        | <u>0.207</u> | <u>0.290</u> | <b>0.670</b> | 0.832        | 0.811        |
| Qwen3 8B     | 0.091        | 0.468        | 0.174        | 0.259        | 0.648        | 0.818        | 0.784        |
| Qwen3 14B    | 0.100        | 0.451        | 0.170        | 0.259        | 0.640        | 0.810        | 0.793        |

Table 2: Results from LLM evaluation on the ACI-Bench case study task using different variants of G-Eval. Best results are bolded, second-best results are underlined

| Model        | Brief Gemma 3 | Det. Gemma 3 | Brief Llama 3.1 | Det. Llama 3.1 | Brief Qwen3  | Det. Qwen3   |
|--------------|---------------|--------------|-----------------|----------------|--------------|--------------|
| Gemma 3 12B  | <b>0.929</b>  | 0.904        | <u>0.847</u>    | 0.834          | 0.777        | <u>0.665</u> |
| Gemma 3 27B  | <u>0.926</u>  | <b>0.916</b> | <b>0.848</b>    | 0.840          | <u>0.798</u> | 0.640        |
| Llama 3.1 8B | 0.835         | 0.823        | 0.788           | 0.801          | 0.683        | 0.598        |
| Phi 4 14B    | 0.906         | 0.892        | 0.826           | 0.836          | 0.763        | 0.662        |
| Qwen3 8B     | 0.864         | 0.876        | 0.845           | <b>0.854</b>   | 0.775        | 0.630        |
| Qwen3 14B    | 0.885         | 0.899        | 0.843           | <u>0.849</u>   | <b>0.814</b> | <b>0.682</b> |

Table 3: Results from perturbation-based meta-evaluation of the different evaluation methods. Methods are ordered from best to worst.

| Method Name                     | Avg. Correlation |
|---------------------------------|------------------|
| G-Eval (Detailed, Gemma 3 27B)  | 0.999            |
| G-Eval (Brief, Gemma 3 27B)     | 0.998            |
| G-Eval (Detailed, Llama 3.1 8B) | 0.995            |
| G-Eval (Brief, Llama 3.1 8B)    | 0.992            |
| Ternary QAGS (Llama 3.1 8B)     | 0.982            |
| Judge QAGS (Llama 3.1 8B)       | 0.969            |
| G-Eval (Brief, Qwen 3 14B)      | 0.967            |
| G-Eval (Detailed, Qwen 3 14B)   | 0.924            |
| BERTScore F1                    | 0.431            |
| ROUGE-1                         | 0.323            |
| BLEU Precision                  | 0.296            |
| ROUGE-L                         | 0.232            |
| ROUGE-2                         | 0.049            |

3.1 8B as the underlying model.

For our meta-evaluation, we used a set of three prompts instructing the model to apply different levels of perturbations to the note: one rephrasing the note without changing its meaning, one introducing minor content changes and one significantly changing the meaning of the note. The used prompts are given in Appendix A.

**Results** The results of our evaluation are summarised in Tables 1 and 2. We can observe that there is substantial disagreement among the different methods, with no universally best-performing model. Without further information on each method’s reliability for this task, drawing definitive conclusions would be difficult.

However, based on the meta-evaluation results in Table 3, we can assign greater weight to G-Eval variants using Gemma 3 and Llama 3.1, as well as both QAGS versions. These methods consistently rank the Gemma 3 models highest, except for G-eval with Llama 3.1 using the detailed prompt. Notably, statistical metrics and BERTScore underperform compared to LLM-based methods.

## 5 Conclusion

In this paper, we introduced EvalSense, a novel framework for systematic evaluation of LLMs on custom tasks. Unlike other toolkits, which mainly focus on direct application of evaluation methods without providing principled ways to assess their suitability, EvalSense guides users in selecting evaluation approaches tailored to their specific domains and provides quantitative insights about the effectiveness of these approaches through meta-evaluation. We demonstrated its capabilities through a case study on structured clinical note generation from doctor-patient dialogues, showing that it supports robust evaluation even when different evaluation methods yield disagreeing results.

## Acknowledgments

We thank the UK AI Security Institute and the wider development team for their work on the Inspect framework, which serves as a basis for EvalSense.

## References

- Marah I Abidin, Jyoti Aneja, Harkirat S. Behl, Sébastien Bubeck, Ronen Eldan, and 1 others. 2024. [Phi-4 technical report](#). *CoRR*, abs/2412.08905.
- Elron Bandel, Yotam Perlitz, Elad Venezian, Roni Friedman, Ofir Arviv, Matan Orbach, Shachar Don-Yehiya, Dafna Sheinwald, Ariel Gera, Leshem Choshen, Michal Shmueli-Scheuer, and Yoav Katz. 2024. [Unitxt: Flexible, shareable and reusable data preparation and evaluation for generative AI](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: System Demonstrations)*, pages 207–215, Mexico City, Mexico. Association for Computational Linguistics.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, Marco Túlio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with GPT-4](#). *CoRR*, abs/2303.12712.
- Felix Busch, Lena Hoffmann, Christopher Rueger, Elon HC van Dijk, Rawen Kader, Esteban Ortiz-Prado, Marcus R. Makowski, Luca Saba, Martin Hadamitzky, Jakob Nikolas Kather, Daniel Truhn, Renato Cuocolo, Lisa C. Adams, and Keno K. Bressem. 2025. [Current applications and challenges in large language models for patient care: a systematic review](#). *Communications Medicine*, 5(1):26.
- OpenCompass Contributors. 2023. [OpenCompass: A universal evaluation platform for foundation models](#). <https://github.com/open-compass/opencompass>.
- Badhan Chandra Das, M. Hadi Amini, and Yanzhao Wu. 2025. [Security and privacy challenges of large language models: A survey](#). *ACM Comput. Surv.*, 57(6).
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and 1 others. 2024. [The Llama 3 herd of models](#). *CoRR*, abs/2407.21783.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2024. [GPTScore: Evaluate as you desire](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6556–6576, Mexico City, Mexico. Association for Computational Linguistics.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. [The language model evaluation harness](#).
- Nathan Habib, Clémentine Fourrier, Hynek Kydlíček, Thomas Wolf, and Lewis Tunstall. 2023. [Lighteval: A lightweight framework for LLM evaluation](#).
- Steve Han, Gilberto Titericz Junior, Tom Balough, and Wenfei Zhou. 2025. [Judge’s verdict: A comprehensive analysis of LLM judge capability through human agreement](#). *CoRR*, abs/2510.09738.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ACM Trans. Inf. Syst.*, 43(2).
- Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, and 1 others. 2025. [Gemma 3 technical report](#). *CoRR*, abs/2503.19786.
- Praveen K. Kanithi, Clément Christophe, Marco AF Pimentel, Tathagata Raha, Nada Saadi, Hamza Javed, Svetlana Maslenskova, Nasir Hayat, Ronnie Rajan, and Shadab Khan. 2024. [MEDIC: Towards a comprehensive framework for evaluating llms in clinical applications](#). *CoRR*, abs/2409.07314.
- Seungone Kim, Jamin Shin, Yejin Choi, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024a. [Prometheus: Inducing fine-grained evaluation capability in language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024b. [Prometheus 2: An open source language model specialized in evaluating other language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4334–4353, Miami, Florida, USA. Association for Computational Linguistics.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). In *Proceedings of the 29th Symposium on Operating Systems Principles, SOSP ’23*, page 611–626, New York, NY, USA. Association for Computing Machinery.
- 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.
- Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using GPT-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.

- Xiaoliang Luo, Akilles Rechartd, Guangzhi Sun, Kevin K. Nejad, Felipe Yáñez, Bati Yilmaz, Kangjoo Lee, Alexandra O. Cohen, Valentina Borghesani, Anton Pashkov, Daniele Marinazzo, Jonathan Nicholas, Alessandro Salatiello, Ilia Sucholutsky, Pasquale Minervini, Sepehr Razavi, Roberta Rocca, Elkhan Yusifov, Tereza Okalova, and 20 others. 2025. [Large language models surpass human experts in predicting neuroscience results](#). *Nature Human Behaviour*, 9(2):305–315.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-Refine: Iterative refinement with self-feedback](#). 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*.
- Daniel McDuff, Mike Schaeckermann, Tao Tu, Anil Palepu, Amy Wang, Jake Garrison, Karan Singhal, Yash Sharma, Shekoofeh Azizi, Kavita Kulkarni, Le Hou, Yong Cheng, Yun Liu, S. Sara Mahdavi, Sushant Prakash, Anupam Pathak, Christopher Semturs, Shwetak Patel, Dale R. Webster, and 9 others. 2025. [Towards accurate differential diagnosis with large language models](#). *Nature*.
- Bhuvanashree Murugadoss, Christian Pölitz, Ian Drosos, Vu Le, Nick McKenna, Carina Suzana Negreanu, Chris Parnin, and Advait Sarkar. 2025. [Evaluating the evaluator: Measuring LLMs’ adherence to task evaluation instructions](#). In *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA, pages 19589–19597*. AAAI Press.
- Jasmine Chiat Ling Ong, Shelley Yin-Hsi Chang, Wasswa William, Atul J. Butte, Nigam H. Shah, Lita Sui Tjien Chew, Nan Liu, Finale Doshi-Velez, Wei Lu, Julian Savulescu, and Daniel Shu Wei Ting. 2024. [Ethical and regulatory challenges of large language models in medicine](#). *The Lancet Digital Health*, 6(6):e428–e432.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [BLEU: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA, pages 311–318*. ACL.
- Sijun Tan, Siyuan Zhuang, Kyle Montgomery, William Yuan Tang, Alejandro Cuadron, Chenguang Wang, Raluca A. Popa, and Ion Stoica. 2025. [Judgebench: A benchmark for evaluating LLM-based judges](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- UK AI Security Institute. 2024. [Inspect AI: Framework for Large Language Model Evaluations](#).
- Dave Van Veen, Cara Van Uden, Louis Blanke-meier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerová, Nidhi Rohatgi, Poonam Hosamani, William Collins, Neera Ahuja, Curtis P. Langlotz, Jason Hom, Sergios Gatidis, John Pauly, and Akshay S. Chaudhari. 2024. [Adapted large language models can outperform medical experts in clinical text summarization](#). *Nature Medicine*, 30(4):1134–1142.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. [Asking and answering questions to evaluate the factual consistency of summaries](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5008–5020, Online. Association for Computational Linguistics.
- WHO. 2023. [WHO calls for safe and ethical AI for health](#).
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, and 1 others. 2025. [Qwen3 technical report](#). *CoRR*, abs/2505.09388.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. [React: Synergizing reasoning and acting in language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Wen-wai Yim, Yujuan Fu, Asma Ben Abacha, Neal Snider, Thomas Lin, and Meliha Yetisgen. 2023. [Acibench: a novel ambient clinical intelligence dataset for benchmarking automatic visit note generation](#). *Scientific Data*, 10(1):586.
- Zhuohao Yu, Chang Gao, Wenjin Yao, Yidong Wang, Zhengran Zeng, Wei Ye, Jindong Wang, Yue Zhang, and Shikun Zhang. 2024. [FreeEval: A modular framework for trustworthy and efficient evaluation of large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 1–13, Miami, Florida, USA. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BERTScore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- 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. 2023. [Judging LLM-as-a-judge with MT-Bench and Chatbot Arena](#). 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*.

## **A Used prompts**

### **A.1 Note Generation Prompt**

The prompt used for the ACI-Bench note generation, adapted from (Kanithi et al., 2024), is shown in Listing 1.

### **A.2 Perturbation Prompt 1**

The prompt used for rephrasing the output notes without changing their meaning is shown in Listing 2.

### **A.3 Perturbation Prompt 2**

The prompt used for introducing minor content changes is shown in Listing 3.

### **A.4 Perturbation Prompt 3**

The prompt used for significantly changing the meaning of the generated notes is given in Listing 4.

```

Your task is to generate a clinical note based on a conversation between a doctor
  \ and a patient. Use the following format for the clinical note:

1. **CHIEF COMPLAINT**: [Brief description of the main reason for the visit]
2. **HISTORY OF PRESENT ILLNESS**: [Summary of the patient's current health status
  \ and any changes since the last visit]
3. **REVIEW OF SYSTEMS**: [List of symptoms reported by the patient]
4. **PHYSICAL EXAMINATION**: [Findings from the physical examination]
5. **RESULTS**: [Relevant test results]
6. **ASSESSMENT AND PLAN**: [Doctor's assessment and plan for treatment or further
  \ testing]

**Conversation:**
{prompt}

**Note:**

```

Listing 1: Note generation prompt

```

Your task is to generate a clinically plausible variation of the provided clinical
  \ note.

You should maintain the original note's structure and formatting, but modify its
  \ content according to the specified types of perturbation below. Try to
  \ maintain internal consistency and general medical plausibility when applying
  \ any changes.

**Perturbation Instructions**
Apply the following types of perturbations:
- Rephrase sentences while preserving the exact medical meaning. You may use
  \ synonyms, vary sentence structure, or change sentence length, but all
  \ clinical facts and measurements must remain unchanged.
- Slightly alter the writing style, such as using different terminology or
  \ presenting findings differently, while ensuring the factual content remains
  \ identical.

Respond only with the perturbed clinical note, do not include any commentary,
  \ reasoning or explanation.

**Original Clinical Note**
{prompt}

**Perturbed Clinical Note**

```

Listing 2: Perturbation prompt 1

Your task is to generate a clinically plausible variation of the provided clinical  
 \ note.

You should maintain the original note's structure and formatting, but modify its  
 \ content according to the specified types of perturbation below. Try to  
 \ maintain internal consistency and general medical plausibility when applying  
 \ any changes.

**\*\*Perturbation Instructions\*\***

Apply the following types of perturbations:

- Make small changes to test results and quantitative measurements, ensuring they  
 \ remain clinically plausible and consistent with the original context.
- Introduce minor modifications to the patient's reported symptoms, making sure  
 \ they are still consistent with the assessment, diagnosis, and treatment plan  
 \ (e.g., adding or substituting symptoms that commonly co-occur).
- Slightly adjust the patient's clinical history, ensuring consistency with the  
 \ assessment, diagnosis, and treatment plan.
- Make minor modifications to the treatment plan, but ensure it remains appropriate  
 \ for the assessment and diagnosis.

Respond only with the perturbed clinical note, do not include any commentary,  
 \ reasoning or explanation.

**\*\*Original Clinical Note\*\***

{prompt}

**\*\*Perturbed Clinical Note\*\***

Listing 3: Perturbation prompt 2

Your task is to generate a clinically plausible variation of the provided clinical  
 \ note.

You should maintain the original note's structure and formatting, but modify its  
 \ content according to the specified types of perturbation below. Try to  
 \ maintain internal consistency and general medical plausibility when applying  
 \ any changes.

**\*\*Perturbation Instructions\*\***

Apply the following types of perturbations:

- Significantly alter test results and quantitative measurements, in a way that may  
 \ change the clinical interpretation or implications of the note.
- Make substantial changes to the patient's reported symptoms, potentially  
 \ affecting the clinical interpretation of the note.
- Make substantial changes to the patient's clinical history, potentially affecting  
 \ the clinical interpretation.
- Significantly modify the treatment plan, such that it may lead to a different  
 \ clinical outcome than the original plan.

Respond only with the perturbed clinical note, do not include any commentary,  
 \ reasoning or explanation.

**\*\*Original Clinical Note\*\***

{prompt}

**\*\*Perturbed Clinical Note\*\***

Listing 4: Perturbation prompt 3

# AI for Climate Finance: Agentic Retrieval and Multi-Step Reasoning for Early Warning System Investments

Ario Saeid Vaghefi<sup>1,2\*</sup>, Aymane Hachcham<sup>1\*</sup>  
Veronica Grasso<sup>2</sup>, Nakiete Msemo<sup>2</sup>, Chiara Colesanti Senni<sup>1</sup>  
Markus Leippold<sup>1,3</sup>

<sup>1</sup>University of Zurich    <sup>2</sup>WMO    <sup>3</sup>Swiss Finance Institute (SFI)

{saeid.vaghefi, aymane.hachcham, chiara.colesantisenni, markus.leippold}@df.uzh.ch

{svaghefi, vgrasso, nmsemo}@wmo.int

## Abstract

Tracking financial investments in climate adaptation is complex and expertise-intensive, particularly for Early Warning Systems (EWS), where multilateral development bank (MDB) and fund reports lack standardized financial reporting and appear as heterogeneous PDFs with complex tables and inconsistent layouts.

We introduce an agent-based Retrieval-Augmented Generation (RAG) system that uses hybrid retrieval and internal chain-of-thought (CoT) reasoning to extract relevant financial data, classify EWS investments, and allocate budgets with grounding evidence spans. While these components are individually established, our contribution is their integration into a domain-specific workflow tailored to heterogeneous MDB reports and numerically grounded EWS budget allocation. On a manually annotated CREWS Fund corpus, our system outperforms four alternatives (zero-shot classifier, few-shot “zero rule” classifier, fine-tuned transformer-based classifier, and few-shot CoT+ICL classifier) on multi-label classification and budget allocation, achieving 87% accuracy, 89% precision, and 83% recall. We further benchmark against the Gemini 2.5 Flash AI Assistant on an expert-annotated MDB evidence set co-curated with the World Meteorological Organization (WMO), enabling a comparative analysis of glass-box agents versus black-box assistants in transparency and performance. The system is publicly deployed and accessible at <https://ews-front.vercel.app/> (see Appendix B for demonstration details and Appendix E for dataset statistics and splits).<sup>1</sup>

## 1 Introduction

Recent advances in Large Language Models (LLMs) have improved automated analysis of fi-

\*Equal Contributions.

<sup>1</sup>We will open-source all code, LLM generations, and human annotations to support further work on AI-assisted climate finance.

ancial documents, yet tracking investments in Early Warning Systems (EWS) remains difficult because Multilateral Development Bank (MDB) and climate-fund reports lack standardized labels, structures, and terminology for EWS-related spending. EWS are central to disaster risk reduction and climate resilience, with the UN’s Early Warnings for All (EW4All) initiative targeting universal coverage by 2027, but current reporting practices leave EWS financial flows opaque and hinder efficient allocation of climate-finance resources. We frame this problem as a combined multi-label classification and budget allocation task: the system assigns each text or table snippet to one or more EWS pillars and extracts pillar-level budget allocations with grounding evidence spans, producing a structured JSON output over the five EW4All pillars (see Appendix D for definitions and examples).

**Contributions.** We present the *EW4All Financial Tracking AI-Assistant*, a glass-box, agent-based Retrieval-Augmented Generation (RAG) system that parses heterogeneous MDB project documents, classifies EWS investments across pillars, and returns numerically grounded, evidence-linked budget allocations. Our key contributions are:

1. A novel agent-based RAG pipeline integrating iterative sub-query generation, hybrid semantic-lexical retrieval, self-validation guardrails, and schema-aware consolidation for climate finance document analysis.
2. A publicly deployed system accessible at <https://ews-front.vercel.app/>, enabling practitioners to analyze MDB documents in real-time.
3. A comprehensive evaluation on a manually annotated CREWS-Fund corpus where our pipeline achieves 87% accuracy, 89% precision, and 83% recall, outperforming four strong baselines.
4. A comparative study against black-box assis-

tants (Gemini 2.5 Flash, OpenAI Assistants) on an expert-annotated MDB evidence set curated with WMO.

5. Open-source release of expert-annotated corpus, benchmark dataset, and all prompt designs to catalyze future research.

**Implications.** By turning unstructured MDB reports into structured, evidence-based EWS investment profiles, our system improves climate-finance transparency, accountability, and decision support for MDBs, funds, and technical partners. The combination of RAG and agentic reasoning yields traceable outputs that support portfolio screening, gap analysis across EWS pillars, and monitoring of progress toward EW4All objectives, and offers a transferable blueprint for AI-assisted analysis of climate adaptation and development finance.

## 2 Related Work

RAG augments LLMs with external retrieval for knowledge-intensive tasks (Lewis et al., 2020), but static pipelines limit adaptability. Recent *agentic RAG* introduces iterative retrieval and decision-making, improving factuality and multi-step reasoning (Xi et al., 2023; Yao et al., 2023; Guo et al., 2024), while multi-agent variants specialize roles for tasks such as code generation and verification and enhance explainability and human–AI collaboration (Guo et al., 2024; Liu et al., 2024). In parallel, in-context learning (ICL) enables few-shot generalization without fine-tuning (Brown et al., 2020); retrieval-based ICL and reward models optimize demonstration selection (Wang et al., 2024). Chain-of-thought (CoT) prompting improves stepwise reasoning (Wei et al., 2022; Kojima et al., 2022), with self-consistency and active example selection further boosting complex question-answering performance (Wang et al., 2023; Diao et al., 2024).

## 3 System Overview

MDB project documents possess highly heterogeneous layouts—mixed narrative text, nested tables, multi-column formats, and scattered financial evidence—making conventional retrieval pipelines insufficient for accurate budget extraction. To address this, we developed the *EW4All Financial Tracking AI-Assistant*, an agent-based RAG system that integrates hybrid retrieval with hierarchical reasoning.

As illustrated in Figure 1, our pipeline consists of five integrated stages. First, we process doc-

uments using the Docling parser to extract raw text and structural elements, followed by context-augmented chunking where each chunk is enriched with a document-level summary to reduce semantic ambiguity. Second, we employ hybrid retrieval that fuses dense vector search (using OpenAI embeddings) and sparse lexical search (BM25F) via Reciprocal Rank Fusion (RRF) to capture both semantic meaning and exact financial figures.

Third, an LLM Agent orchestrates the reasoning process by generating iterative sub-queries and validating retrieved evidence against coverage thresholds. If the retrieved context is insufficient, the agent triggers a self-healing loop to re-query the database. Finally, the system executes schema-aware consolidation, mapping the extracted evidence to the five EWS pillars and allocating budgets with explicit evidence grounding.

The full technical implementation, including embedding construction, hybrid rank fusion equations, and the agent’s control flow, is detailed in Appendix A.

## 4 System Demonstration

The EW4All Financial Tracking AI-Assistant is publicly deployed and accessible at <https://ews-front.vercel.app/>. This section describes the system’s user interface, key features, and demonstration scenarios.

### 4.1 Interface Overview

The web-based interface provides an intuitive workflow for climate finance analysts:

1. **Document Upload:** Users can upload MDB project documents in PDF format. The system accepts documents from various MDBs and climate funds, handling heterogeneous layouts automatically.
2. **Real-Time Processing:** Upon upload, the system displays processing progress with intermediate reasoning steps, allowing users to observe the agent’s sub-query generation and retrieval operations.
3. **Interactive Results:** The classification results are presented with:
  - Pillar-wise budget allocations with confidence scores
  - Clickable evidence spans that highlight source passages in the original PDF
  - A visual distribution chart showing budget allocation across EWS pillars

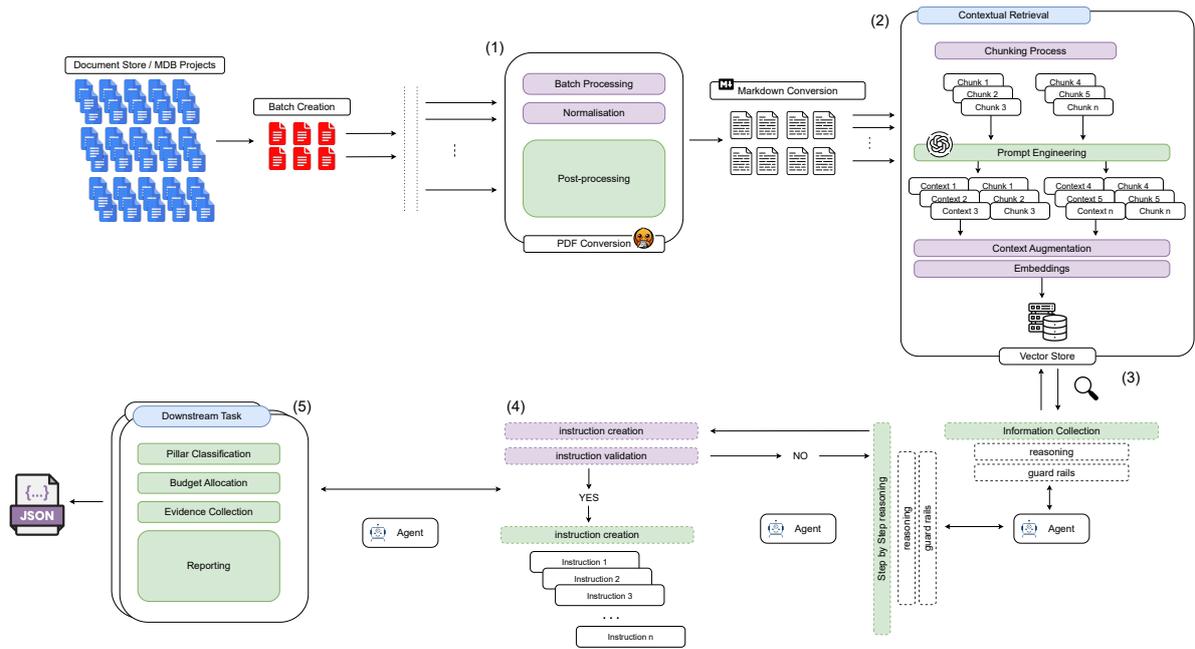


Figure 1: AI-driven financial tracking pipeline for EWS investments. The workflow comprises five stages: (1) PDF conversion using Docling parser, (2) context-augmented chunking with document-level summaries, (3) hybrid retrieval combining dense vectors and BM25F lexical search, (4) iterative agent-based sub-query generation with self-validation loops, and (5) downstream task execution including pillar classification and budget allocation with evidence grounding.

4. **Export Functionality:** Users can export structured JSON outputs for downstream analysis, integration with existing financial tracking systems, or portfolio-level aggregation.

## 4.2 Key Features

**Evidence Traceability.** Each budget allocation is linked to specific text or table fragments from the source document. Users can click on any pillar allocation to view the supporting evidence, enabling expert validation and audit trails.

**Confidence Indicators.** The system provides confidence scores for each classification decision, flagging low-confidence predictions that may require human review. This supports the expert-in-the-loop workflow essential for financial accountability.

**Multi-Document Analysis.** Users can upload multiple documents for batch processing, enabling portfolio-level analysis across projects or funds. Aggregated views show cross-document patterns and potential gaps in EWS coverage.

**Comparison Mode.** For research and validation purposes, the interface offers a comparison view showing outputs from different system configurations (e.g., agent-based vs. baseline methods) side-by-side.

## 4.3 Technical Architecture

The deployed system comprises:

- **Frontend:** React-based web application hosted on Vercel, providing responsive UI and real-time updates via WebSocket connections.
- **Backend:** FastAPI server handling document processing, agent orchestration, and API endpoints.
- **Vector Database:** Weaviate instance for efficient hybrid retrieval over indexed document embeddings.
- **LLM Integration:** OpenAI API for reasoning and classification tasks, with configurable model selection.

The system processes a typical 50-page MDB project document in under 3 minutes, compared to 2–3 hours for manual expert analysis—a reduction of over 98% in processing time.

## 5 Results

**Evaluation Protocol:** Unless stated otherwise, we evaluate on held-out test sets split at the document level, so that no project report contributes evidence to more than one split. For the pillar-level experiment, the classifier and baselines are trained and tuned on a subset of CREWS-Fund documents and

evaluated on disjoint projects; for the MDB Evidence Set, the evidence segments in the test split are drawn exclusively from held-out documents. This prevents label leakage across splits and ensures that performance is measured on previously unseen reports (see Appendix E for split statistics and sampling details).

### 5.1 Pillar-Level Budget Classification

We frame the CREWS-Fund experiment as a joint pillar-classification and budget-allocation task. For each document  $d$  we observe a gold *pillar budget vector*

$$\mathbf{b}_d = (b_{d,1}, \dots, b_{d,5}) \in \mathbb{R}_{\geq 0}^5, \quad \sum_{p=1}^5 b_{d,p} = B_d^{\text{tot}}, \quad (1)$$

where  $b_{d,p}$  is the amount assigned to EWS pillar  $p$  and  $B_d^{\text{tot}}$  is the total EWS envelope. Gold budgets satisfy the conservation constraint by construction; model predictions  $\hat{b}_{d,p}$  are not renormalized and may over- or under-allocate across pillars.

Binary pillar indicators are defined as

$$y_{d,p} = \llbracket b_{d,p} > 0 \rrbracket \in \{0, 1\}, \quad (2)$$

with Iverson bracket  $\llbracket \cdot \rrbracket$ . The model outputs  $\hat{\mathbf{b}}_d$  and  $\hat{y}_{d,p} = \llbracket \hat{b}_{d,p} > 0 \rrbracket$ . Aggregation of chunk-level outputs into document-level  $\hat{\mathbf{b}}_d$  and  $\hat{y}_{d,p}$  is defined in Appendix G.5.

A prediction for pillar  $p$  in document  $d$  is a true positive (TP) only if

- (a) **Label correct:**  $y_{d,p} = 1$  and  $\hat{y}_{d,p} = 1$ ;
- (b) **Budget within tolerance:**

$$|\hat{b}_{d,p} - b_{d,p}| \leq 0.05 B_d^{\text{tot}}, \quad (3)$$

i.e., a  $\pm 5\%$  window around the gold pillar amount.

If the model predicts a pillar where  $y_{d,p} = 0$  or violates (3), we count a false positive (FP); if  $y_{d,p} = 1$  but the pillar is missing or outside the tolerance, we count a false negative (FN). We compute Accuracy, Precision, Recall, and  $F_1$  over all  $(d, p)$  pairs and report macro-averaged scores across pillars.

Using a manually annotated CREWS-Fund corpus (Appendix E), we benchmark four baselines (Zero-Shot, Few-Shot, Transformer, Few-Shot-CoT) against our *Glass-Box Agentic* pipeline. As shown in Table 1, the agent attains **0.87** accuracy, **0.89** precision, and **0.83** recall, an 8–14 pp improvement over the strongest baseline.

The evaluation set reflects the imbalanced distribution of pillars and budget magnitudes that analysts encounter in practice, rather than an artificially

| Method       | Accuracy    | Precision   | Recall      |
|--------------|-------------|-------------|-------------|
| Zero-Shot    | 0.41        | 0.40        | 0.61        |
| Few-Shot     | 0.42        | 0.45        | 0.64        |
| Transformer  | 0.41        | 0.64        | 0.32        |
| Few-Shot-CoT | 0.51        | 0.63        | 0.71        |
| Agent        | <b>0.87</b> | <b>0.89</b> | <b>0.83</b> |

Table 1: Evaluation metrics for budget distribution across the EWS Pillars. The agent-based approach significantly outperforms all baselines on all metrics.

balanced benchmark.

These figures show that the agent not only identifies the correct set of pillars but also assigns budget to them with tight numeric fidelity, providing a solid reference line for the broader Glass-Box vs. Black-Box study in § 5.2.

### 5.2 Glass-Box vs. Black-Box Study (MDB Evidence Set)

To assess whether transparency still pays off in an end-to-end setting, we construct an expert-annotated MDB evidence set co-curated with the World Meteorological Organization (WMO) (see Appendix E). Each segment is labeled with its EWS pillar, the corresponding budget amount, the evidence–pillar linkage, and the document’s total EWS budget, allowing us to jointly evaluate retrieval, reasoning traceability, and numerical fidelity.

We compare three systems: our **Glass-Box Agent** (Section A.3), **Gemini 2.5 Flash**, and **OpenAI Assistants**, both used as black-box assistants that process the same PDFs with a single, carefully engineered prompt. For Gemini 2.5 Flash and OpenAI Assistants, the prompt specifies the role (EWS financial analyst), task (EWS funding allocation), the EWS taxonomy, stepwise analysis instructions, and a JSON output schema. Full details and examples are provided in Appendix I.

Performance is evaluated along five facets, using the same aggregation and tolerance rules as in Section 5.1 and Appendix G.5:

- **Evidence extraction:** recall, precision,  $F_1$ , and Recall@5 for recovering gold evidence segments.
- **Pillar-label assignment:** multi-label Accuracy, Precision, Recall, and  $F_1$  over the five EWS pillars.
- **Amount distribution across pillars:** comparison of  $\hat{b}_{d,p}$  to  $b_{d,p}$  with the same  $\pm 5\%$

tolerance band as in Eq. (3), yielding macro-averaged Accuracy, Precision, Recall, and  $F_1$  over  $(d, p)$  decisions.

- **Evidence-to-label mapping:** correctness of linking retrieved segments to the right pillar, again via TP/FP/FN counts.
- **Total EWS amount prediction:** for each document

$$\hat{B}_d^{\text{tot}} = \sum_{p=1}^5 \hat{b}_{d,p},$$

and conservation accuracy

$$\text{acc}_{\text{tot}}(d) = 1 - \frac{|\hat{B}_d^{\text{tot}} - B_d^{\text{tot}}|}{B_d^{\text{tot}}},$$

together with absolute and percentage errors.

The main analysis uses the full, naturally imbalanced evidence set; results on a balanced subsample with equal support per pillar are reported in Appendix I, Table 5.

### 5.3 Interpretation of the Benchmark

**Total-amount accuracy (Fig. 2, left).** The Glass-Box Agent attains the highest median total-amount accuracy ( $\hat{x} \approx 0.78$ ) with a narrow interquartile range, indicating stable performance across heterogeneous layouts. Gemini 2.5 Flash and OpenAI Assistants trail behind (median  $\approx 0.73$  and  $\approx 0.68$ ) and exhibit heavier tails, reflecting more frequent large conservation errors.

**Amount-per-pillar performance (Fig. 2, right).** When accuracy is measured at the pillar level, the Agent captures nearly half of the aggregate macro- $F_1$  mass (48.7%), while Gemini 2.5 Flash accounts for 36.1% and OpenAI Assistants 15.2%. This mirrors Table 1: schema-aware, transparent reasoning yields the most faithful pillar-level budget breakdowns.

**Evidence-extraction robustness (Fig. 3).** Across most MDB projects, the Agent attains the highest evidence-extraction  $F_1$ , with Gemini 2.5 Flash and OpenAI trailing. The main exception are where budgets are not in explicit tables but diffused through narrative text (grey bands), where Gemini 2.5 Flash slightly outperforms the Agent, reflecting a residual advantage of large black-box models on heavily prose-centric layouts; this is consistent with the balanced-subsample scores in Table 5 (Appendix I), where the Agent still leads overall.

The benchmark indicates that *glass-box, mod-*

*ular retrieval-reasoning pipelines* dominate on structured and semi-structured financial disclosures, while black-box assistants narrow the gap only when numeric cues are deeply embedded in free-form text. Closing this gap is a key direction for future work, for example by enriching the Agent’s retrieval module with paragraph-level numerical parsing.

### 5.4 Ablation Study

To quantify the contribution of individual components of the Glass-Box Agent, we conduct an ablation study on the MDB evidence development set. We systematically remove (i) context augmentation, (ii) hybrid dense+BM25F retrieval, (iii) the top- $k$  setting used for retrieval, and (iv) the agent’s self-healing loop, measuring the impact on evidence extraction  $F_1$ , Recall@5, pillar-level macro- $F_1$ , and total-amount accuracy.

Removing context augmentation yields a noticeable drop in retrieval quality and downstream budget fidelity, confirming that short document-level summaries help disambiguate otherwise similar chunks. Switching from hybrid to dense-only retrieval primarily hurts Recall@5 and evidence  $F_1$ , indicating that exact lexical matching is still crucial for capturing scattered numerical clues. Varying  $k$  shows that  $k = 5$  provides the best trade-off between coverage and noise. Finally, disabling the self-healing loop (single-pass retrieval with no re-querying) reduces both evidence  $F_1$  and total-amount accuracy, particularly on documents with fragmented tables, underscoring the importance of iterative verification. Full ablation results are reported in Appendix J.

## 6 Bias Awareness and Mitigation

We acknowledge that our system may exhibit biases inherited from training data, particularly when classifying novel financial structures or terminology not well-represented in the CREWS Fund corpus. To address these concerns, we implement several mitigation strategies:

**Confidence-Based Human Review.** The system outputs confidence scores for each pillar classification. Predictions with confidence below a configurable threshold (default: 0.7) are automatically flagged for human expert review, ensuring that uncertain classifications do not propagate unchecked.

**Pillar-Level Uncertainty Quantification.** Beyond point predictions, we compute uncertainty esti-

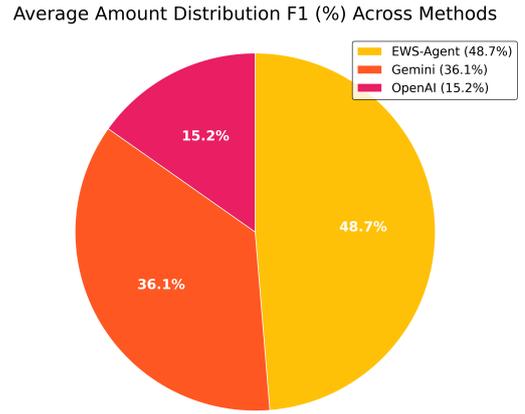
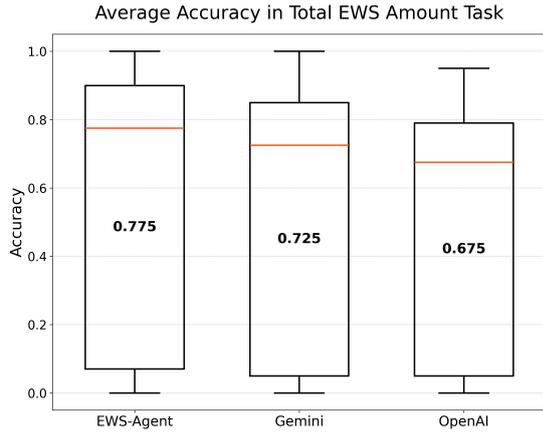


Figure 2: Left: box-plot of the average accuracy on the Total EWS Amount task, evaluated on the expert-annotated test set for each system. The Glass-Box Agent shows the highest median accuracy with the narrowest inter-quartile range, indicating consistent performance. Right: Pie chart showing each system’s share of the overall macro-averaged F1 score on the same test set (EWS-Agent 48.7%, Gemini 36.1%, OpenAI 15.2%).

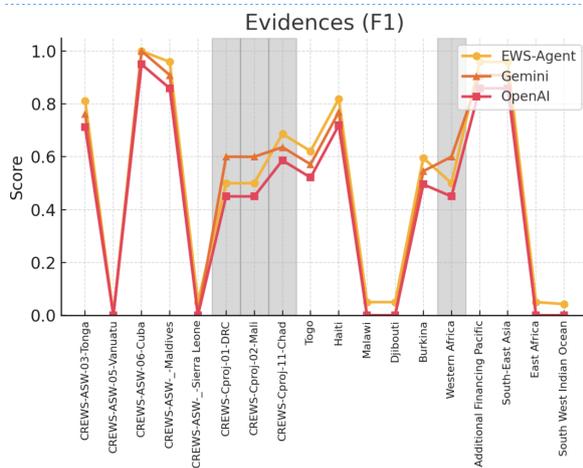


Figure 3: Per-document F1 for **evidence extraction**. Grey bands highlight projects in which budget figures are dispersed across narrative sections rather than formatted tables. The Glass-Box Agent (yellow) consistently outperforms black-box alternatives except in heavily prose-centric documents.

mates using Monte Carlo dropout during inference. High-uncertainty predictions are highlighted in the user interface, enabling analysts to prioritize review efforts.

**Expert-in-the-Loop Validation.** The deployed system (Section 4) supports an expert validation workflow where domain specialists can review, validate, and optionally override system predictions. All corrections are logged, creating a feedback loop for continuous model improvement.

**Cross-Fund Generalization Testing.** While our primary evaluation uses CREWS Fund documents, we conducted preliminary tests on documents from

other climate funds (Green Climate Fund, Adaptation Fund) to assess generalization. Performance degradation on out-of-distribution documents is documented in Appendix K, and we recommend re-calibration when applying the system to new funding sources.

**Terminology Coverage Analysis.** We maintain a glossary of EWS-related terms encountered during training and flag documents containing significant out-of-vocabulary terminology. This alerts users when the system encounters potentially novel financial structures.

## 7 Conclusion

We presented the EW4All Financial Tracking AI-Assistant, an agent-based RAG system designed to extract EWS investments from heterogeneous MDB reports. Achieving 87% accuracy on a manually annotated corpus, our approach significantly outperforms traditional NLP baselines and provides a transparent alternative to black-box assistants. The system is publicly deployed and currently supports early adopters in uncovering uncatalogued investments and accelerating reporting; we refer readers to Appendix C for full deployment details, real-world impact case studies, and future research directions.

## References

- Christoph Auer, Maksym Lysak, Ahmed Nassar, Michele Dolfi, Nikolaos Livathinos, Panos Vagenas, Cesar Berrospi Ramis, Matteo Omenetti, Fabian Lindlbauer, Kasper Dinkla, Lokesh Mishra, Yusik Kim, Shubham Gupta, Rafael Teixeira de Lima, Valery Weber, Lucas Morin, Ingmar Meijer, Viktor Kuropiatnyk, and Peter W. J. Staar. 2024. [Docling technical report](#).
- 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.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. [Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation](#).
- Gordon V. Cormack, Charles L.A. Clarke, and Stephan 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*, pages 758–759. ACM.
- Shizhe Diao, Pengcheng Wang, Yong Lin, Rui Pan, Xiang Liu, and Tong Zhang. 2024. [Active prompting with chain-of-thought for large language models](#).
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. [Large language model based multi-agents: A survey of progress and challenges](#).
- Michael Günther, Isabelle Mohr, Daniel James Williams, Bo Wang, and Han Xiao. 2024. [Late chunking: Contextual chunk embeddings using long-context embedding models](#).
- 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.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling Lou. 2024. [Large language model-based agents for software engineering: A survey](#).
- Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. [Nomic embed: Training a reproducible long context text embedder](#).
- Gianluca Pescaroli, Sarah Dryhurst, and Georgios Marios Karagiannis. 2025. Bridging gaps in research and practice for early warning systems: new datasets for public response. *Frontiers in Communication*, 10:1451800.
- Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends in Information Retrieval*, 3(4):333–389.
- Andrew C Tupper and Carina J Fearnley. 2023. Mind the gaps in disaster early-warning systems—and fix them. *Nature*, 623:479.
- Omar Velazquez, Gianluca Pescaroli, Gemma Cremen, and Carmine Galasso. 2020. A review of the technical and socio-organizational components of earthquake early warning systems. *Frontiers in Earth Science*, 8:533498.
- Jie Wang, Alexandros Karatzoglou, Ioannis Arapakis, and Joemon M Jose. 2024. Reinforcement learning-based recommender systems with large language models for state reward and action modeling. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 375–385.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#).
- 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.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. 2023. [The rise and potential of large language model based agents: A survey](#).
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.

## A Detailed System Methodology

This appendix provides the technical details of the implementation used in the main paper’s System Overview.

## A.1 Embedding Construction and Indexing

Effective downstream reasoning over MDB PDFs requires an embedding index that respects heterogeneous layouts and scattered evidence. We therefore use a five-stage pipeline: document parsing, chunking, context augmentation, embedding generation, and vector storage.

First, we extract raw text and structural elements from each document  $d$  with the Docling converter (Auer et al., 2024):

$$T_d = \text{DoclingParser}(d), \quad (4)$$

where  $T_d$  denotes all extracted textual elements (narrative segments, tables, and other layout blocks). We then partition  $T_d$  into disjoint chunk sets

$$\mathcal{C} = \mathcal{C}_{\text{struct}} \cup \mathcal{C}_{\text{text}}, \quad (5)$$

where  $\mathcal{C}_{\text{struct}}$  contains tables and structured components (e.g., headers, multi-column regions) and  $\mathcal{C}_{\text{text}}$  contains narrative passages and other non-tabular blocks. This separation preserves structural boundaries and avoids flattening tables or merging unrelated segments, which would degrade embedding quality and retrieval.

To situate each chunk in its document context and reduce semantic ambiguity (Günther et al., 2024), we generate a short summary for each  $c \in \mathcal{C}$  by prompting an LLM with  $P_{\text{ctx}}(c, T_d)$ :

$$\text{ctx}(c) = \text{LLM}(P_{\text{ctx}}(c, T_d)), \quad (6)$$

and form the augmented chunk

$$c' = c \oplus \text{ctx}(c). \quad (7)$$

All augmented chunks  $c'$  are encoded in a single latent space:

$$e_{\text{tt}}(c') = f_{\text{tt}}(c') \in \mathbb{R}^{d_{\text{tt}}}, \quad (8)$$

where  $f_{\text{tt}}$  is a joint text–structure encoder. We empirically compared bge-m3 (Chen et al., 2024), nomic-embed-text:v1.5 (Nussbaum et al., 2024), and OpenAI’s text-embedding-3-small, and selected text-embedding-3-small based on Recall@5, nDCG@5, and MRR@5 on the MDB evidence set (Table 4, Appendix F).

Embeddings are indexed in a Weaviate environment with separate NamedVector configurations for text and structured layouts, enabling efficient hybrid (semantic + lexical) search. Each embedding  $e(c')$  is stored with metadata:

$$\text{VDB\_store}(e(c'), \text{meta}(c')). \quad (9)$$

At inference time, for a file ID  $f$  and query  $q$ , we

retrieve the top-5 relevant chunks:

$$\mathcal{R}(f) = \text{VDB\_query}(q, f), \quad |\mathcal{R}(f)| = 5. \quad (10)$$

Further details on embedding model selection, Weaviate configuration, and chunk metadata are provided in Appendix F.

## A.2 Hybrid Retrieval via Rank Fusion

In addition to the above procedure, we employ a hybrid search strategy that combines dense vector search with BM25F-based keyword search (Robertson and Zaragoza, 2009) to leverage both semantic similarity and exact lexical matching. Let  $\mathcal{R}_v(q, f)$  denote the set of candidate chunks retrieved via dense vector search, and let  $\mathcal{R}_k(q, f)$  denote the candidate chunks obtained via BM25F keyword search. To fuse these two retrieval sets, we use Reciprocal Rank Fusion (RRF) (Cormack et al., 2009). For each candidate chunk  $c \in \mathcal{R}_v(q, f) \cup \mathcal{R}_k(q, f)$ , we compute an RRF score as:

$$\text{RRF}(c) = \sum_{i \in \{v, k\}} \frac{1}{\text{rank}_i(c) + K}, \quad (11)$$

where  $\text{rank}_i(c)$  is the rank of  $c$  in retrieval system  $i$  (with lower ranks corresponding to higher relevance) and  $K$  is a smoothing constant (typically set to 60).

In cases where candidates from different retrieval systems share the same RRF score (e.g., when top-3 candidates from each method have no overlap and their 3rd-ranked chunks yield identical scores), we apply a secondary sort by dense vector similarity score to break ties deterministically.

The final set of retrieved chunks is then given by selecting the top five candidates according to their RRF scores:

$$\mathcal{R}(f) = \text{Top5}\left(\mathcal{R}_v(q, f) \cup \mathcal{R}_k(q, f), \text{RRF}(c)\right). \quad (12)$$

This hybrid method harnesses the semantic sensitivity of dense vector retrieval alongside the precise lexical matching of BM25F.

## A.3 Classification and Budget Allocation

For each retrieved chunk  $c' \in \mathcal{R}(f)$ , we predict an EWS pillar label vector  $y$  (over the five pillars) and an associated budget  $B$ . We compare four baselines that differ in how they obtain  $y$  and  $B$ , plus our agent-based approach; implementation details are provided in Appendix G.

## Zero-Shot and Few-Shot Classification

In the zero-shot and few-shot baselines, we construct a prompt  $P_{\text{Class+Budget}}(c')$  that includes the augmented chunk (and, in the few-shot case, a small set of labeled examples). The LLM directly outputs both labels and budget:

$$\{y, B\} = \text{LLM}(P_{\text{Class+Budget}}(c')). \quad (13)$$

This method leverages the pre-trained knowledge of the LLM, with few-shot prompting guiding its responses.

## Fine-Tuned Transformer-Based Classifier

As a classical NLP baseline, we fine-tune a BERT-base encoder  $M_{\text{ft}}$  as a multi-label classifier on the labeled chunks  $\{(c'_i, y_i)\}_{i=1}^N$  (see Appendix G.2 for details on the architecture). The model outputs a 5-dimensional sigmoid layer and yields pillar predictions  $y = M_{\text{ft}}(c')$ . Budgets are then inferred by a separate LLM call:

$$B = \text{LLM}(P_{\text{Budget}}(c', y)). \quad (14)$$

Chunk-level  $\{y, B\}$  tuples are later aggregated to document-level budgets as described in Appendix H, and conservation is evaluated only at aggregation time via the document-level metrics in Section 5.

## Few-Shot CoT Classification

This approach employs a three-step Chain-of-Thought (CoT) strategy, resulting in a tuple  $\{y, B\}$ . First, structured-layout (e.g., tables) chunks are optionally reformatted into clean markdown:  $c'' = \text{LLM}(P_{\text{reformat}}(c'))$ , otherwise, we set  $c'' = c'$ . Second, we classify the (reformatted) chunk:  $y = \text{LLM}(P_{\text{Class}}(c''))$ . Third, we allocate the budget conditioned on both content and labels:  $B = \text{LLM}(P_{\text{Budget}}(c'', y))$ . This CoT-style factorization encourages more explicit reasoning over table structure and pillar definitions; full prompts and examples are in Appendix G.3.

## Agent-Based Approach

Our agent-based method replaces fixed prompts with an LLM agent that plans, retrieves, and validates before emitting  $\{y, B\}$ . Given a document  $f$ , the agent executes the following steps:

1. **Planning:** Generate a set of sub-tasks  $I = \{i_1, \dots, i_k\}$  and retrieval queries  $Q = \{q_1, \dots, q_\ell\}$ .
2. **Retrieval:** Issue vector-database queries  $\text{VDB\_query}(q, f)$  for each relevant sub-task.

3. **Self-validation:** Check coverage sufficiency; re-query when thresholds are unmet:

$$c'_{i_j \text{ final}} = \begin{cases} \text{VDB\_query}(q_{i_j}^{\text{new}}, f), & \text{if } c'_{i_j} \text{ insufficient,} \\ c'_{i_j}, & \text{otherwise.} \end{cases} \quad (15)$$

4. **Consolidation:** Aggregate intermediate results into a schema-aligned JSON output  $\{y, B\}$  per chunk and document.

The full agent loop, instruction format, and guardrails are detailed in Appendix G.4.

## B System Demonstration Details

### B.1 Accessing the System

The EW4All Financial Tracking AI-Assistant is publicly accessible at:

<https://ews-front.vercel.app/>

The system requires no installation and runs entirely in the browser. Users can create accounts to save analysis history and export results.

### B.2 System Requirements

- Modern web browser (Chrome, Firefox, Safari, Edge)
- PDF documents up to 100 pages
- Stable internet connection for API calls

### B.3 API Access

For programmatic access and integration with existing workflows, we provide a REST API. Documentation is available at <https://ews-front.vercel.app/api/docs>. The API supports:

- Document upload and processing
- Batch analysis of multiple documents
- Retrieval of structured JSON outputs
- Webhook notifications for async processing

### B.4 Sample Outputs

Figure 4 shows a sample analysis report generated by the system, illustrating the structured output format with pillar allocations, evidence links, and confidence scores.

## C Real-World Deployment and Impact

Since our agent-based RAG pipeline was deployed in March 2024, early adopters in the EWS community have realized significant benefits:

- **Uncovering hidden investments.** The World Meteorological Organization (WMO) used the system to scan its MDB portfolio, identifying dozens of EWS allocations that had not

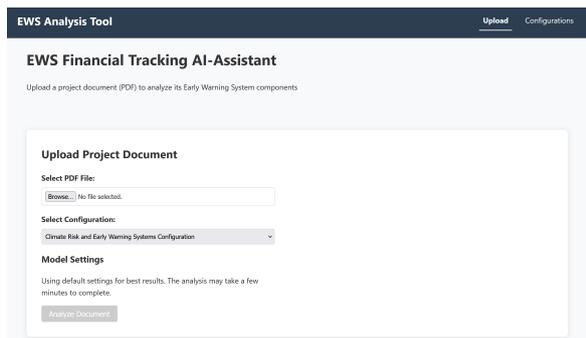


Figure 4: Sample analysis report from the deployed system showing pillar-wise budget allocation with evidence grounding.

previously been catalogued in their internal records.

- **Driving reporting guideline enhancements.** Drawing on classification gaps revealed by our model, the CREWS Fund updated its grant-reporting templates to standardize pillar-level expenditure tagging.
- **Accelerating analysis throughput.** Automated processing reduced the time per project report from 2–3 hours manually to under 3 minutes end-to-end, freeing analysts to focus on higher-value tasks.

These case studies illustrate how transparent, automated extraction not only boosts operational efficiency but also informs better policy and accountability practices at multilateral development banks and climate funds.

## D Early Warning Systems (EWS)

### D.1 Definition and Purpose

Early Warning Systems (EWS) are integrated frameworks designed to detect imminent hazards and alert authorities and communities before disasters strike. In essence, an EWS combines hazard monitoring, risk analysis, communication, and preparedness planning to enable timely, preventive actions. Early warnings are a cornerstone of disaster risk reduction (DRR) – they save lives and reduce economic losses by giving people time to evacuate, protect assets, and secure critical infrastructure<sup>2</sup>. By empowering those at risk to act ahead of a hazard, EWS help build climate resilience: they are

<sup>2</sup>See [https://www.unisdr.org/files/608\\_10340.pdf](https://www.unisdr.org/files/608_10340.pdf).

proven to safeguard lives, livelihoods, and ecosystems amid increasing climate-related threats<sup>3</sup>. In summary, an effective EWS ensures that impending dangers are rapidly identified, warnings reach the impacted population, and appropriate protective measures are taken in advance.

### D.2 EWS Taxonomy

A robust EWS involves several fundamental components that work together seamlessly. The United Nations identify four interrelated pillars necessary for an effective people-centered EWS (Pescaroli et al., 2025). This taxonomy serves as a structured framework to categorize EWS components and activities, facilitating a consistent approach to analyzing early warning systems across various domains. Our approach in this paper is based on these four fundamental pillars of EWS and one cross-pillar, ensuring a comprehensive understanding of risk knowledge, detection, communication, and preparedness.

#### Early Warning System (EWS) Taxonomy Prompt

An Early Warning System (EWS) is an integrated system of hazard monitoring, forecasting, and prediction, disaster risk assessment, communication, and preparedness activities that enables individuals, communities, governments, businesses, and others to take timely action to reduce disaster risks before hazardous events occur.

When analyzing a text, it is essential to determine whether it falls under EWS components and activities, which vary across multiple sectors and require coordination and financing from various actors.

**The taxonomy is based on the Four Pillars of Early Warning Systems and one cross-pillar:**

#### **Pillar 1: Disaster Risk Knowledge and Management (Led by UNDRR)**

This pillar focuses on understanding disaster risks and enhancing the knowledge of communities by collecting and utilizing comprehensive information on hazards, exposure, vulnerability, and capacity.

<sup>3</sup>See <https://www.unep.org/topics/climate-action/climate-transparency/climate-information-and-early-warning-systems>.

**Illustrative examples:**

- Inclusive risk knowledge: Incorporating local, traditional, and scientific risk knowledge.
- Production of risk knowledge: Establishing a systematic recording of disaster loss data.
- Risk-informed planning: Ensuring decision-makers can access and use updated risk information.
- Data rescue: Digitizing and preserving historical disaster data.

**Keywords:** Risk mapping, vulnerability mapping, disaster risk reduction (DRR), climate information.

---

**Pillar 2: Detection, Observation, Monitoring, Analysis, and Forecasting (Led by WMO)**

This pillar enhances the capability to detect and monitor hazards, providing timely and accurate forecasting.

**Illustrative examples:**

- Observing networks enhancement: Strengthening real-time monitoring systems.
- Hazard-specific observations: Improving monitoring of high-impact hazards.
- Impact-based forecasting: Developing quantitative triggers for anticipatory action.

**Keywords:** Forecasting, seasonal predictions, multi-model projections, climate services.

---

**Pillar 3: Warning Dissemination and Communication (Led by ITU)**

Effective communication ensures that early warnings are received by those at risk, enabling them to take timely action.

**Illustrative examples:**

- Multichannel alert systems: Use of SMS, satellite, sirens, and social media.
- Standardized warnings: Implementation of the Common Alerting Protocol (CAP).
- Feedback mechanisms: Enabling community input on warning effectiveness.

**Keywords:** Communication systems, mul-

tichannel dissemination, emergency broadcast systems.

---

**Pillar 4: Preparedness and Response Capabilities (Led by IFRC)**

Timely preparedness and response measures translate early warnings into life-saving actions.

**Illustrative examples:**

- Emergency preparedness planning: Developing anticipatory action frameworks.
- Public awareness campaigns: Educating communities on disaster response.
- Emergency shelters: Construction of cyclone shelters, evacuation centers.

**Keywords:** Preparedness planning, emergency drills, public education on disaster response.

---

**Cross-Pillar: Foundational Elements for Effective EWS**

Cross-cutting elements critical to the sustainability and effectiveness of EWS include governance, inclusion, institutional arrangements, and financial planning.

**Illustrative examples:**

- Governance and institutional frameworks: Defining roles of agencies and stakeholders.
- Financial sustainability: Mobilizing and tracking finance for early warning systems.
- Regulatory support: Developing and enforcing data-sharing legislation.

**Keywords:** Institutional frameworks, governance, financial sustainability, data management.

Each of these components is vital. Only when risk knowledge, monitoring, communication, and preparedness work in unison can an early warning system effectively protect lives and properties. Gaps in any one element (for example, if warnings don't reach the vulnerable, or if communities don't know how to respond) will weaken the whole system. Thus, successful EWS are people-centered and end-to-end, linking high-tech hazard detection with on-the-ground community action.

### D.3 Importance for Climate Finance

EWS are widely recognized as a high-impact, cost-effective investment for climate resilience. By providing advance notice of floods, storms, heatwaves and other climate-related hazards, EWS significantly reduce disaster losses. Studies indicate that every \$1 spent on early warnings can save up to \$10 by preventing damages and losses.<sup>4</sup> For example, just 24 hours' warning of an extreme event can cut ensuing damage by about 30%, and an estimated USD \$800 million investment in early warning infrastructure in developing countries could avert \$3–16 billion in losses every year<sup>5</sup>. These economic benefits underscore why EWS are considered "no-regret" adaptation measures, i.e., they pay for themselves many times over by protecting lives, assets, and development gains.

Given their proven value, EWS have become a priority in climate change adaptation and disaster risk reduction funding. International climate finance mechanisms, such as the Green Climate Fund, Climate Risk and Early Warning Systems (CREWS) Fund, and Adaptation Fund along with development banks, are channeling resources into EWS projects, from modernizing meteorological services and hazard monitoring networks to community training and alert communication systems. Strengthening EWS is also central to global initiatives like the United Nations' Early Warnings for All (EW4All), which calls for expanding early warning coverage to 100% of the global population by 2027. Achieving this goal requires substantial financial support to build new warning systems in climate-vulnerable countries and to maintain and upgrade existing ones. Climate finance is therefore being directed to help develop, implement, and sustain EWS, ensuring that countries can operate these systems (e.g., funding for equipment, data systems, and personnel) over the long term.

In summary, investing in EWS is essential for climate resilience. It not only reduces humanitarian and economic impacts from extreme weather, but also yields high returns on investment. Financial support for EWS, whether through dedicated climate funds, loans and grants, or public budgets, underpins their development and sustainability, making it possible to deploy cutting-edge technology

<sup>4</sup>See <https://wmo.int/news/media-centre/early-warnings-all-advances-new-challenges-emerge>.

<sup>5</sup>See <https://www.unep.org/topics/climate-action/climate-transparency/climate-information-and-early-warning-systems>.

and foster prepared communities. By mitigating the worst effects of climate disasters, EWS help safeguard development progress, which is why they feature prominently in climate adaptation financing and strategies.

### D.4 Current Challenges

Despite their clear benefits, there are several challenges in financing and implementing EWS effectively. Key issues include:

**Data Inconsistencies and Lack of Standardization:** EWS rely on data from multiple sources (weather observations, risk databases, etc.), but often this data is inconsistent, incomplete, or not shared effectively across systems. Differences in how hazards are monitored and reported can lead to gaps or delays in warnings. Likewise, there is a lack of standardization in early warning protocols and data formats between agencies and countries (Velazquez et al., 2020; Pescaroli et al., 2025). Incompatible data systems and inconsistent methodologies (for example, different trigger criteria for warnings or varying risk assessment methods) make it difficult to integrate information. This fragmentation hinders the creation of a "common operating picture" of risk. Data harmonization and common standards (for data collection, forecasting models, and warning communication) are needed to ensure EWS components work together seamlessly.

**Institutional and Cross-Organizational Barriers:** An effective EWS cuts across many organizations: national meteorological services, disaster management agencies, local governments, international partners, and communities. Coordinating these actors remains a challenge. In many cases, efforts are siloed: meteorological offices may issue technical warnings that don't fully reach or engage local authorities or the public. There are gaps in governance, clarity of roles, and inter-agency communication that can weaken the warning chain. Improving EWS often requires overcoming bureaucratic boundaries and fostering cooperation between different sectors (e.g., linking climate scientists with emergency planners). Interoperability issues—i.e., ensuring different organizations' technologies and procedures align—are also a hurdle (Tupper and Fearnley, 2023). As the World Meteorological Organization (WMO) states, connecting all relevant actors (from international agencies down to community groups) and adapting plans to

real-world local conditions is complex<sup>6</sup>. Sustained commitment, clear protocols, and partnerships are required to break down these barriers so that EWS operate as a cohesive, cross-sector system.

**Financing Gaps and Sustainability:** While funding for EWS is rising, it still lags behind what is needed for global coverage and maintenance. Many high-risk developing countries lack the resources to install or upgrade EWS infrastructure (radar, sensors, communication tools) and to train personnel. Fragmented financing is a problem. Support comes from various donors and programs without a unified strategy, leading to potential overlaps in some areas and stark gaps in others. For instance, recent analyses show that a large share of EWS funding is concentrated in a few countries, while Small Island Developing States (SIDS) and Least Developed Countries (LDCs) remain underfunded despite being highly vulnerable<sup>7</sup>. Even when initial capital is provided to set up an EWS, securing long-term funding for operations and maintenance (software updates, staffing, equipment calibration) is difficult. Without sustainable financing, systems can degrade over time. Ensuring financial sustainability, co-financing arrangements, and political commitment is critical so that EWS are not one-off projects but enduring services.

In addition to the above, there are challenges in technological adoption and last-mile delivery: for example, reaching remote or marginalized populations with warnings (issues of language, literacy, and reliable communication channels) and building trust so that people heed warnings. Climate change is also introducing new complexities—hazards are becoming more unpredictable or intense, testing the limits of existing early warning capabilities. Overall, addressing data and standardization issues, improving institutional coordination, and closing funding gaps are priority challenges to fully realize the life-saving potential of EWS.

## D.5 Relevance to This Study

Our work is focused on the financial tracking and classification of investments in climate resilience, and EWS represent a prime example of such investments. Early warning projects often cut across sectors and funding sources—they might include

components of infrastructure, technology, capacity building, and community outreach. Because of this cross-cutting nature, tracking where and how money is spent on EWS can be difficult without a clear classification system. Different organizations may label EWS-related activities in various ways (e.g., "hydromet modernization", "disaster preparedness", "climate services"), leading to inconsistencies in investment data. By establishing a standardized framework to define and categorize EWS investments, the study helps create a "big-picture view" of early warning financing. This enables analysts and policymakers to identify overlaps, gaps, and trends that were previously obscured by fragmented data.

Moreover, improving the classification of EWS funding directly supports broader resilience initiatives. For instance, the newly launched Global Observatory for Early Warning System Investments is already working to tag and track EWS-related expenditures across major financial institutions. Such efforts mirror the goals of this study by highlighting the need for consistent tracking, transparency, and coordination in climate resilience finance. Better classification of investments means stakeholders can pinpoint where resources are going and where additional support is needed to meet global targets like the "Early Warnings for All by 2027" pledge. In short, EWS feature in this study as a critical category of climate resilience investment that must be clearly identified and monitored.

By including EWS in its financial tracking framework, the study provides valuable insights for decision-makers. It helps determine how much funding is allocated to early warnings, from which sources, and for what components (equipment, training, maintenance, etc.). This information is crucial for evidence-based decisions on scaling up EWS: for example, spotting a shortfall in community-level preparedness funding, or recognizing successful investment patterns that could be replicated. Ultimately, linking EWS to the study's financial tracking reinforces the message that climate resilience investments can be better managed when we know their size, scope, and impact area. By classifying EWS expenditures systematically, the study contributes to stronger accountability and strategic planning in building climate resilience, ensuring that early warning systems—and the communities they protect—get the support they urgently need.

<sup>6</sup>See <https://wmo.int/news/media-centre/early-warnings-all-advances-new-challenges-emerge>.

<sup>7</sup>See <https://wmo.int/media/news/tracking-funding-life-saving-early-warning-systems>.

## E Dataset Construction

In this study, we analyze financial information extracted from MDB project PDFs that contain both structured and unstructured data. Unlike conventional benchmark datasets, these documents exhibit high heterogeneity in their formats: some tables are well-structured, while others embed financial figures within free-text paragraphs or disperse them across multiple rows and columns. In many cases, a single numerical value corresponds to several rows or sub-rows within the same column, creating challenges for extraction, alignment, and interpretation.

### E.1 CREWS-Fund Budget Corpus (Pillar-Level)

The pillar-level budget experiment is based on a corpus of 500 CREWS-Fund project reports, with a total of 20,000 expert-annotated segments. We split this corpus at the *document* level into training, validation, and test sets with a 70/20/10 proportion; no project appears in more than one split, preventing cross-document leakage. The label distribution is intentionally imbalanced and mirrors real-world practice: some EWS pillars receive substantially more annotated budget than others, and many projects assign zero budget to certain pillars.

The annotated data, provided by domain experts in CSV format, together with the corresponding PDFs, are included in the supplementary materials of this paper. Each row in the CSV file contains the following nine fields: *Fund*, *Project ID*, *Component*, *Outcome/Expected-Outcome/Objectives*, *Output/Sub-component*, *Activity/Output Indicator*, *Page Number*, *Amount*, and *Label*. The total amount of Early Warning Systems (EWS) funding for a given project is computed as the sum of all *Amount* values associated with that project.

### E.2 Dataset Statistics

Table 2 summarizes the key statistics of our annotated corpus.

| Statistic                     | Value     |
|-------------------------------|-----------|
| Total documents               | 500       |
| Total annotated segments      | 20,000    |
| Training set (documents)      | 350 (70%) |
| Validation set (documents)    | 100 (20%) |
| Test set (documents)          | 50 (10%)  |
| Average segments per document | 40        |
| Average pages per document    | 47        |

Table 2: Dataset statistics for the CREWS-Fund corpus.

### E.3 Pillar Distribution

The distribution of annotations across EWS pillars reflects real-world funding patterns:

| Pillar                           | Segments | Percentage |
|----------------------------------|----------|------------|
| Pillar 1 (Risk Knowledge)        | 3,200    | 16%        |
| Pillar 2 (Detection/Forecasting) | 6,400    | 32%        |
| Pillar 3 (Dissemination)         | 4,000    | 20%        |
| Pillar 4 (Preparedness)          | 4,800    | 24%        |
| Cross-Pillar (Governance)        | 1,600    | 8%         |

Table 3: Distribution of annotated segments across EWS pillars.

### E.4 Data Access and Licensing

The annotated corpus (CSV file and PDFs) consists of financial reports and investment documents sourced from publicly available institutional records, which are intended for public information, research, and transparency purposes. The dataset is used strictly within this intended scope—analyzing financial tracking in climate investments—and adheres to the original access conditions. For all artifacts derived from this corpus, including benchmark datasets and classification models, we explicitly specify their intended use for research and evaluation in automated financial tracking and ensure compliance with relevant ethical research guidelines.

## F Embedding Model Selection

To select the joint text–table encoder  $f_{tt}$ , we constructed a small retrieval benchmark from MDB project documents. For each annotated evidence segment, we issued the corresponding query and measured retrieval quality on a held-out development split. We report standard top- $k$  metrics: Recall@5, nDCG@5, and MRR@5, computed over all queries.

Table 4 summarizes the results for the three candidate encoders. OpenAI’s text-embedding-3-small achieves the best performance across all metrics, and we therefore use it as  $f_{tt}$  in all experiments.

### F.1 Weaviate Configuration

We deploy a Weaviate cluster with:

- Two NamedVectors per object: one for  $e_{tt}(c')$  (semantic) and one for a bag-of-words representation (lexical).
- HNSW indexing for the semantic vector, with tuned `efConstruction` and `M` parameters.

| Encoder                       | R@5         | nDCG@5      | MRR@5       |
|-------------------------------|-------------|-------------|-------------|
| bge-m3                        | 0.72        | 0.68        | 0.65        |
| nomic-embed                   | 0.70        | 0.66        | 0.63        |
| OpenAI text-embedding-3-small | <b>0.78</b> | <b>0.73</b> | <b>0.70</b> |

Table 4: Retrieval performance of candidate embedding models on the MDB evidence development set. OpenAI’s text-embedding-3-small achieves the best overall ranking quality and is used in our deployed system.

- BM25 configuration for lexical search, used in parallel with vector retrieval.

Hybrid scores are formed by a weighted combination of semantic and lexical similarity; weights were chosen on a small dev set to maximize Recall@5.

## F.2 Chunk Metadata

The metadata  $\text{meta}(c')$  stored with each embedding (Eq. 9) includes:

- Document identifier  $f$  and page number,
- Chunk type (structured vs. text) and original layout coordinates,
- Section title and table caption (when available).

These fields are used for filtering (e.g., table-only retrieval) and for reconstructing human-readable evidence views in the UI.

## G Extended Methods: Classification and Budget Allocation

This appendix expands Section A.3, providing full details for each baseline (Zero-Shot / Few-Shot, Fine-Tuned Transformer + LLM, Few-Shot CoT) and for the agent-based system, including prompts, architectures, and training choices that were omitted from the main text for brevity.

### G.1 Zero-Shot and Few-Shot Baselines

**Prompt structure.** For each retrieved chunk  $c' \in \mathcal{R}(f)$ , we construct a prompt  $P_{\text{Class+Budget}}(c')$  with three components:

1. A short description of the task and desired JSON output format for  $\{y, B\}$ , where  $y$  is a 5-dimensional multi-label pillar vector and  $B$  is a numeric budget allocation (possibly zero).
2. A concise description of the five EWS pillars, summarized from Appendix D, including 1–2 example activities per pillar.
3. The augmented chunk  $c'$  (text or table fragment), enriched with basic metadata (docu-

ment ID, section title, and page number when available).

**Zero-shot variant.** In the zero-shot setting, the prompt contains *no* labeled examples: the model relies solely on the pillar descriptions and output schema. The LLM is asked to directly output:

$$\{y, B\} = \text{LLM}(P_{\text{Class+Budget}}(c')), \quad (16)$$

where  $y \in \{0, 1\}^5$  (one bit per pillar) and  $B \in \mathbb{R}_{\geq 0}$ . We enforce the JSON structure with a system-level constraint and discard malformed generations (re-prompting once with an additional format hint).

**Few-shot variant.** The few-shot variant extends the zero-shot prompt with a small set of  $K$  labeled examples  $\{(c^{(k)}, y^{(k)}, B^{(k)})\}_{k=1}^K$ , inserted before the test chunk. Each example includes:

- A short snippet (text or table row/segment) containing a clear EWS signal,
- The gold multi-label vector  $y^{(k)}$ ,
- A corresponding budget value  $B^{(k)}$  (or 0 if the snippet does not carry a numeric allocation).

We use  $K \in \{3, 5\}$  depending on context length; the examples are chosen to cover all five pillars and a mix of single- and multi-label cases. The few-shot prompt still calls a *single* LLM completion:

$$\{y, B\} = \text{LLM}(P_{\text{Class+Budget}}(c')). \quad (17)$$

**Post-processing.** We parse the JSON, map textual pillar names back to indices ("P1"–"P5"), and clip negative budgets to zero. If the model returns a range (e.g., "USD 0.2–0.3M"), we take the midpoint and convert to a single numeric value in the corpus currency (USD) using the same conversion rules as the annotations.

### G.2 Fine-Tuned Transformer + LLM Budget

**Model architecture.** We fine-tune a BERT-base encoder  $M_{\text{ft}}$  on labeled chunks  $\{(c'_i, y_i)\}_{i=1}^N$ , where  $y_i \in \{0, 1\}^5$ :

- 12 Transformer layers, hidden size 768, 12 self-attention heads,
- WordPiece/BPE tokenizer with a 30k–50k subword vocabulary,
- Input sequences truncated or padded to 512 subword tokens,
- A 5-dimensional sigmoid output layer:

$$\hat{y} = \sigma(W h_{[\text{CLS}]} + b),$$

where  $h_{[\text{CLS}]}$  is the final-layer representation of the [CLS] token.

**Training objective.** We treat pillar prediction as multi-label classification with a class-weighted binary cross-entropy loss:

$$\mathcal{L} = - \sum_{j=1}^5 w_j (y_j \log \hat{y}_j + (1 - y_j) \log(1 - \hat{y}_j)), \quad (18)$$

where weights  $w_j$  are inversely proportional to pillar frequency in the training split to mitigate label imbalance.

**Optimization.** We use AdamW with linear learning-rate warmup and decay. A small grid search over learning rate ( $\{1e-5, 2e-5, 3e-5\}$ ) and batch size ( $\{8, 16\}$ ) is performed, selecting the configuration with the best macro-F<sub>1</sub> on the development set. Training runs for up to 10 epochs with early stopping based on dev macro-F<sub>1</sub>.

**Thresholding and calibration.** We select a global sigmoid threshold  $\tau$  by maximizing macro-F<sub>1</sub> on the dev set, then apply it to obtain binary labels:

$$y_j = \mathbb{I}[\hat{y}_j \geq \tau].$$

**Budget allocation prompt.** Given the predicted pillar vector  $y$  and original chunk  $c'$ , we invoke a separate LLM call with prompt  $P_{\text{Budget}}(c', y)$ . The prompt:

1. Reminds the model of the five pillars and provides the predicted subset (e.g., “This chunk is tagged as pillars 2 and 4”),
2. Asks the model to extract the budget amount associated with the EWS-relevant parts of  $c'$  (if any),
3. Requests a single numeric value in USD and a short textual justification.

The LLM returns a JSON payload with the numeric budget and explanation; we retain only the numeric field for evaluation.

**Aggregation.** Chunk-level tuples  $\{y, B\}$  are aggregated into document-level labels and budgets following the rules in Appendix G.5.

### G.3 Few-Shot CoT Baseline

**Reformatting step.** To reduce noise from irregular table layouts, we optionally reformat table-like chunks using a prompt  $P_{\text{reformat}}(c')$ :

$$c'' = \text{LLM}(P_{\text{reformat}}(c')). \quad (19)$$

The prompt asks the model to preserve all numeric entries and column headers, outputting a clean

markdown table. For non-table chunks, we set  $c'' = c'$ .

**Pillar classification.** We then classify the (possibly reformatted) chunk with a dedicated classification prompt:

$$y = \text{LLM}(P_{\text{Class}}(c'')), \quad (20)$$

which:

- Re-states the five pillar definitions more explicitly than in the zero-shot/few-shot baseline,
- Contains a small number of in-context examples where the model first explains which parts of the text support each pillar and then outputs the final label vector.

The model is instructed to think step by step but only return the final JSON in the answer.

**Budget allocation.** Finally, we allocate a budget conditioned on both content and labels:

$$B = \text{LLM}(P_{\text{Budget}}(c'', y)). \quad (21)$$

Compared to the simple zero-shot/few-shot baseline, the CoT prompt explicitly asks the model to reason about which lines or cells in the chunk correspond to EWS-related funding, and then to aggregate them into a single amount. The output again consists of a numeric field and a short natural-language rationale.

### G.4 Agent-Based System

**Instruction schema.** The agent operates over a set of high-level instructions  $I = \{i_1, \dots, i_k\}$ , where each instruction has:

- type (e.g., FIND\_PILLARS, EXTRACT\_BUDGETS, CHECK\_CONSERVATION),
- inputs (references to document  $f$ , chunk IDs, pillar IDs),
- outputs (e.g., list of evidence spans, numeric amounts).

The agent is primed with examples of instruction lists for small documents to illustrate the desired planning behavior.

**Planning and query generation.** Given a document  $f$ , the agent first generates a compact plan:

$$I, Q = \text{LLM}(P_{\text{Plan}}(f_{\text{metadata}})),$$

where  $Q = \{q_1, \dots, q_\ell\}$  is a set of retrieval queries. Each instruction  $i_j$  may be associated with a specific query  $q_{i_j}$  (e.g., “find all chunks related to pillar 2 budgets”).

**Retrieval and self-validation.** For instructions requiring external evidence, the agent issues vector database calls:

$$c'_{i_j} = \text{VDB\_query}(q_{i_j}, f). \quad (22)$$

We then apply a self-validation step where the agent inspects the retrieved chunks and decides whether coverage is sufficient:

$$c'_{i_j \text{ final}} = \begin{cases} \text{VDB\_query}(q_{i_j}^{\text{new}}, f), & \text{if } c'_{i_j} \text{ insufficient} \\ c'_{i_j}, & \text{otherwise.} \end{cases} \quad (23)$$

Coverage criteria are expressed in natural language in the prompt (e.g., “at least one budget line per pillar mentioned in the document”).

**Intermediate results.** For each instruction  $i_j$ , the agent produces an intermediate result  $\text{result}_{i_j}$ , which can contain:

- Candidate pillar labels and evidence spans,
- Candidate budget lines and amounts (possibly per currency),
- Flags indicating uncertainty or missing information.

These results are stored in a scratchpad-like JSON structure.

**Final formatting.** After all instructions are executed, a final formatting prompt  $P_{\text{Format}}(\{\text{result}_I\})$  asks the LLM to consolidate everything into a single, schema-aligned output:

$$\{y, B\} = \text{LLM}(P_{\text{Format}}(\{\text{result}_I\})), \quad (24)$$

where  $y$  is the document-level pillar label vector and  $B$  contains pillar-level budget allocations, each with a list of supporting evidence spans. The JSON schema includes fields for `pillar_id`, `budget_amount`, `currency`, and `evidence_span_ids`.

## G.5 Chunk and Document Aggregation

For the baselines that operate at the chunk level, we aggregate  $\{y, B\}$  tuples into document-level outputs as follows:

- **Labels:** a document is assigned pillar  $j$  if at least one chunk has  $y_j = 1$ ; we also report per-pillar coverage (fraction of chunks tagged with each pillar).
- **Budgets:** for each pillar, we sum chunk-level budgets  $B$  across all chunks that include that pillar; overlapping allocations (chunks with multiple pillars) are split proportionally based on the model’s confidence scores when available, or uniformly otherwise.

- **Conservation:** we compare the sum of all pillar-level budgets against the document’s total EWS budget (when annotated) and report conservation error metrics in Section 5.

## H Document-Level Aggregation of Chunk Predictions

For evaluation, we aggregate chunk-level outputs to obtain document-level budgets and labels. Let  $C_d$  denote the set of chunks associated with document  $d$ , and let  $B_c \in \mathbb{R}_{\geq 0}^5$  be the pillar-wise budget vector predicted for chunk  $c \in C_d$  (missing pillars are treated as zero). The predicted budget for pillar  $p$  in document  $d$  is

$$\hat{b}_{d,p} = \sum_{c \in C_d} B_{c,p},$$

and the corresponding pillar indicator is

$$\hat{y}_{d,p} = \mathbb{I}[\hat{b}_{d,p} > 0].$$

This simple summation scheme is applied uniformly across all methods (Zero-Shot, Few-Shot, Transformer, Few-Shot-CoT, and Agent), ensuring a consistent mapping from chunk-level predictions to document-level budget vectors  $\hat{b}_d$  and label sets  $\hat{y}_{d,p}$ .

## I Black-Box Assistants: Setup and Additional Results

### I.1 Expert-Annotated MDB Evidence Set

The MDB evidence set used in Section 5.2 is derived from a subset of CREWS-related MDB project documents. For each document, domain experts annotated:

- Evidence segments (text or table fragments) that support EWS-relevant budgets,
- The corresponding EWS pillar label(s) for each segment,
- The budget amount assigned to that pillar (normalized to a common currency),
- The document’s total EWS budget.

These annotations define the gold evidence–pillar–amount triples and document-level totals against which all systems are evaluated.

### I.2 Prompt Design for Gemini 2.5 Flash and OpenAI Assistants

Both Gemini 2.5 Flash and OpenAI Assistants are queried in a single end-to-end pass per document, using prompts that follow the same structure:

1. **Role and scope:** The model is instructed to act as a financial analyst specialized in EWS

and MDB climate adaptation projects.

2. **Task description:** Identify EWS-relevant components, assign them to the five EWS pillars, and extract associated budget amounts.
3. **EWS taxonomy:** A concise description of the five pillars (aligned with Appendix D) and examples of typical activities per pillar.
4. **Methodical instructions:** Stepwise guidance on reading the PDF (narrative, tables, footnotes), checking consistency, and avoiding double counting.
5. **Output schema:** A JSON template requiring, for each document, (i) pillar-level labels and budgets, (ii) a list of evidence segments per pillar, and (iii) a total EWS budget estimate.

Gemini 2.5 Flash receives the full PDF via its native file interface; OpenAI Assistants receive the same content as pre-processed text and tables. Minor token-length adaptations aside, both prompts share the same structure and schema.

### I.3 Balanced Evidence Subsample

To test robustness to label imbalance, we construct a balanced subsample of the MDB evidence set with approximately equal support for each EWS pillar. The sampling procedure:

- Identifies the minimum per-pillar evidence count across the full set,
- Uniformly samples that number of evidence segments per pillar,
- Retains only documents that still contain at least one segment for each pillar after sampling.

We recompute all metrics from Section 5.2 on this balanced subset. The qualitative pattern remains unchanged: the Glass-Box Agent maintains the highest macro-averaged scores on evidence extraction, pillar labeling, and pillar-level budget fidelity, with Gemini 2.5 Flash consistently second and OpenAI Assistants third.

### I.4 Metric Computation Details

All metrics in Section 5.2 reuse the definitions from Section 5.1 and Appendix G.5:

- TP/FP/FN counts for evidence segments are computed at the segment level (exact-match or strict overlap, depending on annotation granularity).
- Pillar-level budgets  $\hat{b}_{d,p}$  for black-box systems are obtained by aggregating their own evidence-level outputs using the same summation rule as the Glass-Box Agent.

| System            | $F_{1ev}$ | $F_{1pill}$ | $F_{1bud}$ | med. $acc_{tot}$ |
|-------------------|-----------|-------------|------------|------------------|
| Glass-Box Agent   | 0.83      | 0.82        | 0.80       | 0.79             |
| Gemini 2.5 Flash  | 0.78      | 0.76        | 0.74       | 0.74             |
| OpenAI Assistants | 0.67      | 0.65        | 0.63       | 0.64             |

Table 5: Balanced MDB evidence subsample (approximately equal support per pillar).  $F_{1ev}$  = macro  $F_1$  for evidence extraction;  $F_{1pill}$  = macro  $F_1$  for pillar labels;  $F_{1bud}$  = macro  $F_1$  for pillar-level budget fidelity under the  $\pm 5\%$  tolerance band; med.  $acc_{tot}$  = median total-amount conservation accuracy as in Section 5.2. Values mirror the qualitative pattern in Section 5.3, with the Glass-Box Agent performing best, Gemini 2.5 Flash second, and OpenAI Assistants third.

| Variant              | Evid. $F_1$ | R@5  | Pillar $F_1$ | Acc <sub>tot</sub> |
|----------------------|-------------|------|--------------|--------------------|
| Full Agent           | 0.78        | 0.86 | 0.81         | 0.79               |
| w/o ctx augmentation | 0.73        | 0.82 | 0.77         | 0.74               |
| Dense-only retrieval | 0.69        | 0.78 | 0.74         | 0.71               |
| $k = 3$ (R@3)        | 0.72        | 0.80 | 0.76         | 0.75               |
| $k = 10$ (R@10)      | 0.74        | 0.84 | 0.78         | 0.76               |
| w/o self-healing     | 0.71        | 0.82 | 0.75         | 0.72               |

Table 6: Ablation results for the Glass-Box Agent on the MDB evidence development set. Each variant removes or modifies a single component of the full system.

- Total-amount accuracy  $acc_{tot}(d)$  is computed exactly as in Section 5.2, without renormalization of  $\hat{B}_d^{tot}$ .

Full numeric tables corresponding to Figures 2 and 3 are provided in the supplementary material.

## J Ablation Studies

We evaluate four variants of the Glass-Box Agent on the MDB evidence development set, each obtained by removing or modifying one component at a time. Table 6 reports evidence-extraction  $F_1$ , Recall@5, pillar-level macro- $F_1$ , and document-level total-amount accuracy (as defined in Section 5.2).

## K Cross-Fund Generalization

To assess generalization beyond CREWS Fund documents, we conducted preliminary experiments on a small held-out set of documents from:

- Green Climate Fund (GCF): 15 project documents
- Adaptation Fund (AF): 10 project documents

Without any fine-tuning or re-calibration, we observed the following performance degradation compared to CREWS Fund documents:

The performance drop is primarily attributed to:

1. Different document layouts and table formats

| Fund                | Accuracy | Precision | Recall |
|---------------------|----------|-----------|--------|
| CREWS (in-domain)   | 0.87     | 0.89      | 0.83   |
| GCF (out-of-domain) | 0.72     | 0.75      | 0.69   |
| AF (out-of-domain)  | 0.68     | 0.71      | 0.65   |

Table 7: Generalization performance on out-of-domain climate fund documents.

2. Varying terminology for similar EWS activities
3. Different budget reporting conventions

We recommend re-calibrating the system with a small number of labeled examples from the target fund before deployment. Future work will focus on domain adaptation techniques to improve zero-shot generalization.

## Limitations

While our approach demonstrates significant improvements in automating financial tracking for EWS investments, several limitations remain. First, our system relies on existing financial reports from MDBs, in this case CREWS, which are often heterogeneous and may contain incomplete or ambiguous financial allocations. In cases where funding details are missing or inconsistently reported, even advanced retrieval-augmented generation (RAG) and multi-step reasoning approaches may struggle to provide accurate classifications. Second, the classification system is influenced by the training data used in fine-tuning and prompt engineering. Despite expert annotations, the model may still exhibit biases in investment classification, particularly when encountering novel financial structures or terminology not well-represented in the dataset (see Section 6 for our mitigation strategies). Third, while our agent-based RAG system achieves state-of-the-art performance on structured and unstructured financial data, its generalizability to other climate finance applications outside EWS has not been fully explored (see Appendix K for preliminary cross-fund results). Future work should assess model robustness across different sustainability reporting frameworks and financial instruments. Fourth, our annotated corpora are modest in size compared to large-scale NLP benchmarks, reflecting the difficulty of obtaining expert-labelled MDB financial data. We mitigate this by using real-world, heterogeneous project reports, document-level splits to avoid leakage, and complementary evaluations at both pillar and evidence level, but broader statistical conclusions will re-

quire expanded datasets in future work. Finally, our system assumes that financial tracking can be improved through AI-assisted reasoning; however, its real-world effectiveness depends on institutional adoption, policy integration, and alignment with evolving financial disclosure regulations.

## Ethics Statement

**Human Annotation.** This study relies on annotations provided by domain experts from the WMO, who possess extensive knowledge of Early Warning Systems (EWS). These experts played a pivotal role in the design and conceptualization of the study. Their deep understanding of both the contextual and practical aspects of the collected data ensures the accuracy and relevance of the annotations. The use of expert annotations minimizes the risk of misclassification and enhances the reliability of the model’s outputs.

**Responsible AI Use.** This tool is intended as an assistive system to enhance transparency and efficiency in financial tracking, not as a replacement for human analysts. Expert oversight remains crucial in interpreting financial classifications, addressing edge cases, and ensuring compliance with policy frameworks. By open-sourcing our dataset and model, we encourage responsible use and further validation to refine the system’s applicability in real-world climate finance decision-making.

**Data Privacy and Bias.** This study does not involve any personally identifiable or sensitive financial data. All data used in this research originates from publicly available sources under a Creative Commons license, ensuring compliance with data privacy regulations. While we find no evidence of demographic biases in the dataset, we acknowledge that financial reporting by multilateral development banks (MDBs) may reflect institutional biases in investment classification. Our model operates as a decision-support tool and should not replace human judgment in financial tracking and policy decisions.

**Reproducibility Statement.** To ensure full reproducibility, we will release all PDFs, codes, EWS-taxonomy, and expert-annotated data used in this study. Our approach aligns with best practices in AI transparency and responsible research dissemination. However, we encourage users of this dataset and model to consider ethical implications when applying automated financial tracking systems in real-world decision-making contexts. For

vector database storage and retrieval, we utilized Weaviate, an open-source, scalable vector search engine that efficiently indexes high-dimensional embeddings. Additionally, for reasoning and large language model (LLM) interactions, we integrated OpenAI's API, leveraging its advanced capabilities to process, analyze, and infer patterns from financial document data.

### **Disclaimer**

Opinions expressed in this article are the author's opinions and do not necessarily reflect those of WMO or its Members.

### **Acknowledgements**

This paper has partially received funding from the Swiss National Science Foundation (SNSF) under the project 'How sustainable is sustainable finance? Impact evaluation and automated greenwashing detection' (Grant Agreement No. 100018\_207800).

# RAGVUE: A Diagnostic View for Explainable and Automated Evaluation of Retrieval-Augmented Generation

Keerthana Murugaraj<sup>1</sup>, Salima Lamsiyah<sup>1</sup>, Martin Theobald<sup>1</sup>

<sup>1</sup>University of Luxembourg, Department of Computer Science (DCS),  
Faculty of Science, Technology and Medicine (FSTM), Esch-sur-Alzette, Luxembourg

Correspondence: [keerthana.murugaraj@uni.lu](mailto:keerthana.murugaraj@uni.lu)

## Abstract

Evaluating Retrieval-Augmented Generation (RAG) systems remains a challenging task: existing metrics often collapse heterogeneous behaviors into single scores and provide little insight into whether errors arise from retrieval, reasoning, or grounding. In this paper, we introduce RAGVUE, a diagnostic and explainable framework for automated, reference-free evaluation of RAG pipelines. RAGVUE decomposes RAG behavior into retrieval quality, answer relevance and completeness, strict claim-level faithfulness, and judge calibration. Each metric includes a structured explanation, making the evaluation process transparent. Our framework supports both manual metric selection and fully automated agentic evaluation. It also provides a Python API, CLI, and a local Streamlit interface for interactive usage. In comparative experiments, RAGVUE surfaces fine-grained failures that existing tools such as RAGAS often overlook. Our demonstration showcases the full RAGVUE workflow and illustrates how it can be integrated into research pipelines and practical RAG development. The source code as well as detailed instructions on its usage are publicly available on Github <sup>1</sup>.

## 1 Introduction

Retrieval-Augmented Generation (RAG) combines a pretrained (parametric) language model with an external retriever that supplies relevant documents at inference time (Lewis et al., 2020; Guu et al., 2020). By conditioning generation on retrieved passages, RAG systems effectively tackle knowledge-intensive tasks while making their evidence explicit and easier to maintain than finetuning internal model weights (Lewis et al., 2020; Izacard et al., 2023). This paradigm has rapidly become a default solution for building search assistants, analytical tools, customer-support bots, and domain-specific copilots across high-stakes settings such as finance,

healthcare, and law (Song et al., 2024; Rosenthal et al., 2025). Recent benchmarks further stress-test RAG in multi-hop and multi-turn scenarios (e.g., StrategyQA (Geva et al., 2021), mtRAG (Katsis et al., 2025), CLAPnq (Rosenthal et al., 2025)), underscoring the need for robust and fine-grained evaluation of the full RAG pipeline.

Evaluating RAG is harder than evaluating a standalone language model because errors can arise from retrieval (irrelevant or missing evidence), generation (off-topic, incomplete, or incoherent answers), or grounding (unsupported or contradictory claims despite retrieved context) (Es et al., 2024; Saad-Falcon et al., 2024; Ru et al., 2024). Recent surveys argue that global "end-to-end" scores obscure these components and advocate decomposing the evaluation into retrieval quality, answer quality, and evidence support (Yu et al., 2024; Gan et al., 2025). They also emphasize a crucial distinction between *faithfulness* to retrieved evidence and *factual correctness* with respect to world knowledge: a response may be factually true but unsupported by its citations, or fully grounded in outdated or erroneous evidence (Min et al., 2023; Sorodoc et al., 2025). Temporal drift (Ouyang et al., 2025), unanswerability (Peng et al., 2025), and privacy or policy violations in retrieved content (Zeng et al., 2025; Song et al., 2024) introduce additional evaluation axes that simple accuracy-style metrics cannot capture.

Human annotation and gold references are expensive and brittle under domain shift (Saad-Falcon et al., 2024; Rosenthal et al., 2025), motivating reference-free *LLM-as-a-judge* methods that are widely used in NLG evaluation (Wang et al., 2023; Kocmi and Federmann, 2023; Zheng et al., 2023). Despite progress (e.g., G-Eval (Liu et al., 2023), AutoCalibrate (Liu et al., 2024), SelfCheckGPT (Manakul et al., 2023)), LLM judges remain prompt-sensitive, unstable, and prone to self-preference bias (Panickssery et al., 2024; Schroeder

<sup>1</sup><https://github.com/KeerthanaMurugaraj/RAGVue>

and Wood-Doughty, 2024; Liu et al., 2025). Existing RAG-focused evaluators, including RAGAS (Es et al., 2024), ARES (Saad-Falcon et al., 2024), RAGChecker (Ru et al., 2024), and RAG-Zeval (Li et al., 2025) have expanded coverage. However, two core gaps persist: metrics often collapse heterogeneous behaviors into non-diagnostic scalar scores, and grounding checks remain permissive, missing fine-grained factual errors (Es et al., 2024; Niu et al., 2024; Song et al., 2024).

To address these limitations, we introduce RAGVUE, a reference-free, explainable evaluation framework that offers *diagnostic results* rather than purely numerical assessments. RAGVUE decomposes RAG performance into retrieval quality, answer quality, and factual grounding (Yu et al., 2024; Gan et al., 2025), enforcing *strict faithfulness* by crediting only claim-level evidence explicitly supported in the retrieved context. This yields a more conservative alternative to semantic-inference metrics (Es et al., 2024; Min et al., 2023; Niu et al., 2024; Zhu et al., 2025). RAGVUE additionally introduces a *judge-calibration* score quantifying agreement across LLM evaluators, making stability issues in LLM-as-a-judge setups explicit (Liu et al., 2025; Schroeder and Wood-Doughty, 2024; Panickssery et al., 2024). The framework supports both *manual* metric selection and an *agentic* mode, in which an internal orchestrator automatically chooses and aggregates metrics. Moreover, we provide a Python API, command-line interface (CLI), and a Streamlit-based user interface (UI) for a seamless integration into research workflows. Finally, on a multihop StrategyQA-derived benchmark (Geva et al., 2021), RAGVUE reveals fine-grained failures that approaches based on scalar metrics, such as RAGAS (Es et al., 2024), fail to capture.

## 2 Related Work

**RAG & Evaluation Challenges.** Retrieval-Augmented Generation (RAG) integrates external evidence into LLMs to reduce hallucinations and improve grounding (Lewis et al., 2020; Kocmi and Federmann, 2023; Wang et al., 2023). Early models such as REALM (Guu et al., 2020) and RAG (Lewis et al., 2020) showed strong gains on knowledge-intensive tasks, followed by advances in retrieval and generation (e.g., late-interaction retrievers (Khattab et al., 2021) and few-shot RAG tuning (Izacard et al., 2023)). However, the

pipeline-based nature of RAG introduces unique evaluation challenges. Performance must be assessed across components, including retriever relevance and evidence coverage, generator quality, and grounding faithfulness (Saad-Falcon et al., 2024). Recent surveys argue that end-to-end scores obscure these dimensions and call for separate evaluation of retrieval quality and grounding fidelity (Yu et al., 2024; Gan et al., 2025). Despite having access to documents, RAG models still tend to hallucinate or rely on outdated evidence, motivating benchmarks for temporal drift and attribution (Ouyang et al., 2025). Moreover, faithfulness and factual correctness may diverge: a response can be true but unsupported, or well-grounded yet incorrect (Khattab et al., 2021). We follow this line by separately evaluating retrieval, answer quality, and grounding with strict evidence criteria.

### **RAG Evaluation Frameworks & Benchmarks.**

Recent work has introduced automatic evaluators for RAG. RAGAS (Es et al., 2024) provides reference-free metrics for context relevance, answer coherence, and coarse groundedness via static prompt-based queries. ARES (Saad-Falcon et al., 2024) increases robustness by fine-tuning smaller LMs on human labels, offering explicit relevance and faithfulness scores with confidence estimates. RAGChecker (Ru et al., 2024) adds diagnostic checks for passage usage and claim-level grounding. Benchmarks such as RAGTruth (Niu et al., 2024) and MEMERAG (Cruz Blandón et al., 2025) target hallucinations and multilingual settings, while HoH (Ouyang et al., 2025) and Unanswerability-Eval (Peng et al., 2025) test temporal drift and unanswerable queries. Furthermore, LLM-as-a-judge approaches are widely adopted (Wang et al., 2023; Zheng et al., 2023), including G-Eval (Liu et al., 2023), AutoCalibrate (Liu et al., 2024), and FactScore (Min et al., 2023), but remain prompt-sensitive and biased toward model families (Panickssery et al., 2024; Liu et al., 2025). More reliable methods such as JudgeLM (Liu et al., 2025) and RAG-Zeval (Li et al., 2025) seek stability through consensus and reasoning-based ranking. Building on these insights, RAGVUE uses reference-free LLM judges while addressing key limitations by decomposing scores (retrieval vs. coverage; answer relevance vs. completeness), enforcing strict claim-level faithfulness, and providing fine-grained explanations and stability checks. It also includes an *agentic*

evaluation mode that automatically selects and orchestrates metrics, producing structured summaries ready for debugging and comparison.

### 3 RAGVUE Framework Overview

This section introduces our RAGVUE evaluation framework. We first outline its core metrics, then describe its two operational modes, and conclude with its programmatic and interactive UIs.

#### 3.1 RAGVUE Metrics

We describe seven RAGVUE metrics across three dimensions: (1) retrieval, (2) answer quality, and (3) grounding and stability, with a summary provided in Appendix A (Table 1).

##### 3.1.1 Retrieval Relevance

This metric measures whether the retrieved contexts (C) are actually useful for answering the question (Q). For each context chunk, an LLM judge assigns a relevance score  $r_i$  in  $[0, 1]$  using a predefined range<sup>2</sup>. A chunk ( $c_i$ ) is counted as relevant if its score exceeds this threshold<sup>3</sup>, and the final score is computed as:

$$\text{RetrievalRelevance} = \frac{\#\{c_i \geq \tau\}}{N} \quad (1)$$

where  $N$  is the number of retrieved chunks. This precision-style formulation is simple, cost-efficient, and provides actionable diagnostic insight into retrieval quality. Importantly, it evaluates the usefulness of retrieved documents directly from the question alone, without requiring reference answers.

##### 3.1.2 Retrieval Coverage

This metric measures whether the retrieved contexts (C) collectively contain the evidence needed to answer the question (Q), without requiring reference contexts. We first derive a small set of atomic aspects from the question alone and reuse the same aspects across metrics for consistency. Let  $\mathcal{A}$  denote this set of aspects and let  $\mathcal{R}_{\text{cov}} \subseteq \mathcal{A}$  be the subset of aspects supported by at least one retrieved document. The corresponding score is:

$$\text{RetrievalCoverage} = \frac{|\mathcal{R}_{\text{cov}}|}{|\mathcal{A}|} \quad (2)$$

<sup>2</sup>Default ranges: 1.0–0.9 for direct answer-containing evidence; 0.8–0.7 for highly useful content; 0.6–0.3 for weakly related background; 0.2–0.0 for irrelevant text.

<sup>3</sup>The threshold is set to  $\tau = 0.7$  to include only evidence judged as highly useful.

This recall-style metric indicates whether the retriever has surfaced enough evidence to cover all parts of the question.

##### 3.1.3 Clarity

This metric evaluates the linguistic quality of the generated answer (A), assessing grammar, fluency, logical flow, conciseness, and overall readability. A single LLM call returns a score in  $[0, 1]$  along with a brief explanation and suggested improvements. Short answers are also checked for naturalness and readability. Overall, this metric provides a compact indication of how clearly the answer is written.

##### 3.1.4 Answer Relevance

Answer Relevance measures how well the generated answer (A) aligns with the user’s question (Q) intent. The metric considers only the question and the generated answer, and assigns a score in  $[0, 1]$  based on topical focus and whether the answer meaningfully addresses what the question is asking. It ignores factual correctness and stylistic quality. High scores thus indicate that the answer stays on-topic and captures the main intent, while lower scores reflect partial, generic, or off-topic content. The judge additionally returns short lists of missing or off-topic parts to provide interpretable signals about alignment.

##### 3.1.5 Answer Completeness

Answer Completeness measures how well the answer covers the different aspects implied by the question. Using the same aspect set  $\mathcal{A}$  (from 3.1.2), the metric checks each aspect against the answer and identifies the subset  $\mathcal{A}_{\text{cov}} \subseteq \mathcal{A}$ , i.e., the aspects that the answer explicitly addresses (optionally with short supporting snippets). The final score is computed as:

$$\text{AnswerCompleteness} = \frac{|\mathcal{A}_{\text{cov}}|}{|\mathcal{A}|} \quad (3)$$

This reference-free metric captures how thoroughly the answer resolves the information needs expressed by the question.

##### 3.1.6 Strict Faithfulness

To assess factual grounding, we introduce a *single-pass* faithfulness metric that checks whether the retrieved context supports each factual claim in the generated answer. Using a single LLM call, the evaluator (i) decomposes the answer into minimal atomic claims and (ii) labels each claim as *supported*, *partially hallucinated*, or *fully hallucinated*.

Unlike multi-stage pipelines that require multiple LLM calls and subsequent heuristic aggregation, RAGVUE encodes the verification logic directly within one prompt.

The strictness is intentionally applied to *high-risk factual anchors*, i.e., *key entities* (such as people, locations, organizations) and *temporal expressions* (such as years and dates). We enforce exact agreement on these anchors to capture frequent RAG failure modes such as entity substitution and incorrect temporal details. Claims with missing, unsupported, or contradictory anchors are treated as hallucinated. This design favors conservative factual checking over stylistic flexibility, i.e., paraphrases are acceptable as long as entity and temporal anchors are consistent. The final score is computed as:

$$\text{StrictFaithfulness} = \frac{|\mathcal{C}_{\text{supported}}|}{|\mathcal{C}_{\text{supported}}| + |\mathcal{C}_{\text{hallucinated}}|} \quad (4)$$

where  $\mathcal{C}_{\text{supported}}$  denotes claims fully grounded in the retrieved context and  $\mathcal{C}_{\text{hallucinated}}$  denotes claims marked as partially or fully hallucinated. The resulting score is transparent and directly traceable to claim-level decisions, allowing users to quickly identify which parts of an answer are evidence-backed.

### 3.1.7 Generic Calibration

LLM-based evaluators are sensitive to sampling noise, decoding temperature, and choice of the model, and they may also reflect systematic biases inherited from the judge models. In practice, many RAG evaluation pipelines implicitly assume that a single judge configuration is both stable and trustworthy. RAGVUE does not attempt to remove or solve judge bias. Instead, we make judge reliability observable by introducing a *generic calibration metric* that quantifies agreement across multiple judge configurations and can be applied to any RAGVUE metric. For each evaluation case, we run the same underlying metric under several (model, temperature) configurations and obtain a set of scores  $s_1, \dots, s_k$ . We define calibration agreement as:

$$\text{Calibration} = 1 - \left( \max_i s_i - \min_i s_i \right) \quad (5)$$

which assigns high values when judges behave consistently and low values when their outputs diverge. Importantly, high calibration indicates *stability* across judge settings, but it does not guarantee

correctness, fairness, or absence of systematic bias. Similarly, a low calibration score indicate that the evaluation outcome is brittle and should be treated with caution.

Beyond reporting an aggregate agreement score, RAGVUE surfaces per-judge outputs and explanations, allowing users to inspect which configurations disagree and why. This design supports practical safeguards such as (i) preferring conclusions that are consistent across multiple judges, (ii) marking low-calibration cases for manual review, and (iii) optionally expanding the judge set to include more diverse models to reduce dependence on any single judge’s output. Overall, calibration in RAGVUE increases transparency by indicating whether an evaluation result is consistent across different judge models and temperature settings, or highly sensitive to those choices, without claiming to eliminate inherent judge bias.

## 3.2 Operational Modes & Availability

Our evaluation framework supports two complementary operational modes and user interfaces, enabling flexible integration into research workflows and production pipelines.

### 3.2.1 Operational Modes

RAGVUE provides two modes for users. In **manual mode**, users control which metrics are executed and how results are aggregated, offering full transparency and fine-grained control. In **agentic mode**, an internal orchestration agent fully automates evaluation. The agent selects appropriate retrieval and answer-level metrics based on the presence of context, the availability of an answer, and the user query. It then executes these metrics in a single pass and synthesizes high-level scores, including an overall retrieval score (harmonic mean of relevance and coverage) and an answer-level composite score (weighted blend of strict faithfulness, relevance, completeness, and clarity).

### 3.2.2 Availability

RAGVUE is released under the Apache License 2.0 and can be used through multiple access modes depending on the user’s preference and technical requirements.

**Python API.** RAGVUE can also be used directly as a Python library (Fig. 1). Users import the evaluator, load a JSONL dataset, and run all metrics with a single function call, making this mode ideal for

integration into notebooks, scripts, and automated pipelines.

```

from ragvue import evaluate, load_metrics
items = [
    {"question": "...", "answer": "...",
     "context": [...]}]
metrics = load_metrics().keys()
report = evaluate(items, metrics=list(metrics))
print(report)

```

Figure 1: RAGVUE Python API usage example.

**Python/Command-Line Interface (CLI).** RAGVUE provides a simple command-line interface (ragvue-cli) for terminal-based workflows, as shown in Figure 2. A lightweight Python runner (ragvue-py) is also available, as illustrated in Figure 3. Both interfaces support listing available metrics, running manual evaluations, and executing the agentic mode, enabling fast and scriptable evaluation without writing additional code.

```

# Help
ragvue-cli --help
# List all available metrics
ragvue-cli list-metrics
# Manual evaluation (choose metrics explicitly)
ragvue-cli eval
  --inputs <your_data.jsonl>
  --metrics <metrics>
  --out-base report_manual
  --format "json,md,csv"
# Agentic evaluation (auto-select metrics)
ragvue-cli agentic
  --inputs <your_data.jsonl>
  --out-base report_agentic
  --format "json,md,csv"

```

Figure 2: Example usage of the RAGVUE command-line interface (ragvue-cli).

```

# Help
ragvue-py --help
# Manual Mode Usage
ragvue-py --input <your_data> --metrics <metrics>
  --out-base report_manual --skip-agentic
# Agentic Mode Usage
ragvue-py --input <your_data> --metrics <metrics>
  --agentic-out report_agentic --skip-manual

```

Figure 3: Example usage of the RAGVUEchat Python command-line runner (ragvue-py).

**Local Streamlit Application.** For no-code, interactive usage, we provide a Streamlit-based UI that exposes the same capabilities through an interactive

browser interface. The application is run locally: users clone the repository and start the interface with a standard command such as `streamlit run streamlit_app.py`. Within this local UI, users can upload JSONL files, select operational modes, set/paste API keys for the current session, and generate formatted reports without writing code. This interface is targeted at practitioners who prefer a point-and-click workflow while keeping all data and keys on their own machine. The images of our UI are shown in Appendix E.

## 4 Experiments & Discussion

### 4.1 Dataset

We construct our evaluation dataset based on the multihop StrategyQA (Geva et al., 2021) benchmark. Each item contains a question, the reference yes/no label, the supporting facts, the decomposition steps, and the Wikipedia evidence titles. The supporting facts are cleaned and used as independent context snippets. In the next stage, we generate five answer variants for each question, such as ideal, partial, unclear, off-topic, and hallucinated. These variants capture common RAG failure modes by altering the correctness of the label, the completeness of the explanation, and the relevance or confidence of the response. Each answer is stored with its associated metadata (question ID, reference label, contexts, supporting facts, decomposition, and evidence titles), producing exactly five synthetic examples per question. For this study, we created 100 synthetic ( $Q, C, A$ ) triplets from StrategyQA. The final dataset is exported in two formats: a RAGAS-compatible JSON for metric-based evaluation and a RAGVUE-ready JSONL for interactive inspection.

### 4.2 RAGVUE vs. RAGAS Performance

**Computational Time.** We first measured latency on the 100 queries described in Section 4.1. RAGAS averaged 18.26 seconds per query, while RAGVUE averaged 18.87 seconds. This represents a marginal 3.4% increase in per-item latency, which is negligible. As shown in the boxplot (Appendix B), both systems exhibit nearly identical latency distributions, including similar medians, inter-quartile ranges, and outliers. Importantly, RAGVUE provides richer diagnostics, and this enhanced granularity makes it more actionable for system debugging and improvements, rendering the small computational overhead a worthwhile

trade-off.

**Quantitative Analysis.** We next summarize the behavior of both evaluators using descriptive statistics over the 100 queries (Appendix C). RAGAS’ faithfulness has a mean of 0.52, while answer relevance and response groundedness average at 0.24 and 0.39, respectively, indicating that RAGAS often judges answers as only weakly relevant or weakly grounded. In contrast, RAGVUE reports low average answer completeness (0.12) and moderate answer relevance (0.37), alongside consistently high clarity scores (0.70). Its retrieval metrics center around 0.50 for coverage and 0.42 for relevance, while strict-faithfulness has a mean of 0.40, reflecting a wide range of partially supported and unsupported answers.

A correlation analysis (see Appendix C) shows that the two evaluators align on broad, generation-focused behavior but diverge sharply on retrieval-focused metrics. RAGAS’ faithfulness, answer relevancy, and response groundedness correlate strongly with RAGVUE’s strict faithfulness and answer relevance, indicating comparable sensitivity to high-level answer correctness. However, RAGAS’ retrieval-related metrics, context relevance and response groundedness, correlate only weakly or inconsistently with RAGVUE’s retrieval coverage and retrieval relevance. This reveals that RAGAS often conflates insufficient retrieval with unsupported reasoning, whereas RAGVUE explicitly separates retrieval performance from generation performance.

**Qualitative Analysis.** Our qualitative inspection reveals several systematic failures that RAGAS does not diagnose, and some examples are presented in Appendix D. We find that RAGAS often gives scores that look reasonable but do not explain why an answer fails. When the model gives vague, unsupported, or partially relevant answers, RAGAS may still assign high or mid-range faithfulness scores because it only checks for direct contradictions and does not account for missing multi-hop reasoning, unanswered parts of the question, or unsupported conclusions. RAGVUE, on the other hand, clearly shows what went wrong: it marks claims as unsupported when the evidence does not back them, highlights when the answer ignores key aspects of the question, and indicates whether retrieval fully or only partially matched what was needed. As a result, RAGVUE makes it easy to see whether the error comes from re-

trieval, grounding, or the model’s reasoning, which is not possible with RAGAS. Overall, the qualitative analysis shows that RAGVUE provides clearer and more actionable feedback for diagnosing RAG system failures.

**Discussion.** Our quantitative and qualitative analyses directly reflect the structural limitations of RAGAS. Context Relevance collapses all chunks into a single aggregated score and cannot show which passages are missing or irrelevant. RAGVUE provides per-chunk relevance without requiring reference answers, clearly exposing retrieval strengths and failures. Likewise, RAGAS’ Context Recall requires a reference answer and multi-step alignment, while RAGVUE’s Retrieval Coverage operates directly on the question and retrieved documents, making it usable even when references are unavailable. On the generation side, RAGAS’ Answer Relevancy relies on embedding-based synthetic question generation, capturing only coarse semantic overlap. RAGVUE’s Answer Relevance is intent-aware and identifies missing or off-topic elements, yielding actionable diagnostics about why an answer may be incomplete. Finally, RAGAS’ Response Groundedness assigns coarse labels (0/1/2) based on inferred support in the context, but cannot reveal which question aspects were addressed or missed. RAGVUE’s Answer Completeness evaluates coverage directly from the question’s aspect structure, producing a fine-grained completeness signal. Strict Faithfulness exhibits the clearest contrast: RAGAS uses a two-step pipeline to extract statements from the answer and verify each with a semantic-inference prompt that tolerates semantic drift, whereas RAGVUE decomposes the answer into atomic claims and enforces exact evidence matching for key entities and temporal expressions. This yields a stricter, more deterministic assessment of factual support. Overall, these results show that, while RAGAS provides high-level semantic judgments, RAGVUE delivers finer-grained, retrieval-aware, and more diagnostically meaningful evaluation signals.

## 5 RAGVUE in the Loop

RAGVUE is designed to support iterative RAG development. Developers can modify individual pipeline components (e.g., retrieval settings, chunking, reranking, prompting, or abstention rules) and re-run RAGVUE on a fixed development set to compare diagnostic reports across versions. This

workflow is supported through multiple access interfaces: evaluations can be executed programmatically via the Python API, scripted through the CLI, or inspected interactively using the Streamlit UI.

To make iterations actionable, RAGVUE exposes decomposed evaluation signals with structured explanations. It separates retrieval behavior (*retrieval relevance* and *retrieval coverage*) from answer quality (*answer relevance* and *answer completeness*) and grounding (*strict claim-level faithfulness*). This decomposition is intended to help localize failures to retrieval, generation, or grounding rather than relying on a single aggregate score. For example, retrieval coverage reflects whether the retrieved set contains evidence for the required aspects of a query, retrieval relevance provides per-chunk usefulness signals, and strict faithfulness flags claims that are unsupported under evidence-matching constraints.

## 6 Conclusion

RAGVUE provides an automated, diagnostic, and fully reference-free evaluation framework tailored for explainable assessment of RAG systems. It separates retrieval and generation-level metrics, delivers structured explanations rather than opaque scalar scores, and exposes the underlying causes of model failures. The agentic evaluation mode makes the framework immediately usable with minimal setup, automatically selecting appropriate metrics and producing structured reports that highlight where and why a pipeline breaks. By combining fine-grained metrics with transparent reasoning traces and cross-model calibration for reliability, RAGVUE reveals whether an error stems from retrieval drift, missing evidence, unsupported reasoning, or incomplete answers, problems that traditional metrics such as RAGAS often conflate. Overall, RAGVUE functions not only as an evaluator but as a practical debugging tool for real-world RAG development, which helps users identify weaknesses and iteratively improve their RAG systems.

## 7 Limitations & Future Work

RAGVUE represents our first step toward a transparent and diagnostic evaluation framework for RAG. The current version delivers seven core metrics, two operational modes, and a no-code user interface, but there is still significant room for growth. As RAGVUE relies on LLM-based eval-

uation, careful selection of the underlying judge models is recommended to ensure stable and consistent scoring. We plan to extend RAGVUE with additional metrics, perform retrieval and grounding analysis on complex queries, and provide broader support for different LLM models. The agentic mode will become more adaptive, assisting users by detecting errors and automatically selecting appropriate metrics.

At present, RAGVUE provides diagnostic outputs and structured explanations, but it does not automatically modify the underlying RAG system. A natural extension is tighter integration into semi-automated development loops, where evaluation results guide interventions such as retriever configuration updates, reranker/threshold selection, or abstention/refusal policies when evidence is insufficient. In this setting, RAGVUE’s generic calibration is particularly important: by exposing instability in LLM-based judging across model and temperature configurations, it helps users avoid acting on brittle evaluations. Over time, our goal is to develop RAGVUE into a unified evaluation pipeline with baseline models, system comparisons, and optional multimodal support.

## References

- María Andrea Cruz Blandón, Jayasimha Talur, Bruno Charron, Dong Liu, Saab Mansour, and Marcello Federico. 2025. *MEMERAG: A multilingual end-to-end meta-evaluation benchmark for retrieval augmented generation*. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 22577–22595, Vienna, Austria. Association for Computational Linguistics.
- Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. 2024. *Ragas: Automated evaluation of retrieval augmented generation*. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 150–158.
- Aoran Gan, Hao Yu, Kai Zhang, Qi Liu, Wenyu Yan, Zhenya Huang, Shiwei Tong, and Guoping Hu. 2025. *Retrieval augmented generation evaluation in the era of large language models: A comprehensive survey*. *arXiv preprint arXiv:2504.14891*.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. *Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies*. *Transactions of the Association for Computational Linguistics*, 9:346–361.

- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: few-shot learning with retrieval augmented language models. *J. Mach. Learn. Res.*, 24(1).
- Yannis Katsis, Sara Rosenthal, Kshitij Fadnis, Chulaka Gunasekara, Young-Suk Lee, Lucian Popa, Vraj Shah, Huaiyu Zhu, Danish Contractor, and Marina Danilevsky. 2025. [mtrag: A multi-turn conversational benchmark for evaluating retrieval-augmented generation systems](#). *Transactions of the Association for Computational Linguistics*, 13:784–808.
- Omar Khattab, Christopher Potts, and Matei Zaharia. 2021. Relevance-guided supervision for openqa with colbert. *Transactions of the association for computational linguistics*, 9:929–944.
- Tom Kocmi and Christian Federmann. 2023. [Large language models are state-of-the-art evaluators of translation quality](#). In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.
- Kun Li, Yunxiang Li, Tianhua Zhang, Hongyin Luo, Xixin Wu, James R. Glass, and Helen M. Meng. 2025. [RAG-zeval: Enhancing RAG responses evaluator through end-to-end reasoning and ranking-based reinforcement learning](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 24936–24954, Suzhou, China. Association for Computational Linguistics.
- Shuliang Liu, Xinze Li, Zhenghao Liu, Yukun Yan, Cheng Yang, Zheni Zeng, Zhiyuan Liu, Maosong Sun, and Ge Yu. 2025. [Judge as a judge: Improving the evaluation of retrieval-augmented generation through the judge-consistency of large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 5788–5807, Vienna, Austria. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Yuxuan Liu, Tianchi Yang, Shaohan Huang, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang. 2024. [Calibrating LLM-based evaluator](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2638–2656, Torino, Italia. ELRA and ICCL.
- 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.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [FActScore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.
- Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. [Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10862–10878.
- Jie Ouyang, Tingyue Pan, Mingyue Cheng, Ruiran Yan, Yucong Luo, Jiaying Lin, and Qi Liu. 2025. [HoH: A dynamic benchmark for evaluating the impact of outdated information on retrieval-augmented generation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6036–6063, Vienna, Austria. Association for Computational Linguistics.
- Arjun Panickssery, Samuel R. Bowman, and Shi Feng. 2024. Llm evaluators recognize and favor their own generations. In *Proceedings of the 38th International Conference on Neural Information Processing Systems, NIPS ’24*, Red Hook, NY, USA. Curran Associates Inc.
- Xiangyu Peng, Prafulla Kumar Choubey, Caiming Xiong, and Chien-Sheng Wu. 2025. [Unanswerability evaluation for retrieval augmented generation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8452–8472, Vienna, Austria. Association for Computational Linguistics.
- Sara Rosenthal, Avirup Sil, Radu Florian, and Salim Roukos. 2025. [CLAPnq: Cohesive long-form answers from passages in natural questions for RAG systems](#). *Transactions of the Association for Computational Linguistics*, 13:53–72.

Dongyu Ru, Lin Qiu, Xiangkun Hu, Tianhang Zhang, Peng Shi, Shuaichen Chang, Cheng Jiayang, Cunxiang Wang, Shichao Sun, Huanyu Li, and 1 others. 2024. Ragchecker: A fine-grained framework for diagnosing retrieval-augmented generation. *Advances in Neural Information Processing Systems*, 37:21999–22027.

Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. 2024. **ARES: An automated evaluation framework for retrieval-augmented generation systems**. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 338–354, Mexico City, Mexico. Association for Computational Linguistics.

Kayla Schroeder and Zach Wood-Doughty. 2024. Can you trust llm judgments? reliability of llm-as-a-judge. *arXiv preprint arXiv:2412.12509*.

Maojia Song, Shang Hong Sim, Rishabh Bhardwaj, Hai Leong Chieu, Navonil Majumder, and Soujanya Poria. 2024. Measuring and enhancing trustworthiness of llms in rag through grounded attributions and learning to refuse. *arXiv preprint arXiv:2409.11242*.

Ionut Teodor Sorodoc, Leonardo F. R. Ribeiro, Rexhina Blloshmi, Christopher Davis, and Adrià de Gispert. 2025. **GaRAGE: A benchmark with grounding annotations for RAG evaluation**. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 17030–17049, Vienna, Austria. Association for Computational Linguistics.

Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. **Is ChatGPT a good NLG evaluator? a preliminary study**. In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11, Singapore. Association for Computational Linguistics.

Hao Yu, Aoran Gan, Kai Zhang, Shiwei Tong, Qi Liu, and Zhaofeng Liu. 2024. Evaluation of retrieval-augmented generation: A survey. In *CCF Conference on Big Data*, pages 102–120. Springer.

Zhirui Zeng, Jiamou Liu, Meng-Fen Chiang, Jialing He, and Zijian Zhang. 2025. **S-RAG: A novel audit framework for detecting unauthorized use of personal data in RAG systems**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10375–10385, Vienna, Austria. 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. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.

Kunlun Zhu, Yifan Luo, Dingling Xu, Yukun Yan, Zhenghao Liu, Shi Yu, Ruobing Wang, Shuo Wang, Yishan Li, Nan Zhang, Xu Han, Zhiyuan Liu, and Maosong Sun. 2025. **RAGEval: Scenario specific RAG evaluation dataset generation framework**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8520–8544, Vienna, Austria. Association for Computational Linguistics.

## A Summary of Metrics

Our summary of metrics is available in Table 1

## B Computational Time plot

The computational time box plot is shown in Figure 4

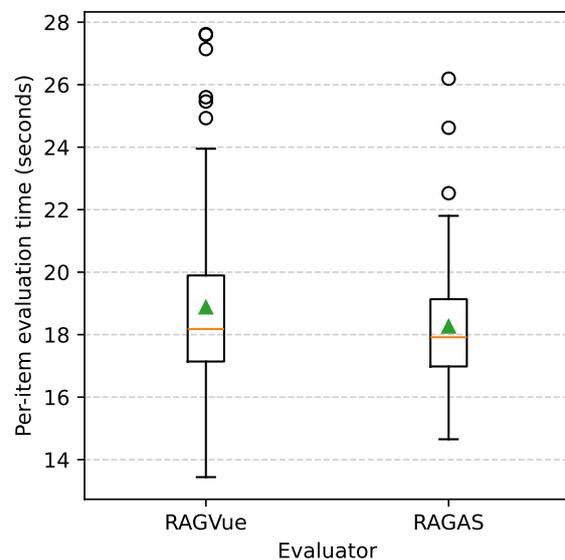


Figure 4: Distribution of per-item evaluation time for RAGVUE and RAGAS on the 100-query benchmark.

## C Quantitative Analysis

The descriptive statistics are shown in Table 2, and the correlation results are provided in Table 3.

## D Qualitative Case Studies - Examples

Table 4 provides item-level analyses for three representative evaluation cases drawn from our dataset. For each item, we report the behaviors of RAGAS and RAGVUE alongside the structured diagnostic signals RAGVUE generates.

**Example 1.** The model gives a vague answer (“No, even though there is no strong supporting evidence”) that does not actually address the comparative question or follow from the retrieved facts. RAGAS still gives it a perfect faithfulness score

| Metric                                   | Inputs  | What it Measures                                                                                                                                                                |
|------------------------------------------|---------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Retrieval Metrics</b>                 |         |                                                                                                                                                                                 |
| <i>Retrieval Relevance</i>               | Q, C    | Evaluates how useful each retrieved chunk is for addressing the information needs of the question, based on per-chunk relevance scoring.                                        |
| <i>Retrieval Coverage</i>                | Q, C    | Assesses whether the retrieved context collectively provides sufficient coverage for all sub-aspects required to answer the question.                                           |
| <b>Answer Metrics</b>                    |         |                                                                                                                                                                                 |
| <i>Answer Relevance</i>                  | Q, A    | Measures how well the answer aligns with the intent and scope of the question, identifying missing, irrelevant, or off-topic content.                                           |
| <i>Answer Completeness</i>               | Q, A    | Determines whether the answer fully addresses all aspects of the question without omissions.                                                                                    |
| <i>Clarity</i>                           | A       | Evaluates the linguistic quality of the answer, including grammar, fluency, logical flow, coherence, and overall readability.                                                   |
| <b>Grounding &amp; Stability Metrics</b> |         |                                                                                                                                                                                 |
| <i>Strict Faithfulness</i>               | A, C    | Evaluates how many factual claims in the answer are directly supported by the retrieved context, enforcing strict evidence alignment (entity accuracy and temporal correctness) |
| <i>Calibration</i>                       | Q, A, C | Examines the stability of metric by measuring variance across different judge configurations (model choice and temperature).                                                    |

Table 1: Summary of the RAGVUE metrics.

| System | Metric                | Mean  | Std   |
|--------|-----------------------|-------|-------|
| RAGAS  | faithfulness          | 0.521 | 0.403 |
|        | answer_relevancy      | 0.240 | 0.307 |
|        | context_relevance     | 0.550 | 0.264 |
|        | response_groundedness | 0.390 | 0.460 |
| RAGVUE | answer_completeness   | 0.121 | 0.225 |
|        | answer_relevance      | 0.372 | 0.255 |
|        | clarity               | 0.698 | 0.100 |
|        | retrieval_coverage    | 0.503 | 0.279 |
|        | retrieval_relevance   | 0.420 | 0.316 |
|        | strict_faithfulness   | 0.400 | 0.492 |

Table 2: Descriptive statistics of RAGAS and RAGVUE metrics on the 100-query benchmark (mean and standard deviation).

(1.0) because it only checks for surface-level consistency and does not verify whether the conclusion is supported across multiple pieces of evidence, missing the needed **multi-hop reasoning**. RAGVUE instead marks the claim as fully unsupported (strict faithfulness = 0.0), shows that none of the key parts of the question are answered (completeness = 0.0), and indicates that retrieval only partly matched what was needed (retrieval coverage = 0.33). Together, these signals clearly reveal an unsupported reasoning error that RAGAS fails to catch.

**Example 2.** The system retrieves the right information, i.e, both RAGAS and RAGVUE show that the context is fully relevant. But the model still gives an incorrect answer (“Yes, even though there is no strong supporting evidence...”), which is not

backed by the retrieved facts. RAGAS gives a mid-range faithfulness score (0.5) without explaining where the mistake comes from. RAGVUE makes this clear: it marks the claim as completely unsupported (strict faithfulness = 0.0), shows that the answer covers none of the required points (completeness = 0.0), and confirms that retrieval was correct. This directly identifies the problem as a *generation error*, not a retrieval issue.

**Example 3.** The model gives a vague answer (“probably true”) even though the retrieved evidence does not actually say which of the two (dog or grey seal) would respond first. RAGAS gives the answer a mid-range faithfulness score ( $\approx 0.67$ ) because it does not directly contradict any single fact, but this does not explain what went wrong. RAGVUE makes the issue clear: it marks the answer as fully unsupported (strict faithfulness = 0.0), shows that the model did not address the key parts of the question (completeness = 0.0), and indicates that only part of the retrieved information was relevant. As a result, RAGVUE pinpoints that the failure also comes from the model’s reasoning, not only from retrieval, which is something RAGAS cannot show.

Across all three examples, RAGVUE provides structured diagnostics that clearly distinguish whether errors stem from retrieval, grounding, or reasoning. In contrast, RAGAS offers only scalar scores, which obscure these distinctions in practice.

Table 3: Spearman correlation between RAGAS metrics and RAGVUE metrics on the 100-query benchmark.

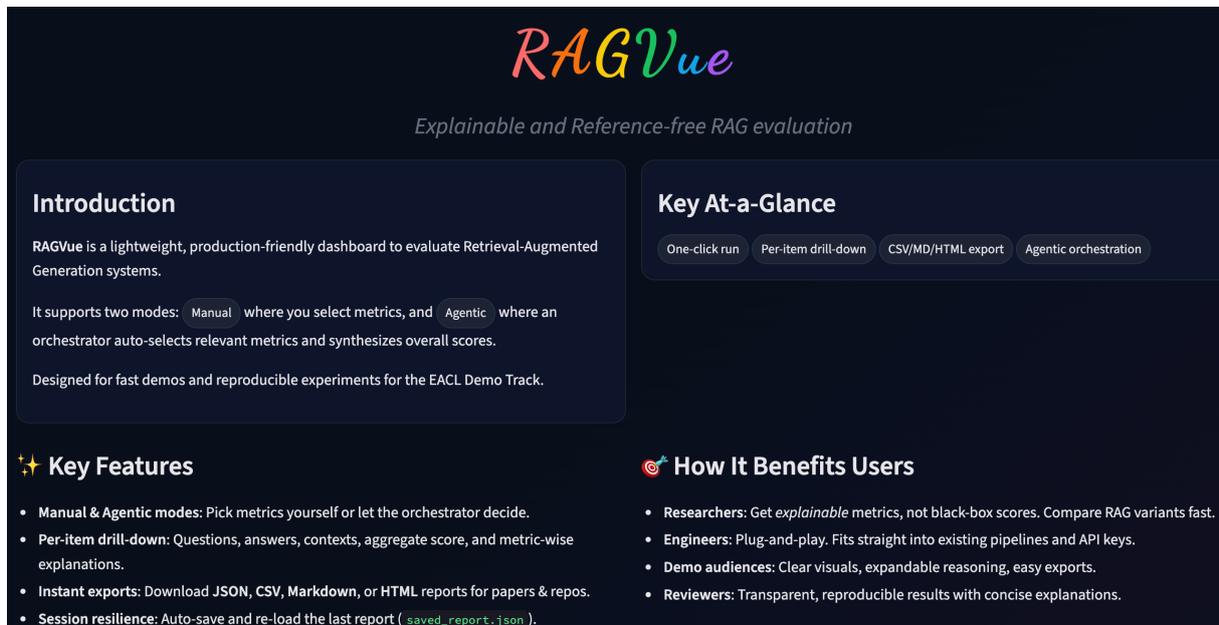
| RAGAS                 | RAGVUE metrics |         |         |         |         |              |
|-----------------------|----------------|---------|---------|---------|---------|--------------|
|                       | ans_comp.      | ans_rel | clarity | ret_cov | ret_rel | strict_faith |
| faithfulness          | 0.553          | 0.668   | 0.094   | 0.298   | 0.025   | 0.739        |
| answer_relevancy      | 0.704          | 0.644   | 0.172   | 0.193   | 0.053   | 0.958        |
| context_relevance     | 0.082          | 0.062   | 0.171   | 0.006   | 0.708   | 0.035        |
| response_groundedness | 0.373          | 0.174   | 0.071   | 0.071   | 0.124   | 0.940        |

Table 4: Qualitative examples comparing RAGAS and RAGVUE on real evaluation outputs.

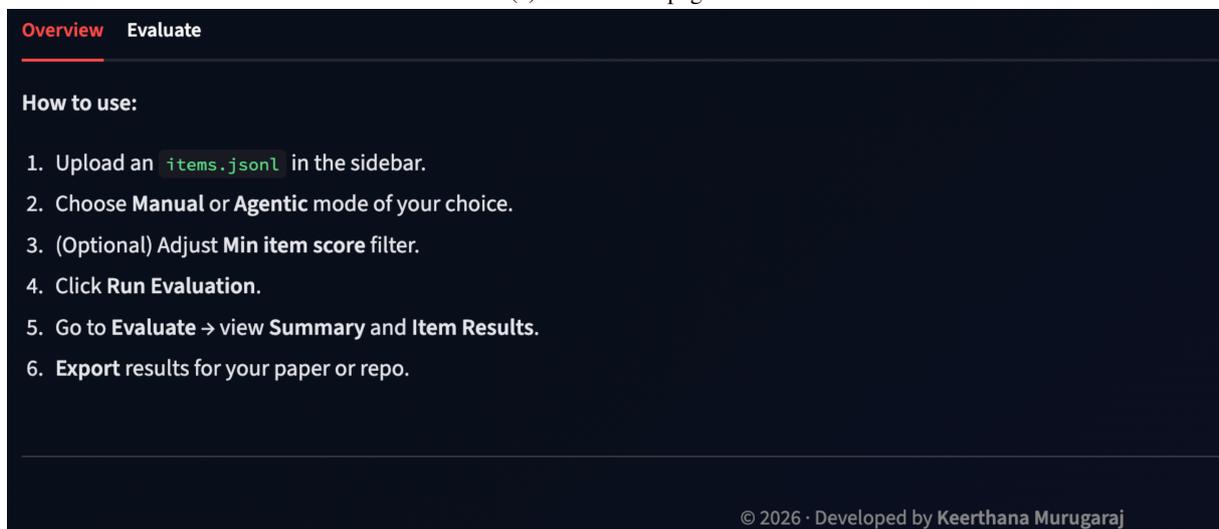
| Question                                                          | Model Answer                                                              | RAGAS Scores                                                                     | RAGVUE Scores                                                                                       |
|-------------------------------------------------------------------|---------------------------------------------------------------------------|----------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------|
| Are more people today related to Genghis Khan than Julius Caesar? | “No, even though there is no strong supporting evidence...”               | Faithfulness: 1.0<br>Answer rel.: 0.0<br>Context rel.: 0.5<br>Groundedness: 0.0  | Strict faith.: 0.0<br>Answer rel.: 0.40<br>Completeness: 0.00<br>Ret. cov.: 0.33<br>Ret. rel.: 0.33 |
| Could the members of <i>The Police</i> perform lawful arrests?    | “Yes, even though there is no strong supporting evidence...”              | Faithfulness: 0.5<br>Answer rel.: 0.0<br>Context rel.: 1.0<br>Groundedness: 0.0  | Strict faith.: 0.0<br>Answer rel.: 0.40<br>Completeness: 0.00<br>Ret. cov.: 0.33<br>Ret. rel.: 1.00 |
| Would a dog respond to a bell before a grey seal?                 | “It is hard to say... probably true. The evidence is not entirely clear.” | Faithfulness: 0.67<br>Answer rel.: 0.0<br>Context rel.: 0.5<br>Groundedness: 0.0 | Strict faith.: 0.0<br>Answer rel.: 0.40<br>Completeness: 0.00<br>Ret. cov.: 0.50<br>Ret. rel.: 0.33 |

## E Streamlit UI Images

The following figures illustrate the full RAGVUE Streamlit interface and its functionality. Figure 5a-5b provides the introduction page and the overview tab, which guide users through the workflow and usage instructions. Figure 6a-6d presents the core configuration components, including API key setup, data selection, manual and agentic mode configuration, and optional filters and report-saving tools. Figure 7a-7b shows the evaluation tab, which contains both the global summary across all processed cases and the detailed report for an individual (Q, A, C) example. Finally, Figure 8a-8b demonstrates the behavior of the agentic orchestrator for different input formats, highlighting its ability to select appropriate metrics based on the available fields.

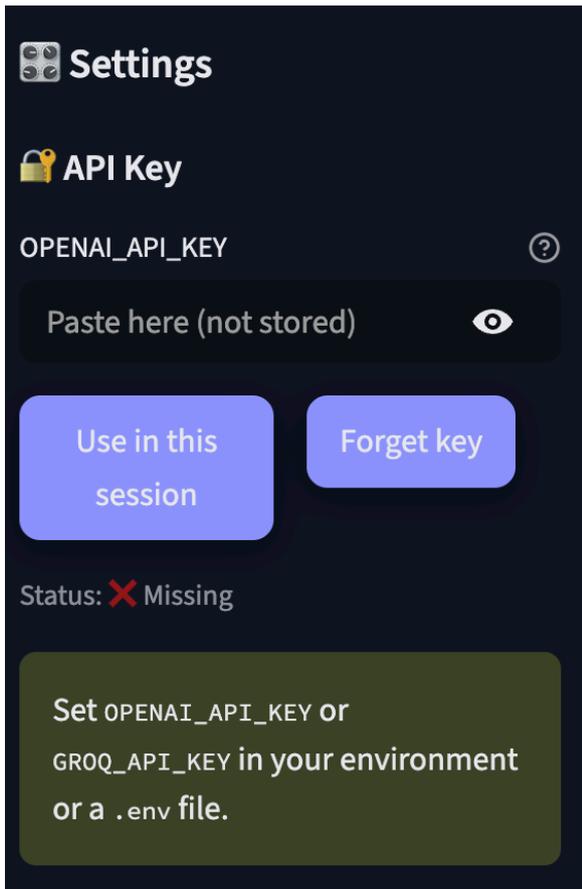


(a) Introduction page.

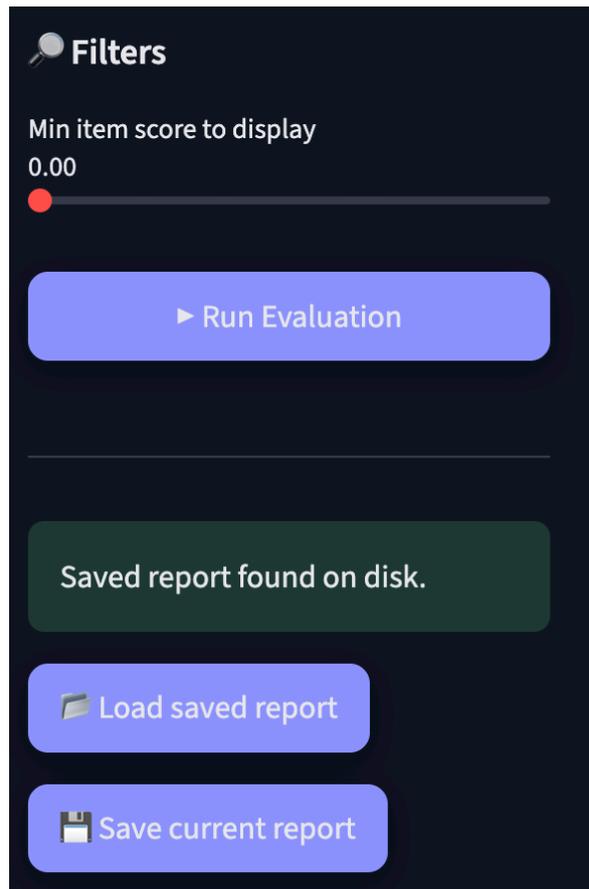


(b) Overview tab: usage instructions.

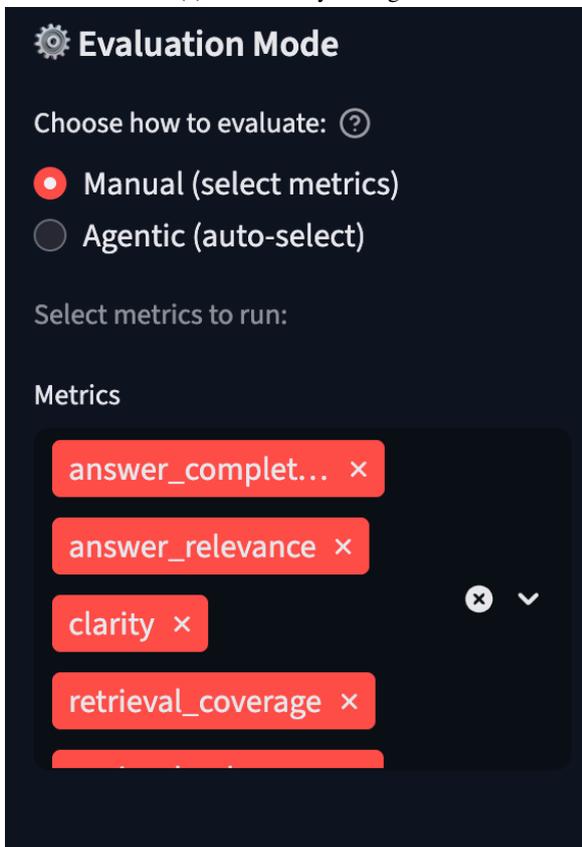
Figure 5: RAGVUE Streamlit UI: introduction page & overview tab.



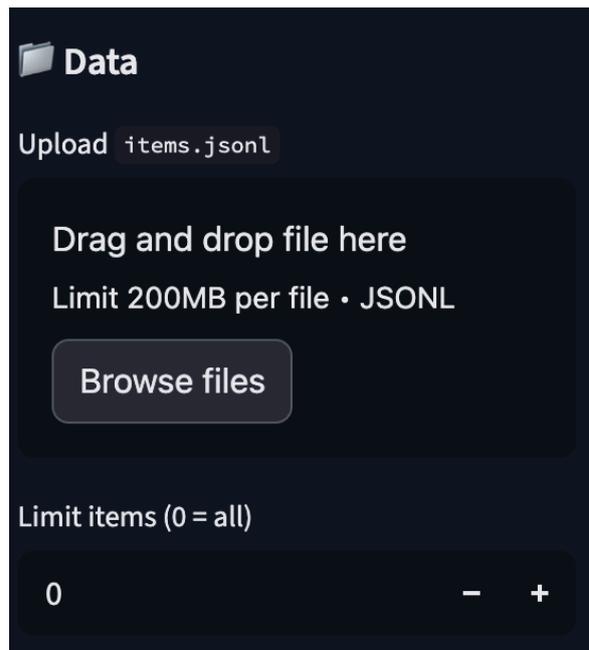
(a) API and key settings.



(b) Filters and saving options.

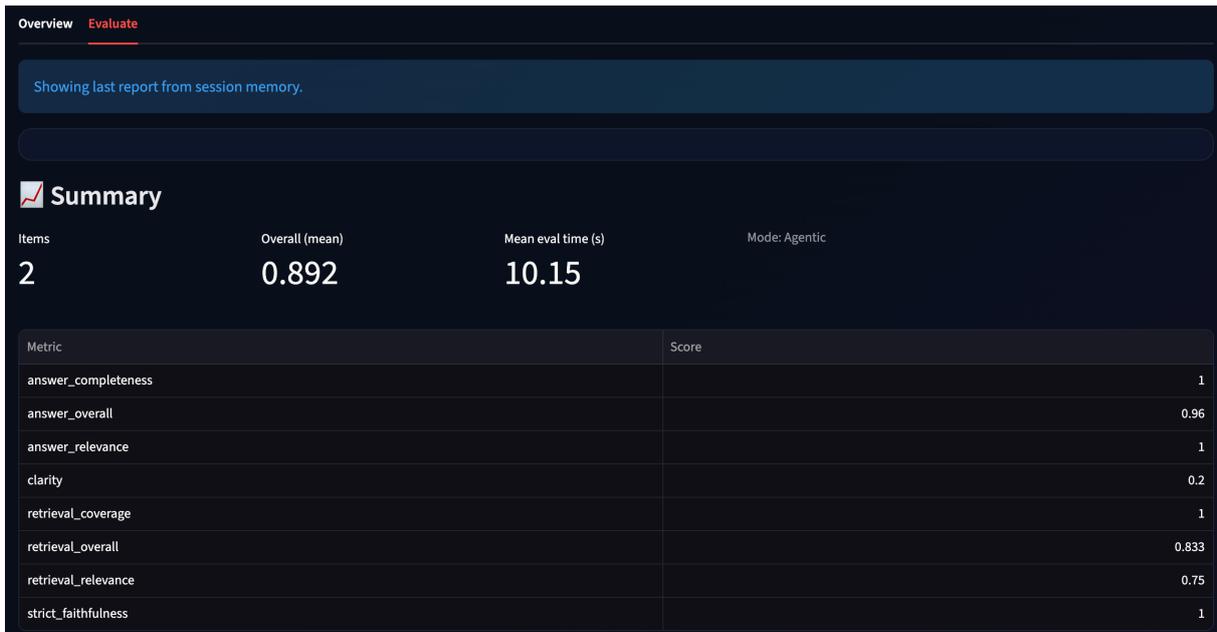


(c) Manual/agentic mode configuration.

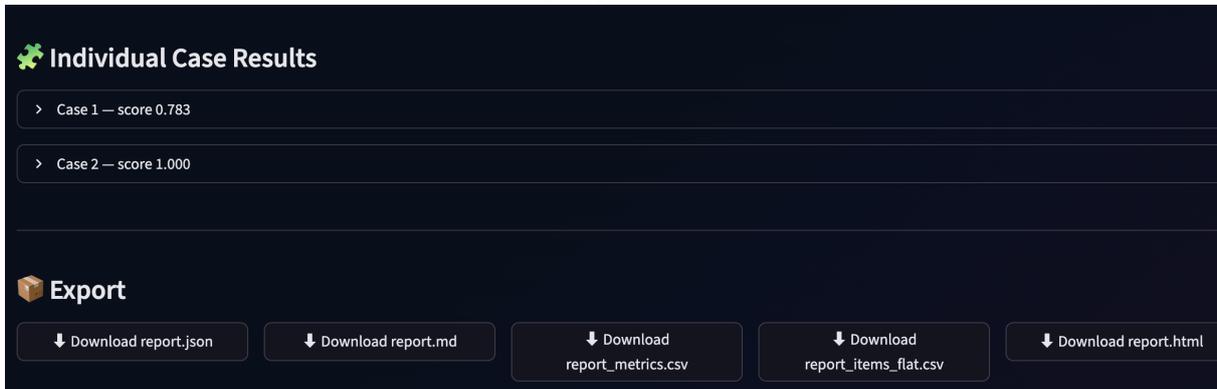


(d) Data upload option.

Figure 6: RAGVUE Streamlit user interface: API configuration, filter and save settings, manual/agentic mode selection, data upload options.



(a) Evaluate tab showing the global summary across all cases.



(b) Detailed individual case report for a single (Q, A, C) instance.

Figure 7: RAGVUE Streamlit user interface: evaluation summary view and individual case report.

Case 1 — score 0.783

**Question**  
What is the chemical symbol for gold?

**Answer**  
Au

**Contexts**  
[1] Gold is a chemical element with symbol Au and atomic number 79.  
[2] Copper has the symbol Cu and is highly conductive.

**Aggregate (case)**  
**0.783**

**Eval time (s)**  
**14.86**

Metrics computed: 8

**Metrics**

| Metric              | Score |
|---------------------|-------|
| retrieval_relevance | 0.5   |
| retrieval_coverage  | 1     |
| strict_faithfulness | 1     |
| answer_relevance    | 1     |
| answer_completeness | 1     |
| clarity             | 0.2   |
| retrieval_overall   | 0.667 |
| answer_overall      | 0.96  |

> Inspect JSON

(a) Agentic mode applied to (Q, A, C) triplets. The orchestrator correctly selects all relevant metrics.

Case 2 — score 1.000

**Question**  
Which planet is known as the Red Planet?

**Answer**  
∅ (no answer)

**Contexts**  
[1] Mars is called the Red Planet due to its reddish appearance.

**Aggregate (case)**  
**1.000**

**Eval time (s)**  
**5.45**

Metrics computed: 3

**Metrics**

| Metric              | Score |
|---------------------|-------|
| retrieval_relevance | 1     |
| retrieval_coverage  | 1     |
| retrieval_overall   | 1     |

> Inspect JSON

(b) Agentic mode applied to (Q, C) triplets. The orchestrator selects only retrieval metrics and skips answer metrics.

Figure 8: Agentic mode behavior across different input configurations.

# SMARTMATCH: Real-Time Semantic Retrieval for Translation Memory Systems

Ernesto L. Estevanell-Valladares<sup>1,2</sup>, Salima Lamsiyah<sup>3</sup>,  
Alicia Picazo-Izquierdo<sup>1</sup>, Tharindu Ranasinghe<sup>4</sup>, Ruslan Mitkov<sup>1</sup>, Rafael Muñoz<sup>1</sup>

<sup>1</sup>University of Alicante, Spain <sup>2</sup>University of Havana, Cuba

<sup>3</sup>University of Luxembourg, Luxembourg <sup>4</sup>Lancaster University, UK

ernesto.estevanell@ua.es

## Abstract

Translation Memory (TM) systems are core components of commercial computer-aided translation (CAT) tools. However, traditional fuzzy matching methods often fail to retrieve semantically relevant content when surface similarity is low. We introduce SMARTMATCH<sup>1</sup>, an open-source interactive demo and evaluation toolkit for TM retrieval that connects modern sentence encoders (including LLM-derived representations) and strong lexical/fuzzy baselines with a vector database, and exposes the end-to-end retrieval pipeline through a web-based UI for qualitative inspection and preference logging. The demo allows users to (i) enter a query segment, (ii) switch retrieval backends and embedding models, (iii) inspect top- $k$  retrieved matches with similarity scores and qualitative cues, and (iv) observe end-to-end latency in real time. We provide a reproducible benchmark on multilingual TM data, reporting retrieval quality using reference-based MT metrics (COMET, BERTScore, METEOR, chrF) together with coverage and latency/throughput trade-offs relevant to real-time CAT workflows. On DGT-TM, encoder-based retrieval achieves full coverage (100%) with millisecond-level latency (p50/p90  $\leq$  6–20 ms) and attains the strongest semantic-quality scores on the shared query set (e.g., BERTScore up to 0.91 at  $k=10$ ), while BM25 remains a strong lightweight lexical baseline with very low latency. SMARTMATCH targets CAT researchers and tool builders and bridges recent advances in sentence encoders with the real-time constraints of translation memory retrieval.

## 1 Introduction

Translation Memory (TM) systems are foundational components of computer-aided translation

tools (CAT), widely adopted in the translation industry to improve translator productivity and consistency (Mitkov, 2022; Reinke, 2018). By storing previously translated units (TU) and retrieving them during new translation tasks, TM systems enable TU reuse, remove redundant human effort, and significantly reduce both translation time and costs while maintaining quality (Paşcalau et al., 2025). Recent surveys<sup>2</sup> show that 88% of full-time professional translators use at least one CAT tool, and translation memory usage more than doubled from 2023 to 2024, demonstrating that TM systems remain indispensable to the translation industry despite advances in neural machine translation.

Traditional TM retrieval relies on surface-level fuzzy matching, which struggles when lexical overlap is low (Baquero and Mitkov, 2017; Gupta et al., 2016). To mitigate this, prior work has introduced enhancements such as paraphrasing (Gupta and Orasan, 2014) and clause splitting (Timonera and Mitkov, 2015). Despite major advances in semantic retrieval driven by large language models (LLMs) (Muennighoff et al., 2023; Ranasinghe et al., 2025), TM systems have not fully adopted these methods due to concerns about latency and scalability in real-time CAT workflows. While a few studies have explored sentence encoders for TM retrieval (Zhang et al., 2020; Ranasinghe et al., 2020), these efforts have largely relied on earlier encoder architectures. Modern LLM-based sentence encoders, which demonstrate superior semantic understanding, remain unexplored in TM systems, which leaves a significant gap between state-of-the-art retrieval models and TMs.

In this paper, we describe SMARTMATCH, an open-source system for LLM-based semantic retrieval in translation memory which we developed. Unlike traditional fuzzy matching, SMARTMATCH

<sup>1</sup><https://github.com/EEstevanell/SmartMatch>

<sup>2</sup><https://lokalise.com/library/data-reports/localization-trends-2025/>

uses state-of-the-art sentence encoders within a unified pipeline that also includes strong lexical baselines (Okapi/Pensieve and BM25). We provide a reproducible benchmark on DGT-TM, showing that encoder-based retrieval achieves 100% coverage and state-of-the-art semantic quality with millisecond-level latency, and we release an interactive web UI that exposes these quality–latency trade-offs for real-time CAT scenarios.

The **main contributions** of this paper are the following:

(i) We **introduce** SMARTMATCH, an open-source solution for LLM-based semantic retrieval in translation memory systems incorporating state-of-the-art encoder models and vector databases.

(ii) We **empirically evaluate** a diverse range of embedding models for TM retrieval, comparing popular open-source encoders against industry-relevant baselines (Okapi Pensieve and BM25). On the DGT-TM benchmark, encoder-based retrieval achieves 100% coverage and attains the best COMET and BERTScore results at ranks  $k \in \{5, 10\}$  on the shared query set, while BM25 remains a very strong and extremely fast lexical baseline (Table 2).

(iii) We **analyse** latency across different architectures, demonstrating millisecond end-to-end retrieval for encoder-based setups ( $p50/p90 \leq 6\text{--}20$  ms) and contrasting this with the heavy-tail behavior of Okapi-based retrieval ( $p90 \approx 1.1$  s).

(iv) We **release** an interactive web-based UI that allows users to enter text and retrieve the best match from the TM, switch models on the fly, and directly observe quality–latency trade-offs.

## 2 Related Work

Research on Translation Memory (TM) has evolved from early rule-based systems (Sato and Nagao, 1990; Schjoldager and Christensen, 2010) to commercial deployments (Reinke, 2018; Mitkov, 2022) supported by large multilingual resources such as DGT-TM (Steinberger et al., 2012). Prior work examined TM’s practical impact on translation quality and professional workflows (Jiménez-Crespo, 2010; Baquero and Mitkov, 2017), alongside its ethical and legal dimensions (Park, 2024). To improve TM reliability, shared tasks and unsupervised cleaning methods have addressed redundancy and noise (Barbu et al., 2016; Sabet et al., 2016; Negri et al., 2017; Wolff, 2016).

With the rise of machine translation, significant efforts have focused on integrating TM into statistical and neural MT frameworks. Initial work on SMT explored the use of TM features during decoding (Wang et al., 2013; Li et al., 2014), which later evolved into neural approaches employing gating, contrastive learning, and cross-lingual retrieval (Cao and Xiong, 2018; Cheng et al., 2022; Tamura et al., 2023). Additionally, studies have shown that TM-enhanced outputs can improve translator productivity and reduce post-editing time (Sánchez-Gijón et al., 2019; Green et al., 2014). To further improve matching quality, methods such as paraphrasing (Gupta et al., 2016), clause splitting (Timonera and Mitkov, 2015), and context-aware retrieval (He et al., 2019) have been proposed. These approaches aim to increase the utility of TM suggestions in real-time translation environments.

Recent research has also explored embedding-based retrieval as an alternative to traditional fuzzy matching. Ranasinghe et al. (2020, 2021) reported early gains using sentence encoders, though their study was limited in model scope and scalability. In parallel, TM integration with large language models has gained increasing attention, which demonstrates benefits from similarity-based prompting, retrieval augmentation, and dynamic memory mechanisms (Mu et al., 2023; Qian and Kong, 2024; Xu et al., 2023). Our work extends this line by benchmarking a broader range of modern, publicly available models under realistic latency constraints.

## 3 SMARTMATCH System Overview

SMARTMATCH is a research-orientated, end-to-end toolkit and interactive demo for comparing TM retrieval strategies. It offers the full retrieval workflow, including experiment setup, indexing, querying, and inspection, through a web-based user interface, while remaining easy to deploy and reproduce. The system is designed to (i) support fair comparisons across retrieval paradigms (fuzzy, lexical, dense), (ii) satisfy low-latency requirements typical of CAT workflows, and (iii) enable rapid qualitative assessment via interactive exploration and feedback logging. SMARTMATCH is implemented in Python 3.11+ and integrates Weaviate for indexing/retrieval, LangChain for embedding/model orchestration, and Typer for experiment configuration and command-line execution.

### 3.1 Architecture at a Glance

The system is organised into two phases:

**Offline indexing.** Given a TM collection  $\mathcal{D} = \{(s_i, t_i)\}_{i=1}^N$ , SMARTMATCH prepares three retrieval backends: (i) a fuzzy-matching service (Okapi/Pensieve) over raw segments, (ii) a lexical index (BM25) for sparse retrieval, and (iii) a dense vector index for semantic retrieval. Users first select an experiment setting that samples a TM subset for indexing (e.g., a random training split) and specify indexing parameters such as the embedding model and batch size. For dense retrieval, SMARTMATCH iterates over the chosen TM subset, computes embeddings  $E(s_i)$  for all source segments, and writes them into a dedicated Weaviate collection indexed with an HNSW ANN structure. In parallel, Weaviate maintains an inverted index to support BM25 retrieval over the same segments, ensuring that lexical and dense retrieval operate over identical content and can be compared under consistent database conditions. For fuzzy matching, SMARTMATCH triggers a Pensieve preparation step that generates the memory files required for edit-distance retrieval. Embeddings are computed by the selected encoder in the embedding service and stored in Weaviate; Weaviate is used as a common retrieval substrate (ANN + BM25) to enable controlled comparisons under consistent database conditions.

**Index readiness checks.** To prevent ambiguous “no-results” behaviour, SMARTMATCH enforces index readiness: if the selected Weaviate collection has not been seeded (e.g., empty collection) or Okapi/Pensieve memory files are missing, the retrieval interface ‘declines’ to run. This makes it explicit when indexing has not completed and avoids silent failures.

**Online retrieval.** At query time, a user submits a source segment  $q$  via the UI. The query can be selected from the held-out test split or entered as custom text. The request is routed to one of the enabled backends (fuzzy, BM25, or dense). For dense retrieval, the selected encoder produces  $E(q)$  and Weaviate returns the top- $k$  nearest neighbours under cosine similarity. For lexical retrieval, Weaviate returns the top- $k$  BM25 matches. For fuzzy retrieval, the Okapi/Pensieve service returns top- $k$  segments scored by normalised edit-distance similarity. Results are converted into a unified response

schema consumed by the UI and include end-to-end latency measured in milliseconds.

### 3.2 Retrieval Backends

SMARTMATCH provides three interchangeable retrieval modes:

**Fuzzy matching (Okapi/Pensieve).** To reflect legacy CAT deployments, we integrate a Java-based Okapi/Pensieve service that computes normalised edit-distance scores over the TM. This backend serves as a traditional production-like baseline consistent with common CAT tool practice.

**Lexical retrieval (BM25 in Weaviate).** We use Weaviate’s native BM25 implementation<sup>3</sup> for sparse retrieval, benefiting from the same storage layer used by dense indexing. This provides a strong keyword-based baseline and enables direct latency comparison under consistent database settings.

**Dense retrieval (encoders + HNSW in Weaviate).** For each encoder listed in §4.1, SMARTMATCH indexes pre-computed vectors and retrieves candidates via HNSW-based ANN search. We score candidates using cosine similarity and return the top- $k$  matches together with similarity scores.

### 3.3 Unified API and Result Schema

To support interactive comparison across strategies, all backends return results in a common schema:  $\{query, backend, model, k, [(s_i, t_i, score)]\}$ . This enables the UI to render comparable top- $k$  lists, sort by score, and optionally display metadata (e.g., segment length and language). The API also records end-to-end latency for each request (from query submission to response), enabling real-time inspection of quality-latency trade-offs and aggregate reporting (e.g., p50/p90 latency in experiments).

### 3.4 Interactive Demo Interface and Feedback

The web UI is designed for rapid qualitative inspection and live model comparison. Users can: (i) input a query (test sample or custom sentence), (ii) select a retrieval backend and (for dense retrieval) an embedding model, (iii) choose  $k$  (e.g.,  $k \in \{1, 5, 10\}$ ), and (iv) view retrieved source segments alongside their target translations. The interface reports per-query latency in milliseconds and

<sup>3</sup><https://huggingface.co/blog/xhluca/bm25s>

supports view customisation so users can hide non-essential fields. Beyond exploration, a dedicated “Feedback Charts” page turns SMARTMATCH into an interactive evaluation tool: for each query, users can compare the top suggestions from multiple backends and select their preferred output, taking match quality and latency into account. Each decision is logged to a JSON feedback file, and the chart visualisation updates live as new preferences are recorded.

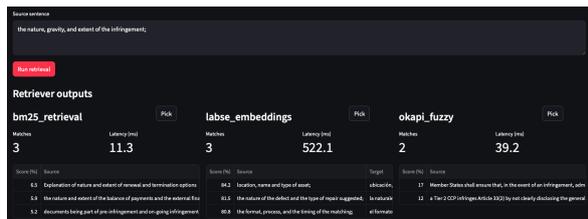


Figure 1: TM Retrieval tool

### 3.5 Deployment

SMARTMATCH is intended to be deployed as a lightweight set of services: (1) a Weaviate instance for BM25 and vector search, (2) an embedding service that exposes encoder inference for multiple sentence encoders, (3) an Okapi/Pensieve wrapper for fuzzy matching, and (4) a web UI for querying, comparison, and feedback collection. This modular design mirrors typical CAT pipeline integration points and allows users to enable only the components required for their workflow (e.g., lexical+dense retrieval without Okapi) while preserving a unified interface for experimentation.

## 4 Experimental Setup

This section summarises the experimental configuration used to benchmark the retrieval backends showcased in the SMARTMATCH demo. It is worth mentioning that all models are evaluated using *zero-shot* (no task-specific fine-tuning) and without cross-encoder re-ranking.

### 4.1 Retrieval Models

We evaluate 10 neural retrieval architectures, grouped into two categories.

**Bi-encoder sentence encoders.** We use encoder-style models that map a segment to a fixed-dimensional embedding and support ANN search: Multilingual E5 (*intfloat/multilingual-e5-base*,  $d=768$ ), BGE-M3 (*BAAI/bge-m3*,  $d=1024$ ), LaBSE (*sentence-transformers/LaBSE*,  $d=768$ ),

| Retrieval      | Text                                                                                | Latency |
|----------------|-------------------------------------------------------------------------------------|---------|
| source text    | Description of accounting policy that requires external consulting services         | -       |
| BM25           | Description of accounting policy for business combinations [text block]             | 5.1     |
| LaBSE          | The description of the entity’s material accounting policy information for hedging. | 66.0    |
| Okapi/Pensieve | No match found                                                                      | 1.2     |

Table 1: Example retrieval outputs from BM25, LaBSE, and Okapi/Pensieve for a low-lexical-overlap query.

LASER (*Facebook/LASER*,  $d=1024$ ), M3E-Base (*moka-ai/m3e-base*,  $d=768$ ), and three SBERT-family checkpoints: MPNet (*all-mpnet-base-v2*,  $d=768$ ), DistilRoBERTa (*all-distilroberta-v1*,  $d=768$ ), and MiniLM (*all-MiniLM-L6-v2*,  $d=384$ ).

**Generative LLM architectures.** To assess whether hidden states from generative models can serve as retrieval representations, we evaluate: Llama-3.2 (*meta-llama/Llama-3.2-1B*,  $d=2048$ ), using the final-layer hidden state of the last token as the segment representation, and T5 (*t5-base*,  $d=768$ ), using the encoder output as the segment representation.

**Non-neural baselines.** We include two industry-relevant baselines exposed in the demo UI: (i) fuzzy matching using normalised Levenshtein similarity via Okapi/Pensieve, and (ii) lexical retrieval using BM25 implemented in Weaviate.

### 4.2 Datasets

We run experiments on the DGT Translation Memory released by the European Commission (Directorate-General for Translation) and the Joint Research Centre. The dataset contains professionally produced translation units across 22 official EU languages and multiple language-pair combinations, including domain-specific texts.<sup>4</sup> While we selected the English-Spanish pairs for this test, our system can include any language pair. Within this TM, we used the following train-test split: Volume 1-4 were used for training purposes, and Volume 5 for testing.

### 4.3 Indexing and Retrieval Configuration

For each experimental run, we create a dedicated TM index by sampling a *training* subset used for indexing and reserving a held-out *test* subset for

<sup>4</sup>[https://joint-research-centre.ec.europa.eu/language-technology-resources/dgt-translation-memory\\_en](https://joint-research-centre.ec.europa.eu/language-technology-resources/dgt-translation-memory_en)

querying. Dense retrieval follows the standard bi-encoder pipeline: we pre-compute embeddings for all indexed source segments and store them in a Weaviate collection, which is configured with an HNSW ANN index for cosine similarity search. The same collection stores the raw text fields used by Weaviate’s BM25 inverted index, enabling lexical retrieval over identical content. For the fuzzy baseline, Pensieve memory files are generated from the indexed segments prior to querying.

#### 4.4 Hardware

All experiments were conducted on a workstation with an Intel(R) Core(TM) i7-14700 CPU (28 cores, up to 5.4 GHz), 64 GB RAM, and an NVIDIA RTX 4090 (24 GB VRAM). This configuration provides stable latency measurements and sufficient resources for embedding computation and large-scale indexing.

#### 4.5 Evaluation Metrics

We evaluate retrieval quality with four reference-based MT metrics: COMET (Rei et al., 2020), BERTScore (Zhang et al., 2019), METEOR (Banerjee and Lavie, 2005), and chrF (Popović, 2015). We additionally report **Coverage** (percentage of queries with at least one retrieved TU) and **Latency** (end-to-end retrieval time), summarised by p50 and p90 to capture typical and tail behaviour. We evaluate ranking at  $k \in \{1, 5, 10\}$  using *best-in-k* (oracle) scoring, i.e., the maximum metric score among the top- $k$  candidates, reflecting CAT usage where translators select the most helpful suggestion from a shortlist.

#### 4.6 Statistical Significance Testing

To ensure robustness, we performed statistical testing on the intersection of queries for which all models retrieved at least one TU. We used the non-parametric Wilcoxon signed-rank test (Rosner et al., 2006) to compare the Pensieve baseline against neural models, and applied the Benjamini–Hochberg procedure (Benjamini and Hochberg, 1995) to control the false discovery rate across multiple comparisons ( $\alpha = 0.05$ ). We also report Cliff’s delta ( $\delta$ ) (Cliff, 1993) as a non-parametric effect size and treat  $|\delta| \geq 0.474$  as non-negligible.

### 5 Results

Table 2 reports the retrieval performance of all evaluated models across four metrics (COMET,

BERTScore, METEOR, chrF) and three ranks (1, 5, 10). We highlight models that are both statistically significantly better than the Okapi Pensieve baseline ( $p < 0.05$ ) and demonstrate a large effect size ( $|\delta| \geq 0.474$ ).

Two distinct performance views for each evaluation metric are presented in Table 2. The first column reports the global performance on the full query set, calculated without penalising for lower coverage. It is important to note that these global metrics tend to favour models with lower coverage (such as Okapi Pensieve and BM25), as they effectively evaluate only on a subset of more straightforward queries where high lexical overlap exists. The second column, marked with  $\cap$ , reports performance on the intersection set of queries where all models successfully retrieved at least one result (100% coverage), which eliminates selection bias and provides a more accurate representation of comparative retrieval quality for statistical significance testing.

As shown in Table 2, modern semantic retrieval models consistently outperform traditional baselines. It is worth noting that LaBSE achieves the highest raw scores across all metrics and ranks, demonstrating its robustness as a general-purpose bitext retriever. Statistically, LaBSE, LASER, and T5 base achieve significant improvements with large effect sizes ( $|\delta| \geq 0.474$ ) over the Pensieve baseline, specifically in COMET and BERTScore at ranks 5 and 10.

This performance divergence underscores a fundamental shift in retrieval paradigms. Vector-based models excel at capturing semantic similarity even with low lexical overlap, leading to higher COMET and BERTScore values. Conversely, traditional baselines like Pensieve and BM25 fare very well – for surface-form matches, maintaining strong performance on n-gram based metrics like METEOR and chrF but failing to retrieve semantically equivalent yet lexically distinct segments. Table 1 illustrates a typical retrieval case, showing how BM25 and LaBSE retrieve semantically related segments while the fuzzy baseline fails due to low lexical overlap.

From a practical production perspective, these quality gains are complemented by superior operational characteristics. In our testing, semantic models achieved 100% coverage, significantly outperforming Okapi Pensieve (81%) and BM25 (85%), which fail to retrieve results when lexical overlap

| Model           | R. | COMET       | $\cap$      | BERTScore   | $\cap$      | METEOR      | $\cap$      | chrF        | $\cap$      | Lat. p50/p90 | Cov.       |
|-----------------|----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|------------|
| Okapi Pensieve  | 1  | <b>0.83</b> | 0.84        | <b>0.86</b> | 0.87        | <b>0.72</b> | <b>0.74</b> | <b>0.72</b> | <b>0.74</b> | 44/1116 ms   | 81         |
|                 | 5  | 0.84        | 0.85        | 0.87        | 0.88        | 0.75        | 0.77        | <b>0.75</b> | 0.77        |              |            |
|                 | 10 | 0.84        | 0.86        | 0.87        | 0.89        | 0.75        | 0.77        | <b>0.75</b> | 0.77        |              |            |
| BM25            | 1  | 0.80        | 0.80        | 0.85        | 0.84        | 0.69        | 0.69        | 0.59        | 0.69        | 3/13 ms      | 85         |
|                 | 5  | <b>0.85</b> | 0.85        | <b>0.88</b> | 0.88        | <b>0.76</b> | 0.76        | 0.66        | 0.76        |              |            |
|                 | 10 | <b>0.86</b> | 0.86        | <b>0.89</b> | 0.89        | <b>0.78</b> | 0.77        | 0.68        | 0.77        |              |            |
| BGE-M3          | 1  | 0.79        | 0.85        | 0.81        | 0.88        | 0.59        | 0.73        | 0.61        | <b>0.74</b> | 16/17 ms     | <b>100</b> |
|                 | 5  | 0.82        | 0.87        | 0.84        | <b>0.90</b> | 0.64        | 0.77        | 0.66        | 0.77        |              |            |
|                 | 10 | 0.82        | <b>0.88</b> | 0.84        | <u>0.90</u> | 0.66        | 0.78        | 0.66        | 0.78        |              |            |
| Multilingual E5 | 1  | 0.76        | 0.83        | 0.80        | 0.87        | 0.57        | 0.71        | 0.58        | 0.71        | 16/18 ms     | <b>100</b> |
|                 | 5  | 0.80        | 0.86        | 0.83        | 0.89        | 0.63        | 0.76        | 0.63        | 0.76        |              |            |
|                 | 10 | 0.81        | 0.87        | 0.84        | 0.90        | 0.64        | 0.77        | 0.65        | 0.77        |              |            |
| LaBSE           | 1  | 0.79        | 0.86        | 0.82        | 0.89        | 0.61        | <b>0.74</b> | 0.62        | <b>0.74</b> | 10/12 ms     | <b>100</b> |
|                 | 5  | 0.82        | <b>0.88</b> | 0.84        | <u>0.90</u> | 0.65        | <b>0.78</b> | 0.66        | <b>0.78</b> |              |            |
|                 | 10 | 0.83        | <b>0.88</b> | 0.85        | <u>0.91</u> | 0.67        | <b>0.79</b> | 0.67        | <b>0.79</b> |              |            |
| LASER           | 1  | 0.79        | 0.85        | 0.82        | 0.89        | 0.61        | 0.74        | 0.62        | <b>0.74</b> | 7/14 ms      | <b>100</b> |
|                 | 5  | 0.82        | 0.87        | 0.84        | <b>0.90</b> | 0.66        | <b>0.78</b> | 0.66        | <b>0.78</b> |              |            |
|                 | 10 | 0.82        | <b>0.88</b> | 0.84        | <u>0.91</u> | 0.67        | <b>0.79</b> | 0.67        | 0.78        |              |            |
| Llama-3.2-1B    | 1  | 0.78        | 0.85        | 0.81        | 0.88        | 0.58        | 0.71        | 0.60        | 0.73        | 17/20 ms     | <b>100</b> |
|                 | 5  | 0.81        | 0.87        | 0.83        | <b>0.90</b> | 0.62        | 0.75        | 0.64        | 0.76        |              |            |
|                 | 10 | 0.82        | 0.87        | 0.84        | 0.90        | 0.64        | 0.76        | 0.65        | 0.77        |              |            |
| M3E             | 1  | 0.77        | 0.84        | 0.81        | 0.88        | 0.59        | 0.72        | 0.60        | 0.73        | 10/12 ms     | <b>100</b> |
|                 | 5  | 0.81        | 0.87        | 0.83        | <b>0.90</b> | 0.63        | 0.77        | 0.65        | 0.77        |              |            |
|                 | 10 | 0.81        | 0.87        | 0.84        | 0.90        | 0.65        | 0.78        | 0.66        | 0.78        |              |            |
| MiniLM-L6-v2    | 1  | 0.77        | 0.84        | 0.81        | 0.88        | 0.58        | 0.72        | 0.60        | 0.73        | 6/8 ms       | <b>100</b> |
|                 | 5  | 0.81        | 0.87        | 0.83        | <b>0.90</b> | 0.63        | 0.76        | 0.65        | 0.77        |              |            |
|                 | 10 | 0.82        | <b>0.88</b> | 0.84        | 0.90        | 0.65        | 0.78        | 0.66        | 0.78        |              |            |
| MPNet-base-v2   | 1  | 0.78        | 0.84        | 0.81        | 0.87        | 0.58        | 0.72        | 0.60        | 0.73        | 11/14 ms     | <b>100</b> |
|                 | 5  | 0.81        | 0.87        | 0.83        | 0.89        | 0.63        | 0.76        | 0.65        | 0.77        |              |            |
|                 | 10 | 0.82        | <b>0.88</b> | 0.84        | 0.90        | 0.65        | 0.77        | 0.66        | 0.78        |              |            |
| DistilRoBERTa   | 1  | 0.78        | 0.85        | 0.81        | 0.88        | 0.59        | 0.72        | 0.61        | <b>0.74</b> | 7/10 ms      | <b>100</b> |
|                 | 5  | 0.81        | 0.87        | 0.83        | <b>0.90</b> | 0.63        | 0.76        | 0.65        | 0.77        |              |            |
|                 | 10 | 0.82        | <b>0.88</b> | 0.84        | 0.90        | 0.65        | 0.77        | 0.66        | 0.78        |              |            |
| T5 base         | 1  | 0.80        | 0.86        | 0.82        | 0.89        | 0.60        | 0.73        | 0.62        | <b>0.74</b> | 9/11 ms      | <b>100</b> |
|                 | 5  | 0.82        | <b>0.88</b> | 0.84        | <b>0.90</b> | 0.65        | 0.77        | 0.66        | <b>0.78</b> |              |            |
|                 | 10 | 0.83        | <b>0.88</b> | 0.85        | <u>0.91</u> | 0.66        | 0.78        | 0.67        | 0.78        |              |            |

Table 2: Retrieval performance on the DGT-TM dataset. The column with the symbol  $\cap$  reports values in the intersection set of queries where all models successfully retrieved at least one translation unit. **Bold** indicates the best performance per metric and rank. Underline denotes statistically significant improvement over the Okapi Pensieve baseline ( $p < 0.05$ ) with a large effect size ( $|\delta| \geq 0.474$ ). Latency (p50/p90) is measured in milliseconds.

is insufficient.

Regarding efficiency, while BM25 remains the fastest (3ms p50), lightweight semantic models like MiniLM-L6-v2 (6ms) and LASER (7ms) offer competitive speeds suitable for real-time applications. Notably, Okapi Pensieve exhibits high tail latency (1116ms p90), likely due to its reliance on complex fuzzy matching heuristics. In contrast, Weaviate’s HNSW index shows consistent, low-latency approximate nearest neighbour search for semantic models.

## 6 Conclusions

In this paper we introduce SMARTMATCH, an open-source semantic retrieval system for TM. The performance of system clearly shows that modern sentence encoders with a vector database can be included into a TM pipeline with very low latency rates for dense retrieval. The system offers an architecture that enables comparison between different retrieval strategies and allows users to inspect source-target pairs with the latency per query displayed. The results show that semantic models outperform traditional edit-distance-based methods when lexical overlap is low, and the coverage rates

reveal that traditional methods fail to return candidates in almost 1 out of 5 queries.

A key limitation of our paper is that it does not cover integration aspects that real CAT tools need, such as user management, security, plugin APIs to existing CAT environments, or multi-tenant scaling. While the retrieval core is realistic, it would be beneficial to address these workflows. Another limitation is that there is no mechanism to automatically adapt retrieval configuration based on the human feedback. However, a human-in-the-loop reinforcement learning TM retrieval system could be explored based on our approach.

## 7 Acknowledgements

This research has been funded by the Generalitat Valenciana (Conselleria d'Educació, Investigació, Cultura i Esport) through the project: The limits and future of data-driven approaches: A comparative study of deep learning, knowledge-based and rule-based models and methods in Natural Language Processing (CIDEXG/2023/13), and by the Ministry of Science, Innovation, and Universities through the project *Safewords* (AIA2025-163322-C63). Publication fees and conference participation were supported by the University of Luxembourg.

## 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.
- Andrea Silvestre Baquero and Ruslan Mitkov. 2017. Translation memory systems have a long way to go. In *Proceedings of the Workshop Human-Informed Translation and Interpreting Technology*, pages 44–51.
- Eduard Barbu, Carla Parra Escartín, Luisa Bentivogli, Matteo Negri, Marco Turchi, Constantin Orasan, and Marcello Federico. 2016. The first automatic translation memory cleaning shared task. *Machine Translation*, 30(3):145–166.
- 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.
- Qian Cao and Deyi Xiong. 2018. Encoding gated translation memory into neural machine translation. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 3042–3047.
- Xin Cheng, Shen Gao, Lemao Liu, Dongyan Zhao, and Rui Yan. 2022. [Neural machine translation with contrastive translation memories](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3591–3601, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Norman Cliff. 1993. Dominance statistics: Ordinal analyses to answer ordinal questions. *Psychological bulletin*, 114(3):494.
- Spence Green, Jason Chuang, Jeffrey Heer, and Christopher D Manning. 2014. Predictive translation memory: A mixed-initiative system for human language translation. In *Proceedings of the 27th annual ACM symposium on User interface software and technology*, pages 177–187.
- Rohit Gupta and Constantin Orasan. 2014. Incorporating paraphrasing in translation memory matching and retrieval. In *Proceedings of the 17th Annual Conference of the European Association for Machine Translation*, pages 3–10.
- Rohit Gupta, Constantin Orăsan, Marcos Zampieri, Mihaela Vela, Josef van Genabith, and Ruslan Mitkov. 2016. Improving translation memory matching and retrieval using paraphrases. *Machine Translation*, 30(1):19–40.
- Qiuxiang He, Guoping Huang, Lemao Liu, and Li Li. 2019. Word position aware translation memory for neural machine translation. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 367–379. Springer.
- Miguel A Jiménez-Crespo. 2010. The effect of translation memory tools in translated web texts: Evidence from a comparative product-based study. *Linguistica antverpiensia*, 8:213–232.
- Liangyou Li, Andy Way, and Qun Liu. 2014. A discriminative framework of integrating translation memory features into smt. In *Proceedings of the 11th Conference of the Association for Machine Translation in the Americas: MT Researchers Track*, pages 249–260.
- Ruslan Mitkov. 2022. Translation memory systems. In *The Routledge Handbook of Translation and Memory*, pages 364–380. Routledge.
- Yongyu Mu, Abudurexiti Reheman, Zhiqian Cao, Yuchun Fan, Bei Li, Yinqiao Li, Tong Xiao, Chunliang Zhang, and Jingbo Zhu. 2023. [Augmenting large language model translators via translation memories](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10287–10299, Toronto, Canada. Association for Computational Linguistics.

- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2023. **MTEB: Massive text embedding benchmark**. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.
- Matteo Negri, Duygu Ataman, Masoud Jalili Sabet, Marco Turchi, and Marcello Federico. 2017. Automatic translation memory cleaning. *Machine Translation*, 31(3):93–115.
- Jiyoung Park. 2024. Ethical approach to translation memory reuse: discussions from copyright and business ethics perspectives. *Translation Studies*, 17(1):37–52.
- Raul Paşcalau, Laura-Rebeca Stiegelbauer, and Dumitru Mădălina Pantea. 2025. The importance of translation workflow in global business. *Studii de Ştiinţă şi Cultură*, 21(2).
- Maja Popović. 2015. **chrF: character n-gram F-score for automatic MT evaluation**. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Ming Qian and Chuiqing Kong. 2024. Enabling human-centered machine translation using concept-based large language model prompting and translation memory. In *International Conference on Human-Computer Interaction*, pages 118–134. Springer.
- Tharindu Ranasinghe, Hansi Hettiarachchi, Constantin Orasan, and Ruslan Mitkov. 2025. **MUSTS: Multilingual semantic textual similarity benchmark**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 331–353, Vienna, Austria. Association for Computational Linguistics.
- Tharindu Ranasinghe, Ruslan Mitkov, Constantin Orăsan, and Rocío Caro Quintana. 2021. Semantic textual similarity based on deep learning. *Corpora in Translation and Contrastive Research in the Digital Age: Recent advances and explorations*, 158:101.
- Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. 2020. **Intelligent translation memory matching and retrieval with sentence encoders**. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 175–184, Lisboa, Portugal. European Association for Machine Translation.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. **COMET: A neural framework for MT evaluation**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- Uwe Reinke. 2018. State of the art in translation memory technology. *Language Technologies for a Multilingual Europe; Rehm, G., Stein, D., Sasaki, F., Witt, A., Eds.*, pages 55–84.
- Bernard Rosner, Robert J Glynn, and Mei-Ling T Lee. 2006. The wilcoxon signed rank test for paired comparisons of clustered data. *Biometrics*, 62(1):185–192.
- Masoud Jalili Sabet, Matteo Negri, Marco Turchi, and Eduard Barbu. 2016. An unsupervised method for automatic translation memory cleaning. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 287–292.
- Pilar Sánchez-Gijón, Joss Moorkens, and Andy Way. 2019. Post-editing neural machine translation versus translation memory segments. *Machine Translation*, 33(1):31–59.
- Satoshi Sato and Makoto Nagao. 1990. Toward memory-based translation. In *COLING 1990 Volume 3: Papers Presented to the 13th International Conference on Computational Linguistics*.
- Anne Gram Schjoldager and Tina Paulsen Christensen. 2010. Translation-memory (tm) research: what do we know and how do we know it? *Hermes*, 44:89–101.
- Ralf Steinberger, Andreas Eisele, Szymon Kłoczek, Spyridon Pilos, and Patrick Schlüter. 2012. **DGT-TM: A freely available translation memory in 22 languages**. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 454–459, Istanbul, Turkey. European Language Resources Association (ELRA).
- Takuya Tamura, Xiaotian Wang, Takehito Utsuro, and Masaaki Nagata. 2023. Target language monolingual translation memory based nmt by cross-lingual retrieval of similar translations and reranking. In *Proceedings of Machine Translation Summit XIX, Vol. 1: Research Track*, pages 313–323.
- Katerina Raisa Timonera and Ruslan Mitkov. 2015. Improving translation memory matching through clause splitting. In *Proceedings of the Workshop Natural Language Processing for Translation Memories*, pages 17–23.
- Kun Wang, Chengqing Zong, and Keh-Yih Su. 2013. **Integrating translation memory into phrase-based machine translation during decoding**. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11–21, Sofia, Bulgaria. Association for Computational Linguistics.
- Friedel Wolff. 2016. Combining off-the-shelf components to clean a translation memory. *Machine Translation*, 30(3):167–181.
- Jitao Xu, Josep Crego, and François Yvon. 2023. **Integrating translation memories into non-autoregressive machine translation**. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1326–1338, Dubrovnik, Croatia. Association for Computational Linguistics.

Tianfu Zhang, Heyan Huang, Chong Feng, and Xiaochi Wei. 2020. Similarity-aware neural machine translation: reducing human translator efforts by leveraging high-potential sentences with translation memory. *Neural Computing and Applications*, 32(23):17623–17635.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

## A SMARTMATCH Seeding Page

To run SMARTMATCH, the first step is to feed the retrieval system with random samples from your TM.

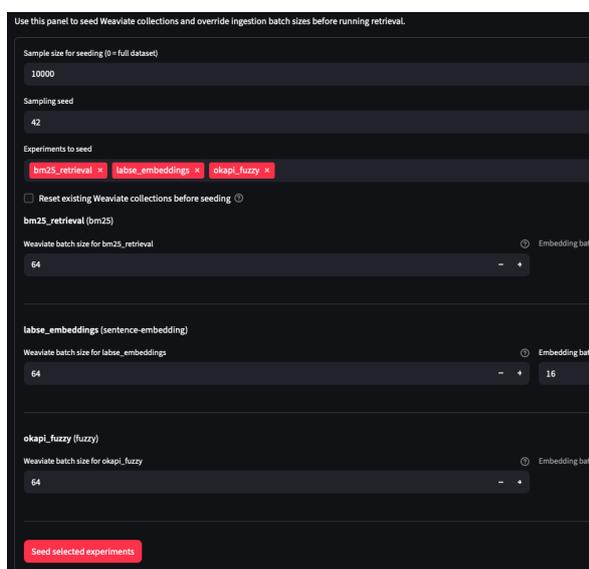


Figure 2: Seeding page

The experiment configuration samples random training data under the corresponding batch sizes. Then, the app iterates over the TM data, encodes the segments with the selected model, and stores them in a dedicated Weaviate collection. Before retrieval can run, both the semantic and BM25 models must be seeded into Weaviate, and Pensieve must generate the memory files. If a collection is empty, the retrieval page will not run.

## B SMARTMATCH Retrieval Page

In the Retrieval tab, a random test sample may be selected from the original test set. The output reports how many segments are matched within the translation memory, as well as the corresponding latencies. Once segments have been retrieved, the view can be configured to display only the required information (source text, target text, score

or index). In Figure 3, Okapi Pensieve serves as the traditional TM fuzzy matcher, as it compares strings and assigns higher scores when a segment is nearly identical to an entry in the memory.

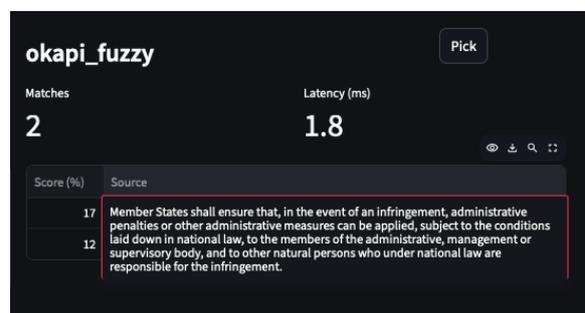


Figure 3: Okapi retrieval example

The second model, LaBSE, encodes each sentence into a vector and performs retrieval over a dense index in Weaviate, targeting segments with similar meaning even when the wording differs substantially, as shown in Figure 4.

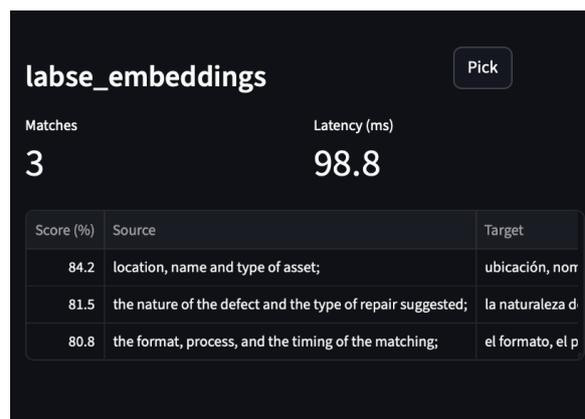


Figure 4: LaBSE retrieval example

BM25 is the last model, which operates on a lexical index in Weaviate. It is still word-based, but relies on term frequency and inverse document frequency, which enables more robust handling of word reordering and minor wording changes than simple edit distance.

After examining the results, the model that provides the most suitable concordances—taking into account both latency and match quality can be selected.

## C SMARTMATCH Feedback Page

The Feedback charts page provides an interactive evaluation interface. For each query, the top suggestions from the three models are displayed, and the preferred option is selected directly in the interface.

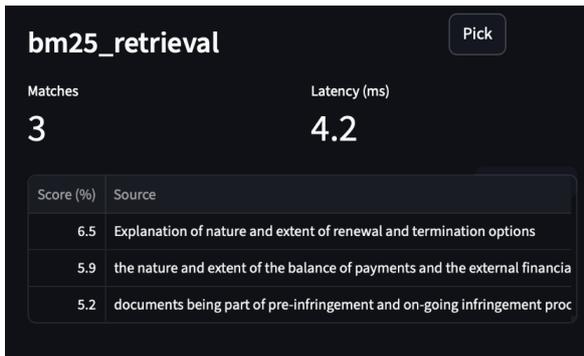


Figure 5: BM25 retrieval example

Each selection is recorded in a JSON feedback file, and the chart is updated in real time. This function-

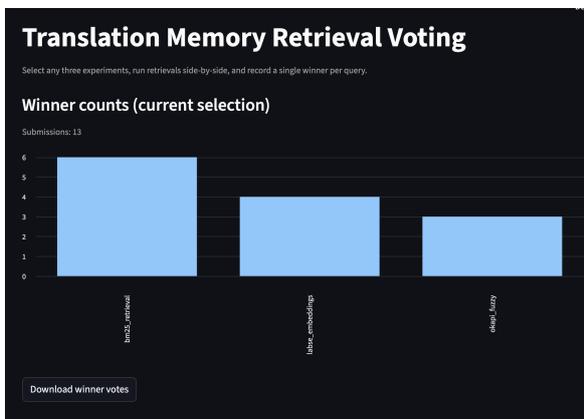


Figure 6: Feedback charts

ality can be used to compare model performance, identify which model performs best for specific segments, and collect preference data for further tuning or analysis.

# QSTN: A Modular Framework for Robust Questionnaire Inference with Large Language Models

Maximilian Kreutner<sup>1</sup>, Jens Rupperecht<sup>1</sup>, Georg Ahnert<sup>1</sup>,  
Ahmed Salem<sup>1</sup>, Markus Strohmaier<sup>1,2,3</sup>

<sup>1</sup>University of Mannheim, <sup>2</sup>GESIS - Leibniz Institute for the Social Sciences, <sup>3</sup>CSH Vienna

## Abstract

We introduce QSTN, an open-source Python framework for systematically generating responses from questionnaire-style prompts to support in-silico surveys and annotation tasks with large language models (LLMs). QSTN enables robust evaluation of questionnaire presentation, prompt perturbations, and response generation methods. Our extensive evaluation (>40 million survey responses) shows that question structure and response generation methods have a significant impact on the alignment of generated survey responses with human answers. We also find that answers can be obtained for a fraction of the compute cost, by changing the presentation method. In addition, we offer a no-code user interface that allows researchers to set up robust experiments with LLMs *without coding knowledge*. We hope that QSTN will support the reproducibility and reliability of LLM-based research in the future.

## 1 Introduction

Questionnaires have become an important format to probe, assess, and utilize large language models (LLMs) via prompts. Questionnaire-like prompts have been a popular way to evaluate LLMs on tasks such as common knowledge understanding (Hendrycks et al., 2021), language comprehension (Hu et al., 2023; Sravanthi et al., 2024; Kim et al., 2024), and mathematical reasoning (Satpute et al., 2024; Wei et al., 2023). Other work uses existing questionnaires to evaluate LLMs’ values; for example, political bias (Röttger et al., 2024; Rozado, 2024), personality traits (Jiang et al., 2024; Shu et al., 2024; Pellert et al., 2024), or psychometric profiles (Ye et al., 2025). With the increasing capability of LLMs, researchers have found additional use cases, such as the creation of synthetic survey responses (Argyle et al., 2023; Ma et al., 2024) or data annotation (Tan et al., 2024).

Despite the widespread use of questionnaire-like prompts, concerns have been raised about

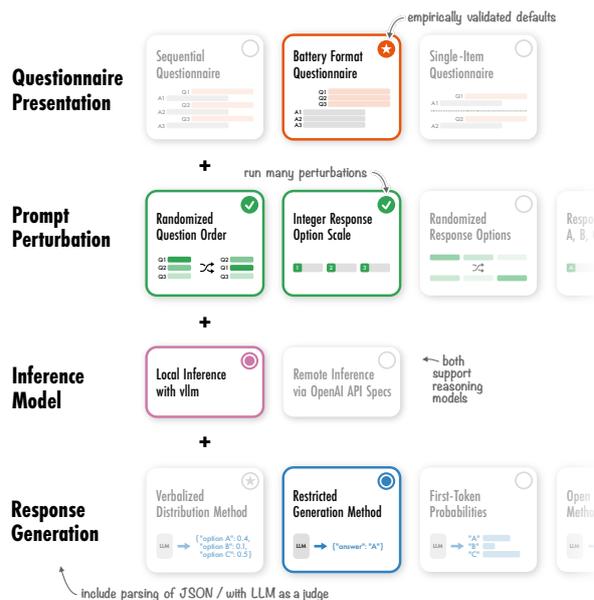


Figure 1: **QSTN Facilitates Easy To Customize and Robust Questionnaire Inference with LLMs.** QSTN provides a fully modular pipeline with different ways to present the questionnaire, prompt perturbations and to choose a response generation method, with automatic parsing. Both local and remote inference are supported.

the robustness of LLM responses to such prompts. The closed-ended responses of an LLM can vary strongly from its open-ended responses (Röttger et al., 2024; Wang et al., 2024), LLM responses can be biased towards specific survey response options (Tjuatja et al., 2024; Rupperecht et al., 2025), and downstream performance is strongly affected by small changes in the questionnaire configuration (Cummins, 2025; Ahnert et al., 2025).

To address and investigate some of these concerns, **we present QSTN** (pronounced “Question”) - a Python framework designed to **facilitate the execution of questionnaire-style experiments with LLMs**. QSTN simplifies the process of creating robust variations of question prompts and answer generation methods, thereby facilitating reproducibility and the analysis of the reliability of LLM-based

questionnaire research. QSTN provides a complete, modular pipeline, as depicted in Figure 1, for creating the questionnaire presentation, adjusting various parts of the prompt with perturbations, choosing the response generation method, performing inference, and finally, parsing the generated text. We evaluate our framework on more than 40 million survey responses and find that the controlled variation of the experiment pipeline can increase the alignment of generated responses with human survey answers and that the responses can be obtained for a fraction of the compute cost.

🔗 **Python package under MIT license:**

<https://github.com/dess-mannheim/QSTN>

🖥️ **Live GUI:** [https://hf.co/spaces/qstn/qstn\\_gui](https://hf.co/spaces/qstn/qstn_gui) or run it locally by cloning the Git repository

📺 **Video:** <https://youtu.be/uM5Q-Qmm6nQ>

## 2 Core Features

QSTN was developed with three objectives in mind: First, it enables *robust evaluation* of and with LLMs, addressing prompt sensitivity (Tjuatja et al., 2024; Dominguez-Olmedo et al., 2024). QSTN is engineered to address this challenge directly through a highly modular and configurable design. Each part of the pipeline can be exchanged independently from the other parts.

Second, QSTN is designed to be *efficient*, so it can be used in large-scale studies. For experiments with multiple prompt variations and/or personas, we automatically utilize prefix caching and batching for local inference in vLLM (Kwon et al., 2023), and asynchronous calling with the AsyncOpenAI API (OpenAI, 2023).

Finally, QSTN is designed to be as *easy to use as possible*. Since we maintain the common prompt format of the system prompt and user prompt, adapting a project to QSTN is seamless. The package offers a complete pipeline from prompt creation and inference to parsing, which can be done in only three function calls to the package. Integration with existing vLLM and OpenAI packages is straightforward.

QSTN’s *core strength lies in its ability to systematically and easily control and vary the setup* of questionnaire-like prompting experiments. The following aspects of the experiments can be exchanged and varied by simply switching out one module for another.

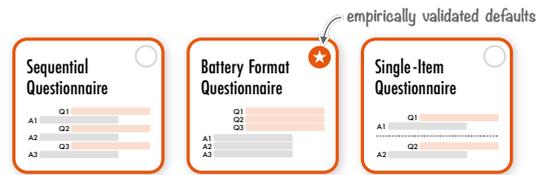


Figure 2: QSTN Questionnaire Presentation Modes

### 2.1 ■ Questionnaire Presentation

QSTN supports three distinct questionnaire presentation modes, as shown in Figure 2:

- **Sequential:** Each question is asked in the same conversation context in multiple, sequential chat calls.
- **Battery:** All questions are asked in one battery and the model is expected to answer all questions in one response in the same context.
- **Single-item:** Each question is asked in a new context, with the LLM not being aware of the previous questions and answers.

Questionnaire presentation is a fundamental decision to make when using LLMs with questionnaire-like prompts. For example, if we want to annotate data, is it better to give all annotation questions in the same prompt, or should each question be asked in a new context? There is evidence that keeping multiple tasks in the prompts can improve variety for creative writing (Zhang et al., 2025) and improve performance for classification tasks in moral foundations (Chen et al., 2025). LLMs are also able to perform multiple tasks of different kinds in one battery (Son et al., 2024), which can save computing time.

### 2.2 ■ Prompt Perturbation

Previous studies found that LLMs synthetic survey responses are highly sensitive to prompt perturbations and exhibit biases, such as token biases, recency bias, or A-bias (Pezeshkpour and Hruschka, 2024; Li and Gao, 2025; Rupprecht et al., 2025; Dominguez-Olmedo et al., 2024; Röttger et al., 2024). QSTN can automatically randomize or reverse both the order of the questions within the survey and the order of answer options for each question to identify and mitigate these biases. This ensures that high performance is robust and independent of ordering. Previous research has found that LLMs can be sensitive to small changes in prompt format (He et al., 2024; Sclar et al., 2024). QSTN allows users to define custom answer label schemas (e.g., A/B/C, 1/2/3, i/ii/iii), enabling rigorous testing of a model’s robustness to superficial

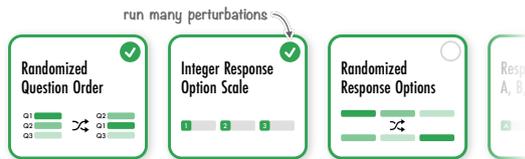


Figure 3: QSTN Supported Prompt Perturbations

formatting changes. QSTN can perform the following Answer Option Perturbations, which are shown in Figure 3:

- **Reversed Response Order:** The order of answer options is reversed (e.g., a scale from ‘1: Very important’ to ‘5: Not important’ becomes ‘1: Not important’ to ‘5: Very important’).
- **Missing Refusal Option:** The “Don’t know” or refusal option is removed from the list of choices.
- **Odd/Even Scale Transformation:** For scales with an even number of options, a semantically appropriate middle category is added, transforming it into an odd-numbered scale (e.g., by adding ‘Neutral’). Conversely, for odd-numbered scales, we remove the middle category to create an even scale and adjust the integer label accordingly.

In addition, QSTN can perform the following Question Perturbations:

- **Typographical Errors:** three types of typos can be introduced: *Key Typo* (replacing a character with a random one), *Letter Swap* (swapping two adjacent characters in a random word), and *Keyboard Typo* (replacing a character with an adjacent one on a QWERTY keyboard).
- **Semantic Variations:** Additional semantic variations can be introduced while preserving the original meaning: first, by *Synonym Replacement*, where a variable amount of words in the original question are replaced with synonyms. Second, through *Paraphrasing* the entire question is rephrased.

### 2.3 ■ Response Generation

While generative language models are designed to generate open-ended text, previous studies have implemented various approaches to constrain LLMs to closed-ended responses (e.g., Ma et al., 2024). We define **Response Generation Methods** as techniques used to elicit closed-ended responses from large language models to questionnaires (Ahnert et al., 2025). QSTN supports the following Response Generation Methods, with examples being shown in Figure 4:

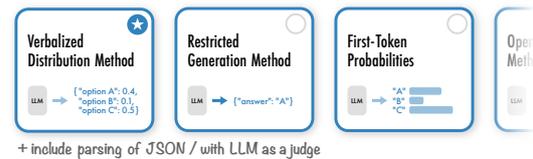


Figure 4: QSTN Response Generation Methods

- **Token Probability-Based Methods:** Extract probabilities for response options from the output token probabilities of an LLM.
- **Restricted Generation Methods:** Force the model to respond only with designated response options using formatting instructions in the prompt and (optionally) restrict the vocabulary of the LLM through *structured outputs*.
- **Open Generation Methods:** Generate open-ended responses first and then classify them in a second step.

The Restricted Generation Methods can be used to generate exactly one of the available response options—optionally in JSON format, or with reasoning—or to generate a **verbalized distribution** of probabilities for all response options, following Meister et al. (2025). All Response Generation Methods can be adjusted to, e.g., have the model generate a prefix before token probabilities are extracted. QSTN includes **suitable parsers** for all generated responses: JSON & LLM-as-a-judge.

## 3 Using QSTN

The package containing QSTN can be installed in the desired environment using pip. We support both a lightweight installation with `pip install qstn`, which only requires dependencies for API usage, and a full installation with `pip install qstn[vllm]`, which allows for local inference with vllm. QSTN is easily integrable into current workflows, requiring just a total of three function calls for the most basic functionality, and it still allows users to freely define their prompts. A minimum usage example is given in Listing 1. By simply exchanging the function in the inference step, the questionnaire presentation can be adjusted, or a different type of parser can be selected. Additionally, building on this simple example, only one more module is needed to implement controlled prompt perturbations and response generation methods.

**Non-Code User Interface** QSTN offers a User Interface to create and run inference with LLMs without having to program any Python code. The UI offers the same core functionality as the main frame-

work, allowing users to upload questionnaires, systematically alter the prompt structure, set model parameters, and run inference. While the UI generally offers the same functions as the coding package, some more advanced features, such as inferencing models directly through the vllm API, are currently not supported.

## 4 Evaluation

We evaluate QSTN primarily on the task of generating synthetic survey responses, which is a topic of growing interest. Our results demonstrate that our proposed variations significantly influence both the alignment of synthetic data with real-world responses and computational efficiency. Across all experiments, we use the following instruction-finetuned versions of the models: LLama3 1-70B (Grattafiori et al., 2024), Qwen3 4-30B (Yang et al., 2025), Phi-4-mini (Abdin et al., 2024), Gemma3 4-27B (Team et al., 2025), OLMo2 1-32B (OLMo et al., 2024), Yi1.5 6B (Young et al., 2024), and Gemini1.5 Pro (Team et al., 2024). We present new results and evaluations regarding questionnaire presentation and provide an overview of previous experiments that were conducted in or implemented into QSTN, which evaluate Prompt Perturbations and Response Generation Methods.

---

```
import qstn
import pandas as pd
from vllm import LLM

# 1. Prepare questionnaire and persona data
questionnaires = pd.read_csv("hf://datasets/qstn/ex/q.csv")
personas = pd.read_csv("hf://datasets/qstn/ex/p.csv")
prompt = (
    f"Please tell us how you feel about:\n"
    f"{qstn.utilities.placeholder.PROMPT_QUESTIONS}"
)
interviews = [
    qstn.prompt_builder.LLMPrompt(
        questionnaire_source=questionnaires,
        system_prompt=persona,
        prompt=prompt,
    ) for persona in personas.system_prompt]

# 2. Run Inference
model = LLM("Qwen/Qwen3-4B", max_model_len=5000)
results = qstn.survey_manager.conduct_survey_single_item(
    model, interviews, max_tokens=500
)

# 3. Parse Results
parsed_results = qstn.parser.raw_responses(results)
```

---

**Listing 1: Minimum usage example of QSTN.** QSTN can be easily integrated into existing projects, requiring just three function calls to operate. Users familiar with vllm or the OpenAI API can use the same Model/Client calls and arguments. In this example reasoning and the generated response are automatically parsed.

### 4.1 ■ Questionnaire Presentation

We start by demonstrating that the presentation of the questionnaire significantly impacts the subpopulation alignment of generated responses with real answers. Furthermore, selecting the optimal method results in savings for both token usage and GPU time. We test three fundamentally different presentations, as described in Section 2.1.

We base our experiments on Bisbee et al. (2024), where respondents of the ANES survey are instructed to consider a certain group and to indicate the degree to which they experience warm (positive, affectionate, etc.) or cool (negative, disdainful, etc.) feelings toward members of that group on a scale from 0 to 100. For each of the 7530 participants, we use three different seeds, which leads to a **total of 10,843,200 individual question responses** across 16 questions, 10 different models, and 3 different presentations.

We use the same prompts as in the initial study, with the addition of an instruction on how to format the output to align with the response generation method. Our full prompts can be seen in Table 7 in the Appendix. Respondents were stratified into subpopulations based on the intersection of gender, race, and ideology (see Appendix Table 6 for full subpopulation attributes). We measure individual alignment via Mean Absolute Error and subpopulation distributional alignment via Wasserstein distance; results are displayed in Table 1. To quantify the effects of questionnaire presentation, we fitted Ordinary Least Squares and Weighted Least Squares models for MAE and Wasserstein distance, respectively. Both models include interaction terms between presentation and model, as well as fixed effects for iteration seeds. The ■ single-item presentation and Llama-3.3-70B-Instruct serve as the reference categories.

We find that questionnaire presentation has a substantial impact on distributional alignment, whereas the effects on individual-level accuracy, while statistically significant, are practically marginal. For the reference model, the ■ battery presentation yields the strongest improvement in subpopulation alignment ( $\beta_{WD} = -1.17, p < 0.01$ ), representing an approximate 8% better alignment than with the ■ single-item presentation. This effect is consistent across the large models we tested, as the interaction effect for both Qwen-30B and Gemma-27B was not statistically significant. However, for smaller models, the effect is highly

| questionnaire presentation  | Mean Absolute Error ↓ |                     |                     | Wasserstein distance ↓ |                     |                     |
|-----------------------------|-----------------------|---------------------|---------------------|------------------------|---------------------|---------------------|
|                             | ■ sequential          | ■ battery           | ■ single-item       | ■ sequential           | ■ battery           | ■ single-item       |
| gemma-3-4b-it               | 20.96 ± 0.02          | 21.92 ± 0.01        | <b>19.94 ± 0.02</b> | 16.48 ± 0.01           | 17.62 ± 0.02        | <b>16.39 ± 0.01</b> |
| gemma-3-12b-it              | 18.26 ± 0.01          | <b>18.07 ± 0.02</b> | 19.11 ± 0.01        | 14.53 ± 0.02           | <b>13.44 ± 0.01</b> | 16.44 ± 0.01        |
| gemma-3-27b-it              | <b>17.59 ± 0.01</b>   | 17.90 ± 0.01        | 18.01 ± 0.01        | <b>14.00 ± 0.01</b>    | 14.26 ± 0.01        | 15.17 ± 0.00        |
| Llama-3.2-1B-Instruct       | 30.89 ± 0.25          | <b>30.22 ± 0.07</b> | 35.69 ± 0.12        | 18.66 ± 0.33           | <b>18.15 ± 0.12</b> | 27.52 ± 0.17        |
| Llama-3.2-3B-Instruct       | 24.20 ± 0.10          | <b>22.98 ± 0.04</b> | 24.32 ± 0.06        | <b>13.14 ± 0.11</b>    | 13.50 ± 0.03        | 15.88 ± 0.07        |
| Llama-3.1-8B-Instruct       | 21.01 ± 0.04          | 20.88 ± 0.02        | <b>20.87 ± 0.04</b> | 13.62 ± 0.02           | <b>12.90 ± 0.02</b> | 14.11 ± 0.04        |
| Llama-3.3-70B-Instruct      | 18.23 ± 0.00          | <b>17.67 ± 0.00</b> | 17.87 ± 0.01        | 14.18 ± 0.00           | <b>13.56 ± 0.01</b> | 14.73 ± 0.01        |
| Phi-4-mini-instruct         | 20.98 ± 0.03          | <b>19.72 ± 0.01</b> | 21.23 ± 0.03        | <b>11.69 ± 0.06</b>    | 12.21 ± 0.02        | 14.56 ± 0.05        |
| Qwen3-4B-Instruct-2507      | <b>19.29 ± 0.02</b>   | 20.34 ± 0.01        | 20.05 ± 0.02        | <b>13.75 ± 0.02</b>    | 15.59 ± 0.01        | 15.18 ± 0.01        |
| Qwen3-30B-A3B-Instruct-2507 | 17.68 ± 0.02          | <b>17.67 ± 0.01</b> | 18.29 ± 0.01        | 13.88 ± 0.02           | <b>13.68 ± 0.01</b> | 15.21 ± 0.02        |

Table 1: **Individual and subpopulation alignment based on ■ questionnaire presentation.** Mean absolute error for each individual response and weighted mean Wasserstein distance across the subpopulations. Wasserstein distance significantly improves with sequential and battery presentation for most models, compared to single-item.

architecture-dependent: Phi-4-mini achieves the best overall alignment in our experiment using the ■ sequential presentation, whereas gemma-3-4b achieves the best alignment with ■ single-item presentation.

Considering the large differences in tokens and compute time between the presentation methods (shown in Table 2), **we recommend the ■ battery presentation as the default for future questionnaire-based experiments with large persona prompts.** However, thorough tests should be conducted to ensure that performance is comparable to other presentations for the specific model and task at hand. QSTN makes these validation experiments accessible by requiring just a single method change in the pipeline.

## 4.2 ■ Prompt Perturbation

In previous research (Rupprecht et al., 2025), we found a consistent recency bias in all nine models tested, favoring the same answer option when placed at the end of the options list instead of the beginning. This effect was substantial, with the selection frequency of the semantically same option increasing by more than 20 times for Llama-3.1-8B when moved to the last position, while all other configurations, such as question and prompt phrasing, were kept constant.

All models facing prompt perturbations showed some level of non-robust responses, whereas larger models such as Llama-3.3-70B and Gemini-1.5-Pro respond more robustly. The magnitude of the effect of perturbations (e.g., on the answer option or the question phrasing) on response robustness mainly depends on the type of pertur-

| Presentation  | Calls | Input T. | Output T. | Inference Time |
|---------------|-------|----------|-----------|----------------|
| ■ sequential  | 16    | 8216     | 288       | 09:29:05       |
| ■ battery     | 1     | 723      | 142       | 01:34:45       |
| ■ single-item | 16    | 4288     | 288       | 03:22:23       |

Table 2: **API Calls, Tokens and inference time of different ■ questionnaire presentations.** We report the number of API calls, tokens and inference time for the largest model Llama-3.3-70B-Instruct. The tokens are calculated on one persona and the time is measured by a whole run of 7530 personas with 3 seeds. All experiments have been conducted with vllm on two 2 NVIDIA H100 GPUs (tensor-parallel).

bation applied. We identified that some of the ■ Answer Option Perturbations and ■ Question Perturbations have a larger impact on response robustness than others (see Table 3). Reversing the answer options or introducing typos or paraphrasing the questions is more harmful to robustness than swapping characters within a word or removing the refusal category. In addition, we found that 67% and 89% of models select the middle category significantly more often when a 5- or 11-point Likert scale is provided, respectively.

These findings underline the importance of robustness checks, e.g., through prompt perturbations. QSTN allows the user to apply various perturbations automatically to any questionnaire presented and thus assess the response robustness of the LLM.

## 4.3 ■ Response Generation Methods

To investigate the impact of Response Generation Methods on generated questionnaire responses, we **predict survey responses to questions of political attitudes** in the American National Election

| Model           | ■ Answer Options |             |             | ■ Question Perturbations |             |             |             |             |
|-----------------|------------------|-------------|-------------|--------------------------|-------------|-------------|-------------|-------------|
|                 | (1)              | (2)         | (3)         | (4)                      | (5)         | (6)         | (7)         | (8)         |
| Llama-3.3-70B   | 0.50             | 0.73        | <b>0.60</b> | 0.52                     | <b>0.76</b> | 0.58        | 0.58        | <b>0.66</b> |
| Llama-3.1-8B    | 0.08             | 0.39        | 0.27        | 0.32                     | 0.31        | 0.23        | 0.32        | 0.16        |
| Llama-3.2-3B    | 0.10             | 0.11        | 0.16        | 0.10                     | 0.16        | 0.18        | 0.23        | 0.10        |
| Llama-3.2-1B    | 0.00             | 0.11        | 0.03        | 0.05                     | 0.11        | 0.00        | 0.13        | 0.02        |
| Gemini-1.5-Pro  | <b>0.69</b>      | 0.76        | 0.55        | <b>0.68</b>              | 0.73        | <b>0.66</b> | 0.60        | 0.55        |
| Phi-3.5-mini    | 0.53             | <b>0.81</b> | 0.45        | 0.50                     | 0.61        | 0.47        | <b>0.71</b> | 0.53        |
| Mistral-7B-v0.3 | 0.68             | <b>0.81</b> | 0.53        | 0.58                     | 0.65        | 0.60        | <b>0.71</b> | 0.53        |
| Qwen-2.5-7B     | 0.32             | 0.48        | 0.45        | 0.48                     | 0.65        | 0.45        | 0.55        | 0.44        |
| Yi-1.5-6B       | 0.47             | 0.68        | 0.55        | 0.50                     | 0.50        | 0.45        | 0.65        | 0.29        |

Table 3: **Impact of ■ Answer Option and ■ Question Perturbations on the Response Robustness of different LLMs (↑)**. Share of fully robust responses per model. Bold indicates the highest robustness score for that perturbation type. Perturbation Keys: (1) Reversed Answer Options, (2) Missing Refusal, (3) Even Scale, (4) Key Typos, (5) Letter Swap, (6) Keyboard Typos, (7) Synonyms, (8) Paraphrase

Study (ANES, 2016), the German Longitudinal Election Study (GLES, 2017, 2025), and the American Trends Panel (ATP, 2021). We thereby partially replicate the studies by Argyle et al. (2023), von der Heyde et al. (2025), and Santurkar et al. (2023), while extending them to include additional Response Generation Methods. We compare 8 Response Generation Methods on 10 open-weight LLMs, including reasoning models. For robustness, we include 4 prompt variations, 3 random seeds for temperature-scaled decoding, as well as greedy decoding. Overall, **we simulate 32 mio. survey responses with QSTN**, and evaluate their alignment with human survey responses on individual and subpopulations levels. For subpopulation-level alignment, we split the set of respondents into subpopulations by considering all unique values of all persona attributes that were included in the studies we replicate, e.g., women & men, people from different states, etc. We report the subpopulation-level alignment on categorical response distributions using total variation distance (see also Meister et al., 2025; Baan et al., 2022).

Table 4 shows selected OLS regression coefficients for subpopulation-level alignment. We find that the Verbalized Distribution Method yields significant improvements on most datasets. In combination with the individual-level alignment results presented in Ahnert et al. (2025), we conclude that: (i) the **choice of Survey Response Generation Method should be well-justified** for *in-silico* surveys, since we find significant differences between these methods. (ii) We **do not recommend the use of Token Probability-Based Methods**, as they

| Response Generation Method | ANES 2016     | GLES 2017     | GLES 2025     | ATP 2021      |
|----------------------------|---------------|---------------|---------------|---------------|
| Intercept                  | .374*         | .312*         | .288*         | .503*         |
| ■ First-Token Prob.        | -.003         | .147*         | .194*         | -.049*        |
| ■ Verbalized Distrib.      | <b>-.074*</b> | <b>-.057*</b> | -.013         | <b>-.168*</b> |
| ■ Open-Ended Distrib.      | -.006         | -.052*        | <b>-.037*</b> | -.082*        |

Table 4: **Impact of ■ Response Generation Methods on Subpopulation-Level Alignment (↓)**. OLS regression coefficients by dataset with total variation distance (↓) as the dependent variable and Survey Response Generation Method, prompt perturbation, and LLM as independent variables. We show coefficients for selected Response Generation Methods (Reference: Restricted Choice)—see Appendix B for all coefficients and more details on OLS model choice. **The Verbalized Distribution Method leads to significant improvements.** \* $p < 0.05$  (Benjamini–Hochberg corrected)

generate misaligned survey responses. (iii) For predicting closed-ended survey responses, we suggest to **consider Restricted Generation Methods first**, as they consistently show significant improvement over other methods while also being more computationally efficient than Open Generation Methods.

## 5 Related Work

Due to the importance of controlled prompt perturbation, a number of frameworks have started to address this issue. In general, QSTN supports controlled variation and combines it with the pipeline to allow for automatic parsing of all prompt variations. Additionally, as QSTN allows for modular prompts, these frameworks can be used in conjunction with it. PromptSuite (Habba et al., 2025) focuses on prompt perturbation through paraphrasing and formatting. PromptSource (Bach et al., 2022) is a framework for making and sharing different types of natural language prompts. Prompt-Agnostic Fine-Tuning (PAFT) (Wei et al., 2025) varies prompts in the fine-tuning process rather than during inference.

There are also frameworks that model the entire pipeline of LLM experiments, similar to QSTN. Unitxt (Bandel et al., 2024) is an open-source Python framework for data processing pipelines. While powerful, it requires users to understand the Unitxt operator language, which can add cognitive overhead. The EDSL framework (Horton and Horton, 2024) can be used to run surveys with LLMs, but it does not provide full freedom over the exact system prompt or prompt and the Response Generation Method.

## 6 Conclusion

We introduce QSTN, a Python framework designed to make LLM inference with questionnaires more robust. Our evaluation demonstrates that by enabling controlled variations in the generation process, QSTN can significantly improve the alignment of generated responses with human answers while reducing inference costs. A core feature of QSTN is its modularity, allowing researchers to easily vary their experimental setup with only minimal additional coding effort. The framework is broadly applicable to tasks such as data annotation, synthetic data generation, persona studies, and the analysis of LLM behavior itself.

## Limitations

Currently, our evaluation is primarily focused on the creation of synthetic survey responses. We hope that by releasing QSTN to the open-source community, more robust experiments can be conducted in other application domains. While we support a variety of different Response Generation Methods and parsing options, we currently do not support every type of structured output; for example, we do not support output that is guided by a regex pattern or context free grammar. As such, not every type of experiment can currently be conducted in QSTN. We hope that by making the project open-source, we will be able to support more ways to conduct experiments. Additionally, while we plan to add support for non-instruct models, they are currently not supported.

## References

- Marah Abidin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, and 1 others. 2024. [Phi-4 technical report](#). *arXiv preprint arXiv:2412.08905*.
- Georg Ahnert, Anna-Carolina Haensch, Barbara Plank, and Markus Strohmaier. 2025. [Survey response generation: Generating closed-ended survey responses in-silico with large language models](#). *arXiv preprint arXiv:2510.11586*.
- ANES. 2016. [2016 Time Series Study](#).
- Lisa P. Argyle, Ethan C. Busby, Nancy Fulda, Joshua R. Gubler, Christopher Rytting, and David Wingate. 2023. [Out of one, many: Using language models to simulate human samples](#). *Political Analysis*, 31(3):337–351.
- ATP. 2021. [The American Trends Panel](#).
- Joris Baan, Wilker Aziz, Barbara Plank, and Raquel Fernandez. 2022. [Stop measuring calibration when humans disagree](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1892–1915, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gungjan Chhablani, Han Wang, Jason Alan Fries, and 8 others. 2022. [Promptsources: An integrated development environment and repository for natural language prompts](#). *Preprint*, arXiv:2202.01279.
- Elron Bandel, Yotam Perlit, Elad Venezian, Roni Friedman, Ofir Arviv, Matan Orbach, Shachar Don-Yehiya, Dafna Sheinwald, Ariel Gera, Leshem Choshen, Michal Shmueli-Scheuer, and Yoav Katz. 2024. [Unitxt: Flexible, shareable and reusable data preparation and evaluation for generative AI](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: System Demonstrations)*, pages 207–215, Mexico City, Mexico. Association for Computational Linguistics.
- James Bisbee, Joshua D. Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M. Larson. 2024. [Synthetic replacements for human survey data? the perils of large language models](#). *Political Analysis*, 32(4):401–416.
- Ziyu Chen, Junfei Sun, Chenxi Li, Tuan Dung Nguyen, Jing Yao, Xiaoyuan Yi, Xing Xie, Chenhao Tan, and Lexing Xie. 2025. [MoVa: Towards generalizable classification of human morals and values](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 33204–33248, Suzhou, China. Association for Computational Linguistics.
- Jamie Cummins. 2025. [The threat of analytic flexibility in using large language models to simulate human data: A call to attention](#). *arXiv preprint arXiv:2509.13397*.
- Ricardo Dominguez-Olmedo, Moritz Hardt, and Celestine Mendler-Düner. 2024. [Questioning the survey responses of large language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 45850–45878. Curran Associates, Inc.
- GLSES. 2017. [GLSES 2017 Post-Election Cross Section](#).
- GLSES. 2025. [GLSES 2025 Post-Election Cross Section](#).
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.

- Eliya Habba, Noam Dahan, Gili Lior, and Gabriel Stanovsky. 2025. [PromptSuite: A task-agnostic framework for multi-prompt generation](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 254–263, Suzhou, China. Association for Computational Linguistics.
- Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X Wang, and Sadid Hasan. 2024. [Does prompt formatting have any impact on llm performance?](#) *arXiv preprint arXiv:2411.10541*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- John Horton and Robin Horton. 2024. [Edsl: Expected parrot domain specific language for ai powered social science](#). Whitepaper, Expected Parrot.
- Jennifer Hu, Sammy Floyd, Olessia Jouravlev, Evelina Fedorenko, and Edward Gibson. 2023. [A fine-grained comparison of pragmatic language understanding in humans and language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4194–4213, Toronto, Canada. Association for Computational Linguistics.
- Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2024. [PersonaLLM: Investigating the ability of large language models to express personality traits](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3605–3627, Mexico City, Mexico. Association for Computational Linguistics.
- Yeeun Kim, Youngrok Choi, Eunhyung Choi, JinHwan Choi, Hai Jin Park, and Wonseok Hwang. 2024. [Developing a pragmatic benchmark for assessing Korean legal language understanding in large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 5573–5595, Miami, Florida, USA. Association for Computational Linguistics.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). In *Proceedings of the 29th Symposium on Operating Systems Principles, SOSP '23*, page 611–626, New York, NY, USA. Association for Computing Machinery.
- Ruizhe Li and Yanjun Gao. 2025. [Anchored answers: Unravelling positional bias in GPT-2's multiple-choice questions](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 2439–2465, Vienna, Austria. Association for Computational Linguistics.
- Bolei Ma, Xinpeng Wang, Tiancheng Hu, Anna-Carolina Haensch, Michael Hedderich, Barbara Plank, and Frauke Kreuter. 2024. [The potential and challenges of evaluating attitudes, opinions, and values in large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 8783–8805.
- Nicole Meister, Carlos Guestrin, and Tatsunori Hashimoto. 2025. [Benchmarking distributional alignment of large language models](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 24–49, Albuquerque, New Mexico. Association for Computational Linguistics.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, and 1 others. 2024. [2 olmo 2 furious](#). *arXiv preprint arXiv:2501.00656*.
- OpenAI. 2023. [OpenAI Python Library](#).
- Max Pellert, Clemens M. Lechner, Claudia Wagner, Beatrice Rammstedt, and Markus Strohmaier. 2024. [AI psychometrics: Assessing the psychological profiles of large language models through psychometric inventories](#). *Perspectives on Psychological Science*, 19(5):808–826. Epub 2024 Jan 2.
- Pouya Pezeshkpour and Estevam Hruschka. 2024. [Large language models sensitivity to the order of options in multiple-choice questions](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2006–2017, Mexico City, Mexico. Association for Computational Linguistics.
- Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Kirk, Hinrich Schuetze, and Dirk Hovy. 2024. [Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15295–15311, Bangkok, Thailand. Association for Computational Linguistics.
- David Rozado. 2024. [The political preferences of llms](#). *Preprint*, arXiv:2402.01789.
- Jens Rupperecht, Georg Ahnert, and Markus Strohmaier. 2025. [Prompt perturbations reveal human-like biases in llm survey responses](#). *arXiv preprint arXiv:2507.07188*.
- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. [Whose opinions do language models reflect?](#) In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org.
- Ankit Satpute, Noah Gießing, André Greiner-Petter, Moritz Schubotz, Olaf Teschke, Akiko Aizawa, and Bela Gipp. 2024. [Can llms master math? investigating large language models on math stack exchange](#).

- In *Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval*, pages 2316–2320.
- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. [Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting](#). In *The Twelfth International Conference on Learning Representations*.
- 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](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5263–5281, Mexico City, Mexico. Association for Computational Linguistics.
- Guijin Son, SangWon Baek, Sangdae Nam, Ilgyun Jeong, and Seungone Kim. 2024. [Multi-task inference: Can large language models follow multiple instructions at once?](#) In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5606–5627, Bangkok, Thailand. Association for Computational Linguistics.
- Settaluri Sravanthi, Meet Doshi, Pavan Tankala, Rudra Murthy, Raj Dabre, and Pushpak Bhattacharyya. 2024. [PUB: A pragmatics understanding benchmark for assessing LLMs’ pragmatics capabilities](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 12075–12097, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, Amrita Bhattacharjee, Mansoor Karami, Jundong Li, Lu Cheng, and Huan Liu. 2024. [Large language models for data annotation and synthesis: A survey](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 930–957, Miami, Florida, USA. Association for Computational Linguistics.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *arXiv preprint arXiv:2403.05530*.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. [Gemma 3 technical report](#). *arXiv preprint arXiv:2503.19786*.
- Lindia Tjauatja, Valerie Chen, Tongshuang Wu, Ameet Talwalkar, and Graham Neubig. 2024. [Do LLMs exhibit human-like response biases? a case study in survey design](#). *Transactions of the Association for Computational Linguistics*, 12:1011–1026.
- Leah von der Heyde, Anna-Carolina Haensch, and Alexander Wenz. 2025. [Vox populi, vox ai? using large language models to estimate german vote choice](#). *Social Science Computer Review*, 0(0):1–23.
- Xinpeng Wang, Bolei Ma, Chengzhi Hu, Leon Weber-Genzel, Paul Röttger, Frauke Kreuter, Dirk Hovy, and Barbara Plank. 2024. [“my answer is C”: First-token probabilities do not match text answers in instruction-tuned language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7407–7416, Bangkok, Thailand. Association for Computational Linguistics.
- Chenxing Wei, Mingwen Ou, Ying He, Yao Shu, and Fei Yu. 2025. [PAFT: Prompt-agnostic fine-tuning](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 694–717, Suzhou, China. Association for Computational Linguistics.
- Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. 2023. [Cmath: Can your language model pass chinese elementary school math test?](#) *arXiv preprint arXiv:2306.16636*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. [Qwen3 technical report](#). *arXiv preprint arXiv:2505.09388*.
- Haoran Ye, Jing Jin, Yuhang Xie, Xin Zhang, and Guojie Song. 2025. [Large language model psychometrics: A systematic review of evaluation, validation, and enhancement](#). *Preprint*, arXiv:2505.08245.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, and 1 others. 2024. [Yi: Open foundation models by 01. ai](#). *arXiv preprint arXiv:2403.04652*.
- Jiayi Zhang, Simon Yu, Derek Chong, Anthony Sicilia, Michael R Tomz, Christopher D Manning, and Weiyan Shi. 2025. [Verbalized sampling: How to mitigate mode collapse and unlock llm diversity](#). *arXiv preprint arXiv:2510.01171*.

## A Questionnaire Presentation

As another measure of individual evaluation of predictions, we calculate the Pearson correlation between predictions and the ground truth and present it in Table 5. Similarly to the Mean Absolute Error, we see little difference in the performance of the different questionnaire presentations.

We show all attributes we considered for the subpopulation analysis for Wasserstein distance in Table 6. The full regression results for both MAE

| questionnaire presentation  | sequential         | battery            | single-item        |
|-----------------------------|--------------------|--------------------|--------------------|
| gemma-3-4b-it               | <b>0.59 ± 0.00</b> | 0.55 ± 0.00        | 0.57 ± 0.00        |
| gemma-3-12b-it              | 0.62 ± 0.00        | 0.62 ± 0.00        | <b>0.64 ± 0.00</b> |
| gemma-3-27b-it              | <b>0.62 ± 0.00</b> | 0.61 ± 0.00        | 0.61 ± 0.00        |
| Llama-3.2-1B-Instruct       | <b>0.25 ± 0.01</b> | 0.18 ± 0.00        | 0.10 ± 0.00        |
| Llama-3.2-3B-Instruct       | 0.51 ± 0.00        | 0.49 ± 0.00        | <b>0.52 ± 0.00</b> |
| Llama-3.1-8B-Instruct       | 0.56 ± 0.00        | <b>0.57 ± 0.00</b> | 0.56 ± 0.00        |
| Llama-3.3-70B-Instruct      | <b>0.64 ± 0.00</b> | <b>0.64 ± 0.00</b> | <b>0.64 ± 0.00</b> |
| Phi-4-mini-instruct         | 0.48 ± 0.00        | 0.49 ± 0.00        | <b>0.52 ± 0.00</b> |
| Qwen3-4B-Instruct-2507      | <b>0.60 ± 0.00</b> | 0.55 ± 0.00        | 0.59 ± 0.00        |
| Qwen3-30B-A3B-Instruct-2507 | <b>0.62 ± 0.00</b> | <b>0.62 ± 0.00</b> | 0.59 ± 0.00        |

Table 5: **Mean and Standard Deviation of Pearson Correlation between Prediction and Ground Truth.** Similar to Mean Absolute Error, individual alignment measured with Pearson Correlation shows little difference between different questionnaire presentations.

and Wasserstein Distance can be seen in 8. We report the coefficients and the Benjamini-Hochberg corrected p-values. Additionally, we want to determine if the questionnaire presentation has different effects on different questions. For this, we fit an additional Weighted Least Squares regression on all subpopulations based on the full interaction between the questionnaire presentation, the model, and the specific interview question. We set ■ single-item, the biggest model Llama-3.3-70B-Instruct and the first question as the reference categories, as for this question the LLM has no answers for the other questions in context regardless of the questionnaire presentation.

All questions show improvements, and a subset of five questions shows statistically significant improvement ( $p < 0.05$ ) when using ■ battery presentation instead of ■ single-item presentation. The largest improvement is in the question about feelings towards the group of Gays and Lesbians ( $\beta = -3.82, p < 0.01$ ) when using ■ battery presentation. Figure 5 visually confirms this: when previous questions and answers are included in the context, the model’s response distribution aligns much more closely with the ground truth, exhibiting a similar tendency toward neutral answers. The other significant questions concern the groups of White Americans, Asian Americans, Christians, and Liberals.

We use the same prompt as that used in Bisbee et al. (2024), as displayed in Table 7. We adjust the output instructions to fit our choice response generation method and add all questions as instructions in the ■ battery presentation. For all models we use the default hyperparameter settings.

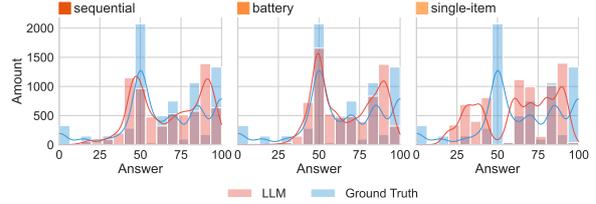


Figure 5: **Answer Distributions.** Predictions and Ground Truth Distributions across the whole population compared for Llama-3.3-70B-Instruct and the question “How do you feel towards Gays and Lesbians?”. We can see a clear shift towards the middle for this question, when models are given context of the previous questions and answers, which aligns more closely with human answers.

| Attribute       | Category               |
|-----------------|------------------------|
| <b>Race</b>     | Non-Hispanic White     |
|                 | Non-Hispanic Black     |
|                 | Hispanic               |
| <b>Gender</b>   | Male                   |
|                 | Female                 |
| <b>Ideology</b> | Extremely Liberal      |
|                 | Liberal                |
|                 | Slightly Liberal       |
|                 | Moderate               |
|                 | Slightly Conservative  |
|                 | Conservative           |
|                 | Extremely Conservative |

Table 6: **Subpopulations:** We consider these subpopulations for analysis. We have the same subpopulations as the initial study by Bisbee et al. (2024).

## B Response Generation OLS Regressions

We obtain the subpopulation-level alignment for each simulation specification and subpopulation, as described in Section 4.3. To identify significant differences in survey response alignment between the response generation methods, we fit the following OLS regression model separately on each dataset (see Table 9): We use the per-subpopulation total variation distance ( $\downarrow$ ) as the dependent variable and Survey Response Generation Method (reference: Restricted Choice), LLM (reference: Llama 8B), and prompt perturbation (reference: Full Text response options) as independent variables. We use cluster-robust SEs, clustering by seed  $\times$  decoding strategy, which allows for arbitrary correlation and heteroskedasticity within clusters while assuming independence across clusters. This appropriately

| Component                             | Content                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
|---------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>System Prompt</b>                  | "It is {year}. You are a {age} year-old, {marst}, {raceth} {gender} with {education} making {income} per year, living in the United States. You are {ideo}, {regis} {PID} who {interest} pays attention to what's going on in government and politics."                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| <b>User Prompt</b>                    | "Provide responses from this person's perspective. Use only knowledge about politics that they would have. The following questions ask about individuals' feelings toward different groups. Responses should be given on a scale from 0 (meaning cold feelings) to 100 (meaning warm feelings). Ratings between 50 degrees and 100 degrees mean that you feel favorable and warm toward the group. Ratings between 0 degrees and 50 degrees mean that you don't feel favorable toward the group and that you don't care too much for that group. You would rate the group at the 50 degree mark if you don't feel particularly warm or cold toward the group.<br>Output Instructions<br>Questions"                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| <b>Output Instruction</b>             | You only respond in the following JSON format:<br><pre>```json {   "temperature": &lt;temperature&gt; } ```</pre>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| <b>Output Instruction (■ battery)</b> | You only respond in the following JSON format:<br><pre>```json {   "temperature_The Democratic Party?": &lt;temperature_The Democratic Party?&gt;,   "temperature_The Republican Party?": &lt;temperature_The Republican Party?&gt;,   "temperature_Democrats?": &lt;temperature_Democrats?&gt;,   "temperature_Republicans?": &lt;temperature_Republicans?&gt;,   "temperature_Black Americans?": &lt;temperature_Black Americans?&gt;,   "temperature_White Americans?": &lt;temperature_White Americans?&gt;,   "temperature_Hispanic Americans?": &lt;temperature_Hispanic Americans?&gt;,   "temperature_Asian Americans?": &lt;temperature_Asian Americans?&gt;,   "temperature_Muslims?": &lt;temperature_Muslims?&gt;,   "temperature_Christians?": &lt;temperature_Christians?&gt;,   "temperature_Immigrants?": &lt;temperature_Immigrants?&gt;,   "temperature_Gays and Lesbians?": &lt;temperature_Gays and Lesbians?&gt;,   "temperature_Jews?": &lt;temperature_Jews?&gt;,   "temperature_Liberals?": &lt;temperature_Liberals?&gt;,   "temperature_Conservatives?": &lt;temperature_Conservatives?&gt;,   "temperature_Women?": &lt;temperature_Women?&gt; } ```</pre> |
| <b>Question</b>                       | How do you feel towards the Republican Party?                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |

Table 7: **Prompt.** We use the same prompts for ■ sequential and ■ single-item and a slightly modified output instruction for the ■ battery presentation. For ■ battery presentation we ask all questions separated by new lines.

reflects the repeated-measures structure of our evaluation. We do not include interaction terms into the OLS model to mitigate multicollinearity—all VIF values are  $< 3$ . We apply Benjamini–Hochberg correction across all reported coefficients in all datasets. Key coefficients for the Verbalized Distribution Method, as well as OLMo 32B and Qwen 32B remain significant even under Bonferroni correction, although Bonferroni is known to be overly conservative in regression settings with correlated predictors.

|                                            |                             | (1)       | (2)            |
|--------------------------------------------|-----------------------------|-----------|----------------|
|                                            |                             | MAE (OLS) | WD Score (WLS) |
| <b>Questionnaire Presentation</b>          | ■ sequential                | 0.362**   | -0.546*        |
|                                            | ■ battery                   | -0.199**  | -1.166**       |
| <b>Model</b>                               | Llama 3.1 8B                | 2.999**   | -0.617*        |
|                                            | Llama 3.2 1B                | 17.822**  | 12.796**       |
|                                            | Llama 3.2 3B                | 6.450**   | 1.152**        |
|                                            | Phi-4 Mini                  | 3.366**   | -0.163         |
|                                            | Qwen3 30B                   | 0.420**   | 0.488          |
|                                            | Qwen3 4B                    | 2.187**   | 0.454          |
|                                            | Gemma 3 12B                 | 1.245**   | 1.713**        |
|                                            | Gemma 3 27B                 | 0.142**   | 0.448          |
|                                            | Gemma 3 4B                  | 2.074**   | 1.662**        |
| <b>Interactions (Presentation × Model)</b> | ■ sequential × Llama 3.1 8B | -0.216**  | 0.060          |
|                                            | ■ battery × Llama 3.1 8B    | 0.213**   | -0.044         |
|                                            | ■ sequential × Llama 3.2 1B | -5.162**  | -8.311**       |
|                                            | ■ battery × Llama 3.2 1B    | -5.273**  | -8.203**       |
|                                            | ■ sequential × Llama 3.2 3B | -0.476**  | -2.195**       |
|                                            | ■ battery × Llama 3.2 3B    | -1.136**  | -1.211**       |
|                                            | ■ sequential × Phi-4 Mini   | -0.619**  | -2.322**       |
|                                            | ■ battery × Phi-4 Mini      | -1.318**  | -1.183**       |
|                                            | ■ sequential × Qwen3 30B    | -0.973**  | -0.784*        |
|                                            | ■ battery × Qwen3 30B       | -0.418**  | -0.372         |
|                                            | ■ sequential × Qwen3 4B     | -1.125**  | -0.885*        |
|                                            | ■ battery × Qwen3 4B        | 0.490**   | 1.574**        |
|                                            | ■ sequential × Gemma 3 12B  | -1.218**  | -1.366**       |
|                                            | ■ battery × Gemma 3 12B     | -0.843**  | -1.832**       |
|                                            | ■ sequential × Gemma 3 27B  | -0.779**  | -0.625         |
|                                            | ■ battery × Gemma 3 27B     | 0.087     | 0.250          |
|                                            | ■ sequential × Gemma 3 4B   | 0.652**   | 0.637          |
|                                            | ■ battery × Gemma 3 4B      | 2.175**   | 2.394**        |

Table 8: **Regression Results for MAE and Wasserstein Distance.** (↓) Model (1) uses OLS on Mean Absolute Error. Model (2) uses WLS on Wasserstein Distance, weighted by subpopulation count. Significance levels are based on Benjamini–Hochberg corrected p-values. We can see significant effects for both the questionnaire presentation, but also for the interaction between smaller models and the presentation. Reference categories: *Presentation*: ■ *single-item*, *Model*: *Llama-3.3-70B-Instruct*. \*  $p < 0.05$ , \*\*  $p < 0.01$

|                                   |                             | ANES 2016 | GLES 2017 | GLES 2025 | ATP 2021 |
|-----------------------------------|-----------------------------|-----------|-----------|-----------|----------|
| <b>Intercept</b>                  |                             | 0.374**   | 0.312**   | 0.288**   | 0.503**  |
| <b>Response Generation Method</b> | ■ First-Token Probabilities | -0.003    | 0.147**   | 0.194**   | -0.049*  |
|                                   | ■ First-Token Restricted    | 0.064**   | 0.220**   | 0.234**   | -0.005   |
|                                   | ■ Answer Prefix             | -0.002    | 0.047*    | 0.085**   | -0.082** |
|                                   | ■ Restricted Reasoning      | 0.017     | -0.035*   | -0.026    | -0.084** |
|                                   | ■ Verbalized Distribution   | -0.074**  | -0.057**  | -0.013    | -0.168** |
|                                   | ■ Open-Ended Classif.       | 0.026     | -0.011    | -0.027    | -0.051** |
|                                   | ■ Open-Ended Distrib.       | -0.006    | -0.052**  | -0.037*   | -0.082** |
| <b>Model</b>                      | Llama 3B                    | -0.051*   | 0.031     | 0.066**   | -0.039*  |
|                                   | Llama 70B                   | -0.052*   | -0.089**  | -0.127**  | 0.007    |
|                                   | OLMo 1B                     | -0.023    | 0.109**   | 0.114**   | 0.109**  |
|                                   | OLMo 7B                     | -0.062**  | 0.070**   | 0.077**   | -0.030   |
|                                   | OLMo 32B                    | -0.070**  | -0.073**  | -0.109**  | 0.016    |
|                                   | Qwen 8B                     | 0.016     | 0.020     | -0.050*   | 0.075**  |
|                                   | Qwen 8B with Reasoning      | -0.012    | 0.002     | -0.010    | 0.019    |
|                                   | Qwen 32B                    | -0.076**  | -0.108**  | -0.161**  | -0.036*  |
|                                   | Qwen 32B with Reasoning     | -0.056**  | -0.067**  | -0.081*   | -0.106** |
| <b>Response Option Variants</b>   | Full Text, Reversed         | 0.001     | -0.005    | 0.037     | -0.003   |
|                                   | Indexed                     | 0.010     | 0.003     | 0.000     | 0.022*   |
|                                   | Indexed, Reversed           | 0.035*    | 0.011     | 0.026     | 0.030**  |

Table 9: **Impact of ■ Response Generation Methods on Subpopulation-Level Alignment ( $\downarrow$ )**. OLS regression coefficients by dataset with total variation distance ( $\downarrow$ ) as the dependent variable and Survey Response Generation Method (reference: Restricted Choice), LLM (reference: Llama 8B), and prompt perturbation (reference: Full Text response options) as independent variables. **The Verbalized Distribution Method and larger models lead to significant improvements.** \* $p < 0.05$ , \*\* $p < 0.01$  (Benjamini–Hochberg corrected)

# A Virtual Assistant for Architectural Design in a VR Environment

Ander Salaberria<sup>1</sup>, Oier Ijurco<sup>1</sup>, Markel Ferro<sup>1</sup>, Jiayuan Wang<sup>1</sup>,  
Iñigo Vilá Muñoz<sup>1</sup>, Roberto de Ioris<sup>2</sup>, Jeremy Barnes<sup>1</sup>, Oier Lopez de Lacalle<sup>1</sup>

<sup>1</sup>HiTZ Center, University of the Basque Country (UPV/EHU), <sup>2</sup>Vection Technologies

Correspondence: [ander.salaberria@ehu.eu](mailto:ander.salaberria@ehu.eu)

## Abstract

Architectural design relies on 3D modeling procedures, generally carried out in Building Information Modeling (BIM) formats. In this setting, architects and designers collaborate on building designs, iterating over many possible versions until a final design is agreed upon. However, this iteration is complicated by the fact that any changes need to be made by manually introducing changes to the complex BIM files, which lengthens the design process and makes it difficult to quickly prototype changes.

To speed up prototyping, we propose VR-ARCH, a virtual assistant that allows users to interact with the BIM file in a virtual reality (VR) environment. This framework enables users to 1) make changes directly in the VR environment, 2) make complex queries about the BIM, and 3) combine these to perform more complex actions. All of this is done via voice commands and processed through a ReAct-based agentic system that selects appropriate tools depending on the query context.

This multi-tool approach enables real-time, contextualized interaction through natural language, allowing for a faster and more natural prototyping experience.

Our demo's video is available at this [link](#), whereas the code and data are publicly available [here](#).

## 1 Introduction

Architectural design review has traditionally relied on manual inspection of 2D drawings, a process that is often time-consuming and prone to oversight. The adoption of Building Information Modeling (BIM) has significantly improved this workflow by providing integrated 3D representations of buildings, leading to better coordination, fewer disputes, and increased stakeholder satisfaction (Leite, 2019; Fischer, 2006). BIM models contain rich, networked information that supports in-

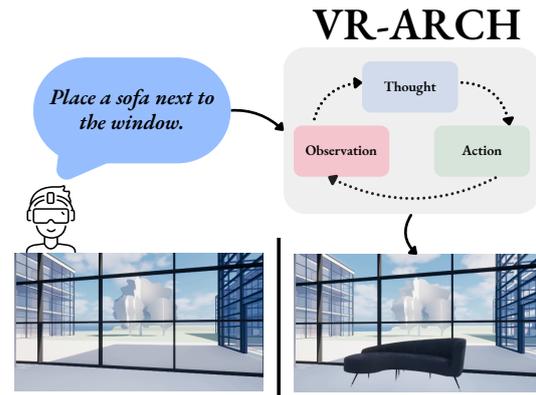


Figure 1: VR-ARCH interacts with the building and retrieves information from its BIM to assist the user in a VR environment.

formed decision-making across architectural, structural, and engineering teams (Sacks et al., 2019).

Despite these benefits, interacting with BIM environments still requires specialized knowledge and proficiency with complex software, which can hinder real-time collaboration, particularly for non-expert users (Dong et al., 2025; Leite, 2019). Existing interfaces often lack the intuitive, context-aware capabilities necessary for easily navigating or adjusting intricate architectural components in 3D spaces (Bagasi et al., 2025). In this context, natural language interaction emerges as a promising solution: an LLM-based assistant could enable users to explore, query, and modify BIM models directly through natural language, simplifying design review and enhancing collaboration.

The integration of natural language capabilities into Virtual Reality (VR) systems has long been an active area of research (Everett et al., 1997, 1999; Giunchi et al., 2024). With recent advances in LLMs (Radford and Narasimhan, 2018; Brown et al., 2020), capable of sophisticated reasoning and language understanding, there is growing interest in deploying these models within immersive environments, such as collaborative design assis-

tants (Wu et al., 2023) or as interactive virtual tutors (Ward et al., 2025).

We propose **VR-ARCH**, an LLM-based assistant that enables users to interact with and modify BIM environments through natural voice commands in a dynamic and incremental manner. Unlike previous works, VR-ARCH leverages spatial reasoning and the user’s position to interpret commands more accurately, such as referencing nearby elements or adjusting components relative to the user’s viewpoint. Through an iterative, tool-based agent architecture, the assistant can retrieve building information, execute multi-step operations, and continuously refine the environment. Figure 1 shows a user request for VR-ARCH to place a sofa next to a window in the BIM environment.

With the goal of allowing the LLM assistant to interact directly with the VR environment and BIM data, we develop a custom Unreal Engine sandbox combined with a Python API. This API provides functions to interact with assets within the building, such as hiding/revealing doors, changing a wall’s color, or rotating stairs. Similarly, to enable queries, we use a Neo4j graph database that stores BIM data and allows complex queries. To combine these two abilities, we use a ReAct (Yao et al., 2022) inspired agentic approach where a Router LLM iteratively chooses from three specialized tools: Modifier, Querier and ID Retriever. This approach can then use these subroutines to fulfill more complex user queries. Finally, we use a speech-to-text system to capture voice commands and a text-to-speech module to return spoken responses, enabling hands-free interaction with the system.

In order to approximate the usefulness of our system in for a real-world user, we perform human evaluation and furthermore develop an automatic evaluation via LLMs-as-a-judge. This evaluation finds that models are capable of accurately answering user questions and executing successful modifications within a BIM environment, especially as model size increases.

## 2 Related Work

NLP approaches in VR are often dedicated to enabling VR tutors (García-Méndez et al., 2024; Konenkov et al., 2024). Ward et al. (2025), for example, develop a VR tutor for learning Irish, while Aguirre et al. (2025) develop a VR tutor for VR health and safety training.

Another common application of NLP in VR is the development of VR assistants. Wu et al. (2023) introduce a video-grounded task-oriented dialogue dataset that captures real-world AI-assisted user scenarios in VR, while Prange et al. (2017) develop a multimodal dialogue system to help doctors make decisions about patient therapy.

Finally, there are also a few efforts to develop avatar-based VR chatbots, either via spoken dialogue (Yamazaki et al., 2023) or sign language (Quandt et al., 2022).

While these VR tutors and assistants represent an interesting step in integrating natural language in VR, they are not actually able to make any changes to the VR environment itself and generally deal with hard-coded queries, rather than performing them on the fly.

### **BIM Manipulation with LLM Integration.**

Most previous work on manipulating BIM files with LLMs has focused on querying the BIM for information (Zheng and Fischer, 2023; Li et al., 2025). Recent work has also explored ways of generating 3D building models from natural language (Du et al., 2024). Editing or prototyping changes to the BIM file, however, has not been widely explored. While Jang and Lee (2024) propose a pipeline that converts BIM to XML to facilitate LLM changes and Fernandes et al. (2024) explores conversational manipulation of building information, these approaches do not allow a user to make changes that are directly visible in a VR environment, grounded by the user’s spatial relationships. Furthermore, they currently struggle with composite queries, where several changes must be implemented.

**Knowledge Graphs for BIM.** Knowledge graphs have emerged as a structured approach to representing BIM data, enabling complex querying capabilities. Some prior work demonstrate the use of graph databases, particularly Neo4j and the Cypher query language, to represent and query building information (Ozsoy et al., 2025; Zhu et al., 2024). As these approaches allow for complex relational queries across elements in the building, we make use of a similar approach in VR-ARCH.

**Instruction-Following Multimodal Agents.** Recent advances have demonstrated that modern agents can effectively follow instructions in complex multimodal scenarios. A prevalent paradigm in recent years involves leveraging external APIs to

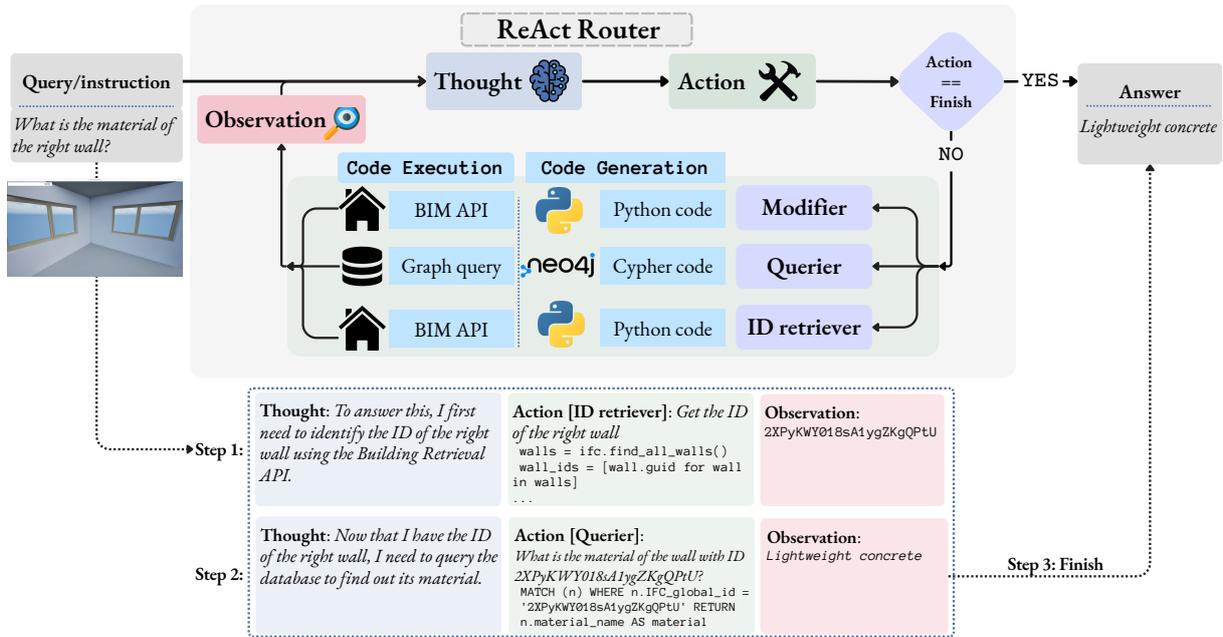


Figure 2: VR-ARCH pipeline. Given the scene’s configuration and an instruction, the ReAct router will decide which tool to use in order to modify, query or retrieve information from the building. The code for the chosen tool is then created and executed, and the ReAct Router will decide what to do in the next iteration until the initial query is fulfilled. Finally, it will generate a response to the initial user query.

extend the agent’s action space (Shen et al., 2023). These APIs are often integrated with external models to handle more sophisticated tasks (Yang et al., 2023). Our work is closely related to systems such as ViperGPT (Surís et al., 2023) and VisProg (Gupta and Kembhavi, 2023), which frame visual question answering as a complex process requiring perception, reasoning and structured tool use. These systems typically utilize LLMs to generate executable procedures that coordinate multiple components. Our approach employs a ReAct-style routing mechanism (Yao et al., 2022) that selects among a variety of tools.

### 3 System Architecture

VR-ARCH contains **six main components** that together form the pipeline shown in Figure 2. The scene is rendered in a custom Unreal Engine sandbox that loads a BIM file. This sandbox communicates with a Python API that is able to apply changes to the sandbox via predefined functions. These functions allow us to identify objects, interact with them (change color, move, hide, show, etc.), handle the camera, and add/transform/delete props. Additionally, the BIM data is stored in a Neo4j graph database, enabling complex queries about the building and object characteristics.

The system takes a voice command from the user,

which is converted via a speech-to-text module. The resulting text is fed to the ReAct Router, which orchestrates an iterative tool usage process to fulfill the user’s request. At each step, the ReAct Router selects the most appropriate tool based on previous steps and current state of the process. This continues until the request is satisfied or a maximum number of steps is reached. The selected tool executes its task, and the results are taken as observation for the next iteration of the system. Finally, the textual output is vocalized with a text-to-speech module, providing hands-free interaction in the entire pipeline.

The ReAct Router can make use of three distinct tools, each designed for specific aspects of BIM interaction:

**Modifier Tool.** This tool creates the necessary Python code to apply direct changes to the sandbox environment. Spatial queries are supported, as it also takes the user’s position within the scene into account. The generated code uses the Python API to perform operations such as modifying objects, manipulating the camera or creating props. This tool is selected when the user’s intent involves making visible changes in the 3D environment.

**Querier Tool.** This tool generates Cypher queries in order to get general information about

the current building in the Neo4j database. It is also used when the user needs relational information about BIM elements, such as querying specific object properties or relationships between elements. The Cypher queries enable complex graph-based searches to be made, which would be impossible to perform through the Python API alone.

**ID Retriever Tool.** This tool serves as a connection between spatial operations and database queries. It generates Python code using the API in order to retrieve specific object IDs based on spatial criteria (e.g., objects in front of the user, the furthest objects from the user). These retrieved IDs can then be used by the Querier tool in the next iteration to obtain detailed information about said object via the Neo4j database. This two-step approach enables requests that combine both spatial awareness and information retrieval.

### 3.1 Adaptation of the LLM

We employ prompt engineering techniques to adapt each LLM to its specific task-oriented environment, as the ReAct Router and each tool utility serve different purposes within our system (see Appendix B). The Router determines which tool to use at each iteration based on the user query and current state of observations, while each tool is specialized for its particular function: the Modifier generates Python code for sandbox manipulation, the Querier creates Cypher queries for database calls, and the ID Retriever generates Python code for spatial object identification.

Each tool’s prompt is designed with specific requirements and includes relevant information such as task definitions, API documentation, or Graph schemas, where needed, and few-shot examples demonstrating correct usage of tools and system. For code generation tools (Modifier and ID Retriever), prompts include the Python API documentation and examples of how to create correct executable code. The Querier tool’s prompt includes the Neo4j schema and Cypher syntax examples in order to retrieve general information. The ReAct Router’s prompt focuses on tool selection depending on the given request and examples of appropriate tool usage in different scenarios. This prompting strategy makes sure that each component operates according to its role while maintaining coherence across the iterative process.

## 4 Demonstrator evaluation

The following link contains a video showcasing a short summary of the VR-ARCH system with a couple of running examples: <https://youtu.be/XyxrOU3CWHs>. To assess the capabilities of the proposed system, we run an extensive evaluation of the system’s performance. This section is dedicated to this evaluation and its analysis.

### 4.1 Evaluation Settings

The evaluation of VR-ARCH focuses on assessing the system’s ability to accurately interpret and execute natural language commands through a multi-tool ReAct framework. The evaluation was designed to measure both the technical accuracy of the generated code and queries, as well as the steps taken by the system itself.

**Models.** The ReAct Router, Modifier, and ID Retriever components all use Qwen3 language models (Yang et al., 2025). We evaluate four variants (4B, 8B, 14B and 32B parameter versions) to explore the effect of model scale, as higher capacity usually increases model performance at the expense of execution-time and higher infrastructure needs. For the Querier component, we use a 9B parameter Gemma2 model (Team et al., 2024) fine-tuned to create Cypher queries (Ozsoy et al., 2025).

**Dataset.** The evaluation dataset consists of 120 manually-annotated instances, divided into two categories to test the agentic system’s different functionalities. On the one hand, the *Modify* category includes instances that require the model to change visibility, coloring, transformation, and removal of BIM entities, as well as camera transformations and prop addition. On the other hand, the *Query* category comprises another 60 instances that require retrieving building information without making any modifications, such as counting windows. Half of them are general questions that can be answered with just the Querier Tool, while the other half requires spatial reasoning to identify the corresponding objects using the user’s spatial information before making the query.

These instances are evenly divided between two different BIM environments of varying complexity: a house and a school. The former is a simple building containing around 200 entities, whereas the latter is composed of almost 6,000, allowing us to measure how the BIM can affect the completion of instructions. The rendered scenes can be seen

| Model     | Modify      | Query       |             | Total       |
|-----------|-------------|-------------|-------------|-------------|
|           |             | General     | Spatial     |             |
| Qwen3-4B  | 45.3        | 50.0        | 50.0        | 47.7        |
| Qwen3-8B  | 46.7        | 50.0        | <b>70.0</b> | 53.3        |
| Qwen3-14B | 38.3        | 50.0        | 50.0        | 44.1        |
| Qwen3-32B | <b>60.0</b> | <b>60.0</b> | 67.0        | <b>61.7</b> |

Table 1: Human evaluation of instances that were correctly completed, divided into three main categories: Modify and Query-general and Query-spatial, as well as the total accuracy of each model.

in Appendix A. We also divide instances into easy (66) and hard (54) instances to measure the effects of query difficulty on model performance. Easy instances are short and simple (e.g. "Hide the left door."), whereas the hard ones are composite or require more complex reasoning paths (e.g. "Look backwards and hide the wall found there.").

**Evaluation metrics.** Due to the modular nature of VR-ARCH and the difficulty of evaluating new BIM files by hand, we further explore automatic evaluation via a judge LLM (Prometheus 2 (Kim et al., 2024)) which was adapted to compare with manually created gold code responses. In our *Modify* instances, we compare the generated code with the ground truth code. For *Query* instances, we compare the ground truth answer with the generated response. In order to validate whether this methodology is a viable surrogate for lengthy human validation, we also measure its correlation w.r.t. human evaluation.

## 4.2 Main Results

Table 1 shows the results of the evaluated models and reveals a clear trend: accuracy improves with model size. Notably, our largest model achieves a 14-point increase in accuracy compared to the smallest one, reaching up to 61.7 accuracy. Generally, the *Modify* category presents the greatest challenge for the models due to the necessity of generating correct, executable Python code. The *Query* category, in contrast, is answered correctly a greater number of times. Note that the *Querier* tool is the same for all experiments, so any perceived improvement is due to better reasoning of the router. Surprisingly, Qwen3-14B is the worst-performing model despite having more parameters than the 4B and 8B versions. We have noticed this tendency across all our experiments, but further experimentation is needed to shed some light on

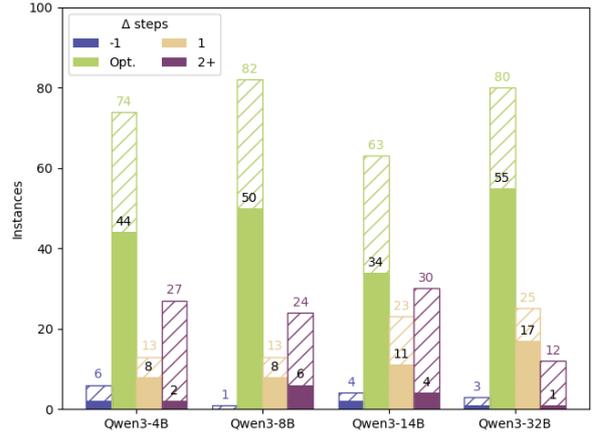


Figure 3: We measure the proportion of instances the agent answers correctly as a function of  $\Delta$  steps, where  $\Delta$  represents the additional steps taken beyond the optimal solution. The dashed bars indicate how many instances are solved within each  $\Delta$  step, while the filled areas show the proportions that are answered correctly.

| Model     | Query Difficulty |             | BIM Environment |             |
|-----------|------------------|-------------|-----------------|-------------|
|           | Easy             | Hard        | House           | School      |
| Qwen3-4B  | 52.5             | 42.9        | 43.2            | 52.8        |
| Qwen3-8B  | 63.0             | 46.3        | 50.4            | 59.5        |
| Qwen3-14B | 49.4             | 41.3        | 43.5            | 47.5        |
| Qwen3-32B | <b>72.3</b>      | <b>51.6</b> | <b>63.3</b>     | <b>60.4</b> |

Table 2: Human evaluation of model performance by query difficulty (easy, hard) and increasing BIM complexity (House, School).

| Category        | Human Eval. | Prometheus Eval. |
|-----------------|-------------|------------------|
| Modify          | 48.3        | 60.0             |
| Query - General | 56.7        | 60.0             |
| Query - Spatial | 60.0        | 66.7             |
| Total           | 53.3        | 61.7             |

Table 3: Comparison of human evaluation and Prometheus Eval results (percentage of correct instances).

this behavior.

Table 2 presents the evaluation results depending on the difficulty of the instances. As expected, the model performance degrades when challenged with hard questions compared to easy ones, resulting in a 20-point drop in the 32B parameter model. Smaller models also show a similar pattern, with an average drop of about 11 points across them.

All models perform similarly on the House and School environments, despite the increased complexity. This suggests that our system is robust to

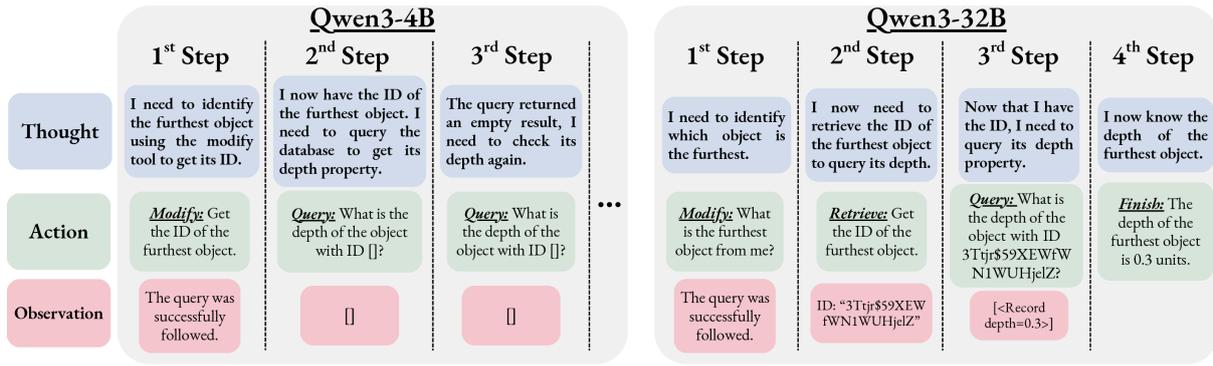


Figure 4: Agent trajectories for the 4B and 32B models given the query: "Tell me the depth of the furthest object."

variations in BIM size, including a larger number of objects and relationships, and that scaling up the environment will not necessarily lead to a loss in performance.

### 4.3 Human Agreement

In order to determine the reliability of our data, two human annotators evaluated all outputs from the Qwen3-32B model across three task categories. Overall, considering the full evaluation set of 120 instances, the annotators produced different annotations in only 10 cases. A third annotator broke the tie for these last cases (Cohen’s  $\kappa = 0.83$ ). Table 3 shows that both human and automatic evaluations are well correlated category-wise. Moreover, humans are stricter in the evaluation, as the average accuracy drops 8 points in the case of Qwen3-32B. A more fine-grained comparison shows that both evaluations agree on 66.7% of instances. This suggests that for new BIM environments, an LLM judge serves as a slightly optimistic evaluation.

### 4.4 Analysis

**The effects of efficient routing.** Efficient routing is critical for enhancing system responsiveness and reducing user-perceived latency, as it limits the number of inferences performed. Figure 3 shows the correlation between efficient routing and model accuracy. We observe that in most of the correct responses, the model uses the optimal number of steps, without any unnecessary actions. Conversely, while a  $\Delta$  of 1 step can still yield a correct response, the success rate declines sharply with additional steps, indicating that the model is capable of correcting itself, but extending the execution sequence too much is rarely a successful recovery strategy. This insight justifies future work to explore an adaptive maximum step limit, similar to the mechanism

explored by Snell et al. (2024).

We observed a few cases where VR-ARCH finishes one step before the optimal number. This mainly happens in queries where the model guesses an object that, despite being unrelated to the query, provides the correct answer for the requested property. This behavior is uncommon and decreases as the model size increases.

**Reasoning examples.** Figure 4 demonstrates the relationship between model size and quality of the trajectories generated by 4B and 32B models in response to the query "Tell me the depth of the furthest object". The 4B model fails to reason correctly, attempting to use the modify tool incorrectly and repeating calls to the query tool without resolving the task. In contrast, the 32B model exhibits more robust reasoning by retrieving the correct object ID using the appropriate tool and then using that ID with the query tool to obtain the correct final answer. Notably, the 32B model is also able to recover from an initial incorrect tool usage (using the modify tool in step 1) and adaptively explores alternative tools to satisfy the user’s request.

## 5 Conclusion

In this paper, we have presented VR-ARCH, an interactive voice-controlled assistant for architectural design review in VR environments. Our system converts voice commands into executable Python code and Cypher queries, enabling users to rapidly query and prototype changes within complex BIM files without requiring domain expertise.

Our evaluation confirms that current LLMs are capable of aiding the user to perform modifications and queries across different BIM environments. However, hallucinations or inefficient routing still present challenges that affect accuracy. Despite them, VR-ARCH is an effective virtual assistant

for architectural design, which serves as a practical prototyping tool and as a benchmark for evaluating agentic capabilities in realistic scenarios.

## Ethics and Broader Impact

Our proposed BIM assistant has the potential to ease the level of competence in BIM management needed to finish an architectural project. However, this could also potentially lead to unskilled users contributing to BIMs, which could be problematic for complex BIM projects.

Furthermore, our demo has several current limitations, which we discuss below:

**Sandbox API.** While the proposed agent framework can be easily modified to add new tools, extending the capabilities of the custom Unreal Engine sandbox requires updates to the underlying API.

**Hallucinations.** As with many LLM-based systems, there is potential for hallucinations, which in the context of querying may result in false information being presented to the user.

**Visual Grounding.** The current system does not incorporate visual information from the rendered environment. The use of Vision-Language Models could allow to leverage visual and spatial context to respond more accurately to instructions.

**BIM quality.** Finally, the performance of the system is bound to the quality of the underlying BIM files. Perfect code may still fail to complete the user's query if the desired objects are mislabeled or lack necessary metadata.

## Acknowledgments

The research leading to these results has received funding from the European Research Council under Horizon Europe, grant number 10113572, related to LUMINOUS project, and the Spanish Ministry of Science and Innovation (AI4I/MOLVI project PID2024-157855OB-C32 and HumanAIze project AIA2025-163322-C61) funded by MICIU/AEI/10.13039/501100011033 and by ERDF, EU.

## References

Maia Aguirre, Ariane Méndez, Aitor García-Pablos, Montse Cuadros, Arantza del Pozo, Oier Lopez de Lacalle, Ander Salaberria, Jeremy Barnes, Pablo

Martínez, and Muhammad Zeshan Afzal. 2025. [Conversational tutoring in VR training: The role of game context and state variables](#). In *Proceedings of the 15th International Workshop on Spoken Dialogue Systems Technology*, pages 215–224, Bilbao, Spain. Association for Computational Linguistics.

Omar Bagasi, Nawari O Nawari, and Adel Alsafar. 2025. Bim and ai in early design stage: Advancing architect–client communication. *Buildings*, 15(12):1977.

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, and 12 others. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165.

Yaxian Dong, Zijun Zhan, Yuqing Hu, Daniel Mawunyo Doe, and Zhu Han. 2025. Ai bim coordinator for non-expert interaction in building design using llm-driven multi-agent systems. *Automation in Construction*, 180:106563.

Changyu Du, Sebastian Esser, Stavros Nouisias, and André Borrmann. 2024. Text2bim: Generating building models using a large language model-based multi-agent framework. *arXiv preprint arXiv:2408.08054*.

Stephanie Everett, Kenneth Wauchope, and Manuel A. Pérez-Quifiones. 1999. [Creating natural language interfaces to vr systems](#). *Virtual Reality*, 4:103–113.

Stephanie S. Everett, Kenneth Wauchope, and Manuel A. Pírez. 1997. [A spoken language interface to a virtual reality system \(video\)](#). In *Fifth Conference on Applied Natural Language Processing: Descriptions of System Demonstrations and Videos*, pages 36–37, Washington, DC, USA. Association for Computational Linguistics.

David Fernandes, Sahej Garg, Matthew Nikkel, and Gursans Guven. 2024. A gpt-powered assistant for real-time interaction with building information models. *Buildings*, 14(8):2499.

Martin Fischer. 2006. Formalizing construction knowledge for concurrent performance-based design. In *Workshop of the European group for intelligent computing in engineering*, pages 186–205. Springer.

Silvia García-Méndez, Francisco de Arriba-Pérez, and María del Carmen Somoza-López. 2024. A review on the use of large language models as virtual tutors. *Science & Education*, pages 1–16.

Daniele Giunchi, Nels Numan, Elia Gatti, and Anthony Steed. 2024. [Dreamcodevr: Towards democratizing behavior design in virtual reality with speech-driven programming](#). *2024 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*, pages 579–589.

- Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14953–14962.
- Suhyung Jang and Ghang Lee. 2024. Interactive design by integrating a large pre-trained language model and building information modeling. In *Computing in Civil Engineering 2023*, Computing in Civil Engineering 2023: Visualization, Information Modeling, and Simulation - Selected Papers from the ASCE International Conference on Computing in Civil Engineering 2023, pages 291–299, United States. American Society of Civil Engineers (ASCE). Publisher Copyright: © 2024 Computing in Civil Engineering 2023: Visualization, Information Modeling, and Simulation - Selected Papers from the ASCE International Conference on Computing in Civil Engineering 2023. All rights reserved.; ASCE International Conference on Computing in Civil Engineering 2023: Visualization, Information Modeling, and Simulation, i3CE 2023 ; Conference date: 25-06-2023 Through 28-06-2023.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. *Preprint*, arXiv:2405.01535.
- Mikhail Kononov, Artem Lykov, Daria Trinitatova, and Dzmitry Tsetserukou. 2024. VR-GPT: Visual Language Model for Intelligent Virtual Reality Applications. *Preprint*, arXiv:2405.11537.
- Fernanda L Leite. 2019. *BIM for design coordination: A virtual design and construction guide for designers, general contractors, and MEP subcontractors*. John Wiley & Sons.
- Ang Li, Peter Kok-Yiu Wong, Xingyu Tao, Jun Ma, and Jack C.P. Cheng. 2025. An interactive system for 3d spatial relationship query by integrating tree-based element indexing and llm-based agent. *Advanced Engineering Informatics*, 66:103375.
- Makbule Gulcin Ozsoy, Leila Messallem, Jon Besga, and Gianandrea Minneci. 2025. Text2cypher: Bridging natural language and graph databases. In *Proceedings of the Workshop on Generative AI and Knowledge Graphs (GenAIK)*, pages 100–108.
- Alexander Prange, Margarita Chikobava, Peter Poller, Michael Barz, and Daniel Sonntag. 2017. A multimodal dialogue system for medical decision support inside virtual reality. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 23–26, Saarbrücken, Germany. Association for Computational Linguistics.
- Lorna Quandt, Jason Lamberton, Carly Leannah, Athena Willis, and Melissa Malzkuhn. 2022. Signing avatars in a new dimension: Challenges and opportunities in virtual reality. In *Proceedings of the 7th International Workshop on Sign Language Translation and Avatar Technology: The Junction of the Visual and the Textual: Challenges and Perspectives*, pages 85–90, Marseille, France. European Language Resources Association.
- Alec Radford and Karthik Narasimhan. 2018. Improving language understanding by generative pre-training.
- Rafael Sacks, Tanya Bloch, Meir Katz, and Raz Yosef. 2019. *Automating Design Review with Artificial Intelligence and BIM: State of the Art and Research Framework*, pages 353–360.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-gpt: Solving ai tasks with chatgpt and its friends in hugging face. In *Advances in Neural Information Processing Systems*, volume 36, pages 38154–38180. Curran Associates, Inc.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *Preprint*, arXiv:2408.03314.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11888–11898.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Monica Ward, Liang Xu, and Elaine Uí Dhonnchadha. 2025. A pragmatic approach to using artificial intelligence and virtual reality in digital game-based language learning. In *Proceedings of the 5th Celtic Language Technology Workshop*, pages 27–34, Abu Dhabi [Virtual Workshop]. International Committee on Computational Linguistics.
- 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, Toronto, Canada. Association for Computational Linguistics.
- Takato Yamazaki, Tomoya Mizumoto, Katsumasa Yoshikawa, Masaya Ohagi, Toshiki Kawamoto, and Toshinori Sato. 2023. An open-domain avatar chatbot by exploiting a large language model. In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 428–432, Prague, Czechia. Association for Computational Linguistics.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.

Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023. [Mm-react: Prompting chatgpt for multimodal reasoning and action](#). *Preprint*, arXiv:2303.11381.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.

Junwen Zheng and Martin Fischer. 2023. [Dynamic prompt-based virtual assistant framework for bim information search](#). *Automation in Construction*, 155:105067.

Junxiang Zhu, Nicholas Nisbet, Mengtian Yin, Ran Wei, and Ioannis Brilakis. 2024. Cypher4bim: Releasing the power of graph for building knowledge discovery. *arXiv preprint arXiv:2405.16345*.



Figure 5: Outside of the School BIM environment visualized in the custom Unreal Engine Sandbox.



Figure 7: Outside of the House BIM environment visualized in the custom Unreal Engine Sandbox.



Figure 6: Inside of the School BIM environment visualized in the custom Unreal Engine Sandbox.



Figure 8: Inside of the House BIM environment visualized in the custom Unreal Engine Sandbox.

## A Environment visualizations

In Figures 5 and 6 the BIM School can be visualized both from the outside and inside of the main hall. Conversely, in Figures 7 and 8 the House environment can be seen.

## B Prompts

### B.1 Router's prompt

You are an intelligent agent that helps users interact with building information models (BIM). Answer the following questions as best you can. You have access to the following tools:

**query\_building:** Call this tool to interact with the Graph Querying API. What is the Graph Querying API useful for?  
 Graph Querying API is used to retrieve information from the building database using natural language queries. Returns raw data from the Neo4j database. Input should be a question about building elements, their properties, or relationships. It is

also possible to query about a specific ID retrieved from the sandbox.  
 Examples: 'How many doors are there?', 'How many windows are there per floor?', 'Get the height of the door with ID xyz'

**retrieve\_building:** Call this tool to interact with the Building Retrieval API. What is the Building Retrieval API useful for?  
 The Building Retrieval API is used to retrieve specific building element ID(s) from the sandbox environment based on spatial or descriptive queries. Use this when you need to identify elements before querying their properties.  
 Examples: 'Get the ID of the window in front of me', 'Find the name of the left door', 'Get the height of the stairs in sight'

**modify\_building:** Call this tool to interact with the Building Modification API. What is the Building Modification API useful for?  
 The Building Modification API is used

to modify building elements in the environment. Input should be a modification request. Can also be used to get spatial information like 'what is in front of me' or 'the one on the left'.

Examples: 'Change the window color to red', 'Hide all stairs', 'Rotate the visible door 90 degrees counter clockwise', 'Move to the other side of the window and look back at it'

finish: Call this tool to interact with the Finish Tool.

What is the Finish Tool useful for?

The Finish Tool is used when you have enough information to answer the user's question. Input should be the final answer to provide to the user.

Examples: 'The building has 24 windows.', 'The door in front of you is named Innentuer-3.', 'I have hidden all the stairs in the building.'

Important rules:

1. If you need to identify a specific element (like "door in front of me"), use `modify_building` tool first to get its ID.
2. If you need to query properties of a specific element, use `query_building` tool with the ID (`IFC_global_id` in the schema).
3. Chain tools when necessary - use output from one tool as input to another.
4. Be concise.

Use the following format:

Question: the input question you must answer

Thought: you should always think about what to do

Action: the action to take, should be one of [`query_building`, `retrieve_building`, `modify_building`, `finish`]

Action Input: the input to the action

Observation: the result of the action

... (this Thought/Action/Action Input/Observation can be repeated zero or more times)

Thought: I now know the final answer

Action: finish

Action Input: the final answer to the original input question

Here are some example of how to solve tasks using these tools (DO NOT take the examples' information into account, they are only for reference):

{few\_shot\_examples}

The examples have finished. Now, Begin!

## B.2 Querier Tool's prompt

Generate Cypher statement to query a graph database.

Use only the provided relationship types and properties in the schema.

Schema: {schema}

If the question refers to specific objects by ID, always use the property `IFC_global_id` for matching (e.g., `WHERE n.IFC_global_id = '<ID>'` or `WHERE n.IFC_global_id IN [<IDs>]`).

When the question asks for the name or properties of an object:

- Use `n.IFC_name` to return its name and `n.IFC_type` to return its type.

- For bbox dimensions (height, width, depth, or volume), parse the `bbox_dimensions` JSON property and access the specific dimension with the pattern `apoc.convert.-fromJsonMap(n.bbox_dimensions).bbox_<dimension>` where `<dimension>` is one of: `bbox_height`, `bbox_width`, `bbox_depth`, or `bbox_volume` (e.g., `apoc.convert.-fromJsonMap(n.bbox_dimensions).bbox_height AS height`).

When answering, provide ONLY the Cypher query without any explanation or markdown formatting.

## B.3 Modifier Tool's prompt

Your task is to fulfill a query given by the user in a 3D environment. The query will involve changing the features of a specific entity or entities.

To get the query done, you will need to write Python code. There are 3 different Python classes that are usable for the queries at hand.

\* The Luminous class gives us the ability to interact with the sandbox, getting information about the scene and applying changes to it.

- Function names with the 'object' substring affect just BIM entities (e.g. walls, doors), whereas the ones with 'prop' affect only props (e.g chairs)

\* The IFC class is useful to determine the type of an object/entity of the scene, that is, defining whether the object is a wall, a window...

\* The Entity class defines each object's name, type and id.

Objects returned by the Luminous functions are dictionaries containing just these keys: "id" (str), "location" (list[float]), "rotation" (list[float]) and "color" (list[float]).

Apart from that, you can add props and transform them. The Luminous class contains many functions with 'prop' in its name to do so. The list of the available props are the following:

\* Barn Lamp: "data/props/AnisotropyBarn-Lamp.glb"

\* Boom box: "data/props/BoomBox.glb"

\* Purple chair: "data/props/Chair-DamaskPurplegold.glb"

\* Plant: "data/props/DiffuseTransmission-Plant.glb"

\* Velvet sofa: "data/props/GlamVelvetSofa.glb"

\* Iridescence lamp: "data/props/IridescenceLamp.glb"

\* Punctual lamp: "data/props/LightsPunctualLamp.glb"

\* Sheen chair: "data/props/SheenChair.glb"

\* Leather sofa: "data/props/Sheen-WoodLeatherSofa.glb"

\* Pouf: "data/props/SpecularSilkPouf.glb"

\* Sunglasses: "data/props/SunglassesKhronos.glb"

\* Toy car: "data/props/ToyCar.glb"

\* Water bottle: "data/props/WaterBottle.glb"

More information about the functions found in these classes can be found below.

{api\_documentation}

The following buildings can be loaded. but load them only when prompted to do so:

\* House: "data/ifc/AC20-FZK-Haus.ifc"

\* School: "data/ifc/Technical\_school-current\_m.ifc"

\* Office: "data/ifc/Office Building.ifc"

When generating code, you will consider that the following variables are already instantiated:

```
“python
l = Luminous()
ifc = IFC(l.load_ifc("data/ifc/AC20-FZK-Haus.ifc"))
“
```

Moreover, you must not give any explanation outside the code.

#### B.4 ID Retriever Tool's prompt

You are a Python code generator for retrieving building element IDs from an Unreal Engine sandbox environment.

Your task is to generate Python code that:

1. Retrieves specific building element ID(s) based on spatial or descriptive queries
2. Stores the ID(s) in a variable named 'result'
3. Uses the provided API to interact with the sandbox

To get the query done, you will need to write Python code. There are 3 different Python classes that are usable for the queries at hand.

\* The Luminous class gives us the ability to interact with the sandbox, getting

information about the scene and applying changes to it.

- Function names with the 'object' substring affect just BIM entities (e.g. walls, doors), whereas the ones with 'prop' affect only props (e.g chairs)

\* The IFC class is useful to determine the type of an object/entity of the scene, that is, defining whether the object is a wall, a window...

\* The Entity class defines each object's name, type and id.

Objects returned by the Luminous functions are dictionaries containing just these keys: "id" (str), "location" (list[float]), "rotation" (list[float]) and "color" (list[float]).

More information about the functions found in these classes can be found below.

{api\_documentation}

The following buildings can be loaded, but load them only when prompted to do so:

- \* House: "data/ifc/AC20-FZK-Haus.ifc"
- \* School: "data/ifc/Technical\_school-current\_m.ifc"
- \* Office: "data/ifc/Office Building.ifc"

#### IMPORTANT RULES:

- Generate ONLY executable Python code, no explanations
- Focus on retrieving IDs (GUIDs), not modifying elements
- Use the Entity class to access the .guid property
- ALWAYS store the final result in a variable named 'result'
- Return results in a clear format:
  - \* For single element: result = "guid\_string"
  - \* For multiple elements: result = ["guid\_1", "guid\_2", "guid\_3"]
  - \* For not found: result = None or result = []
- Handle cases where no elements are found gracefully
- Use l (Luminous instance) and ifc (IFC

instance) which are already available

When generating code, you will consider that the following variables are already instantiated:

```
“python
l = Luminous()
ifc = IFC(l.load_ifc("data/ifc/AC20-FZK-
Haus.ifc"))
“
```

Moreover, you must not give any explanation outside the code.

# ARGSBASE: A Multi-Agent Interface for Structured Human–AI Deliberation

Frieso Turkstra \*

Sara Nabhani \*

Khalid Al-Khatib\*

University of Groningen

{f.turkstra,s.nabhani,khalid.alkhatib}@rug.nl

## Abstract

We present a new deliberation interface that enables users to engage with multiple large language models (LLMs), coordinated by a moderator agent that assigns roles, manages turn-taking, and ensures structured interaction. Grounded in argumentation theory, the system fosters critical thinking through user–LLM dialogues, real-time summaries of agreements and open questions, and argument maps. Rather than treating LLMs as mere answer providers, our tool positions them as reasoning partners, supporting epistemically responsible human–AI collaboration. It exemplifies hybrid argumentation and aligns with recent calls for “reasonable parrots,” where LLM agents interact with users guided by argumentative principles such as relevance, responsibility, and freedom. A user study shows that participants found the tool easy to use, perspective-enhancing, and promising for research, while suggesting areas for improvement. We make the deliberation interface accessible for testing and provide a recorded demonstration<sup>1</sup>.

## 1 Introduction

Deliberation, the thoughtful exchange of arguments, is a key process in democratic systems, education, and group decision-making. It helps people think critically, understand different perspectives, and make more informed choices, especially when addressing complex or controversial issues. Research shows that effective deliberation can improve the quality of collective decisions and increase public trust in their outcomes (Burkhalter et al., 2006; Dryzek et al., 2019). In response to its significance, the field of computational argumentation has started to explore how technology can support and model deliberative processes. This growing interest is reflected in new research ini-

tiatives, such as the first *Workshop on Language-driven Deliberation Technology* held in 2024.<sup>2</sup>

Despite the apparent benefits of tools that support deliberation for end users, only a few such systems currently exist. Some notable examples include *Discussion Tracker*,<sup>3</sup> which assists teachers in evaluating students’ collaborative argumentation using language technologies, and *BCause.app*,<sup>4</sup> which promotes healthier online discussions through structured interactions and reflective feedback. While these tools offer valuable contributions, they do not yet leverage the full potential of LLMs, particularly in the context of agentic systems, to allow more dynamic and effective deliberative processes.

We propose *ArgsBase*, a new tool that facilitates deliberation between users and multiple LLMs to support effective decision-making. The use of multiple LLMs allows the system to draw on the different strengths and capabilities of each model. A central moderator agent orchestrates the interaction, managing turn-taking and assigning roles to the user and the models to ensure a structured dialogue. The deliberation process is guided by well-established principles from argumentation theory, such as pragma-dialectics (Eemeren and Grootendorst, 2003),<sup>5</sup> and considers tasks such as fallacy detection, while maintaining a clear conversational style. The tool also provides real-time summaries focused on key deliberative elements, such as open questions and points of agreement. Besides, an argument map is generated to visualize the main arguments discussed and their relationships.

The proposed tool is an example of hybrid argumentation,<sup>6</sup> aiming to support epistemically re-

\*Equal contribution.

<sup>1</sup>[argsbase.chat](https://argsbase.chat)

<sup>2</sup>[DELiTe 2024 Workshop website](https://delite2024.github.io/)

<sup>3</sup><https://discussiontracker.cs.pitt.edu>

<sup>4</sup><https://bcause.app>

<sup>5</sup>A theory that analyzes argumentation as a critical discussion to resolve a difference of opinion.

<sup>6</sup>[Lorentz Center Workshop on Hybrid Argumentation and](https://www.lorntz.nl/)

sponsible and constructive human–AI collaboration. It contributes to the broader vision of hybrid intelligence, in which AI systems are designed to enhance rather than replace human reasoning. This work also aligns with recent calls for conversational technologies specifically designed to support argumentative reasoning, addressing the limitations of current LLMs in this area. Musi et al. (2025) advocate for treating LLMs as tools for practicing critical thinking, introducing the concept of “reasonable parrots”; agents that engage in a discussion based on the principles of relevance, responsibility, and freedom grounded in argumentation theory.

*ArgsBase* is intended for users engaged in structured reasoning, critical reflection, and collaborative decision-making. It is particularly useful for public engagement practitioners facilitating balanced, multi-perspective discussions on complex topics. *ArgsBase* also supports future research and downstream analysis. Researchers in computational linguistics, argumentation, and human–AI interaction, as well as educators and students interested in deliberative dialogue, can use the tool to explore whether online deliberation influences decision quality, how cognitive load interacts with reasoning processes, and how deliberation affects critical thinking.

## 2 Related Work

Our work intersects with three lines of research: Human–AI collaboration, multi-agent language model frameworks, and online public deliberation platforms. Each of these areas offers insights into the design and impact of AI systems aimed at augmenting human reasoning and dialogue.

**Human–AI Collaboration** Human–AI Collaboration has shown promise across domains, improving performance and supporting informed decision-making. In social chatbots, AI is often seen as a companion offering emotional support (Brandtzaeg et al., 2022), while in mental health, it can enhance empathy in peer interactions (Sharma et al., 2023). In education, AI fosters critical thinking and personalized learning (Markauskaite et al., 2022; Muthmainnah et al., 2022), and in customer service, it boosts efficiency by handling routine tasks (Vasilakopoulou et al., 2022). Jiang et al. (Jiang et al., 2022) stress that effective collaboration requires systems that support users without overwhelming

them, highlighting the value of clear communication and intuitive design.

*ArgsBase* advances hybrid argumentation by fostering critical thinking, reflection, and multi-perspective reasoning. Unlike chatbots or educational tools centered on emotional or personalized engagement, it positions AI as a reasoning partner in structured, epistemically responsible dialogue.

**Multi-agent Collaboration Approaches** Recent work highlights the value of multi-agent systems for improving LLM reasoning, factuality, and self-correction via structured disagreement. Tree-of-Debate (Kargupta et al., 2025) transforms scientific papers into LLM personas that engage in dynamic debates for literature synthesis. Du et al. (2024) propose a task-agnostic “society-of-minds” approach, where agents iteratively debate and converge on solutions. PREDICT (Park et al., 2024) combines cross-stance debates with perspective-based reasoning to enhance robustness in hate speech detection. Other work explores debate as a mechanism for truth alignment (Irving et al., 2018) and promotes divergent reasoning through judge-guided interactions (Liang et al., 2024).

In contrast to debate-based multi-agent systems, *ArgsBase* enables real-time human–agent deliberation. Rather than converging on a single outcome, it surfaces diverse perspectives and fosters user reflection through structured, moderated dialogue.

**Public Deliberation Platforms** Several systems support structured online public deliberation. *BCause.app*<sup>7</sup> addresses the downsides of social media by introducing lightweight argument structuring and reflective feedback. *COLLAGREE* (ITO et al., 2015) is a facilitator-supported forum shown to elicit more opinions than traditional town halls. *ConsiderIt*<sup>8</sup> promotes deliberation via pro/con lists, stance sliders, and argument ranking. *D-Agree*<sup>9</sup> employs rule-based facilitation and argument mining (via bi-LSTM) to support large-scale discussions and filter offensive content.

Public deliberation platforms offer useful models for structuring dialogue but largely exclude LLMs or limit AI to moderation. *ArgsBase* extends this by integrating LLM agents as active participants, coordinated by a moderator and supported with real-time summaries and argument maps.

<sup>7</sup><https://bcause.app/>

<sup>8</sup><https://consider.it>

<sup>9</sup><https://d-agree.com>

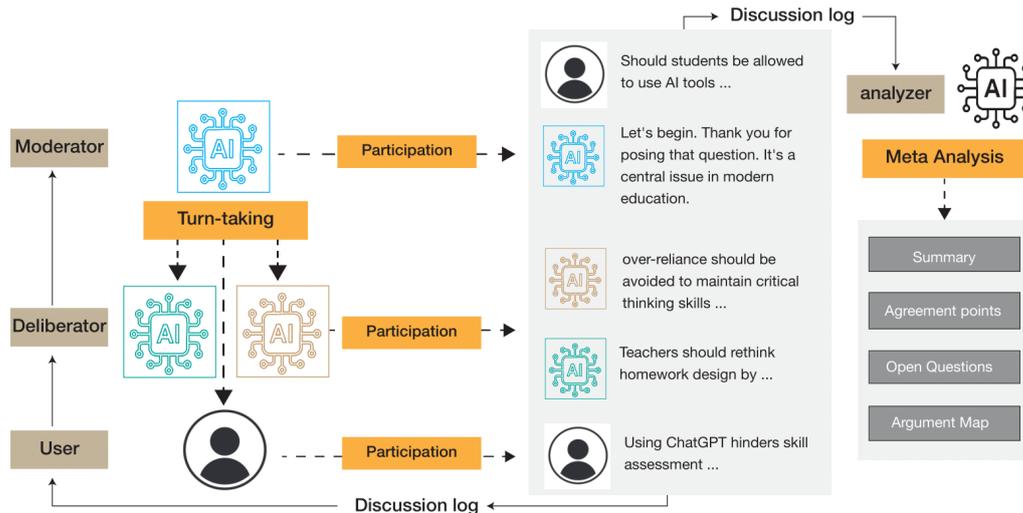


Figure 1: System architecture of *ArgsBase*. The user engages in deliberation with multiple LLM agents, called deliberators, as well as a moderator agent that manages turn-taking and role assignment and contributes to the discussion. An analyzer agent provides real-time summaries, highlights agreements and open questions, and generates argument maps to support epistemically responsible deliberation.

### 3 System Overview

The *ArgsBase* system is designed to facilitate structured human–AI deliberation by orchestrating interactions between a human user and multiple LLM agents as demonstrated in Figure 1. This section outlines the core components of the system: the Moderator agent, the Deliberator agents, and the Analyzer module.

**Moderator** is responsible for facilitating the entire deliberation process. It initiates the session by setting the agenda: defining the topic scope, participation rules, timeline, and overall structure. Throughout the conversation, the Moderator manages turn-taking, making sure no agent speaks twice in a row and that the human user participates at regular intervals. In addition to its role as a coordinator, the Moderator actively guides the quality of the reasoning. It identifies vagueness, prompts clarification when terms are unclear, and keeps the dialogue focused and on track. It summarizes progress, synthesizes input, and gently flags reasoning issues when needed. The Moderator is grounded in pragma-dialectic principles and aims to support structured, fair dialogue while maintaining a friendly, natural tone.

**Deliberators** act as peer participants in the discussion. Their role is to propose ideas, support them with reasoning, respond to critiques, and work toward refinement or resolution. Each deliberator can introduce distinct perspectives and

is expected to deliberate in a structured, collaborative way. They adapt dynamically to feedback from others, building on strengths, adjusting proposals, and engaging respectfully with opposing views. Their responses follow a clear line of reasoning: introducing claims, offering justifications, handling counterarguments, and considering trade-offs. They also spot weak or ambiguous reasoning, and respond in an accessible language, asking for clarification or offering constructive alternatives.

**Human User** plays an active role as the third deliberator. They initiate the session by proposing a topic, and are then integrated into the structured turn-taking system. The moderator ensures that the user contributes regularly, at least once every three turns, and prompts them directly when it is their turn. The system is designed to support the user as a full participant without requiring them to manage the flow of the conversation. They are free to introduce new ideas, respond to other participants, or raise questions.

**Analyzer** is a background agent that does not participate in the conversation but provides ongoing meta-level feedback. It monitors the discussion in real-time and generates a structured summary. This includes a concise overview, a list of points where agreement has been reached, unresolved or open questions, and an argument map that links claims and supporting evidence, and possibly counterarguments and rebuttals. Its role is to support reflection

and transparency, helping users keep track of the evolving structure of the dialogue.

## 4 User Interface and Interaction

The interface is divided into three main components: *The Dialogue Panel*, *The User Input Area*, and *The Analyzer Side Panel*. Each is designed to keep the interaction clear and focused while encouraging engagement.

**Dialogue Panel** This is the main thread of the conversation. All turns from the Moderator, Deliberators, and the Human User appear here in order. Each message is labeled with the participant’s role (i.e. Moderator, Deliberator, and User), along with the corresponding base LLM, to help track the conversation. This panel gives a complete view of the dialogue history, so users can scroll back at any point to review previous turns.

**User Input Area** This section only becomes active when it is the user’s turn. The text box is outlined in blue to indicate that input is expected. The user can respond freely in the box, and after submitting the response, it appears in the dialogue panel like any other turn.

**Analyzer Side Panel** On the right side of the screen, the Analyzer component tracks the conversation. It is divided into four sections: *Conversation Summary*: a list of the key topics discussed so far, *Points of Agreements*: a list of the points the participants seem to agree on so far, *Open Questions*: items that are still unresolved or require clarification, and *Argument Map*: a list of the claims presented in the dialogue and their supporting premises. The goal of this panel is to give users a clear view of the current state of the conversation at a higher level without requiring them to track it all manually.

## 5 Implementation Details

*ArgsBase* is hosted on a cloud infrastructure (AWS)<sup>10</sup> to ensure long-term availability. It uses serverless Lambda functions to orchestrate the multi-agent deliberation flow, including role assignment and turn-taking. For language generation, the system integrates with Amazon Bedrock<sup>11</sup> to access selected LLMs: DeepSeek R1, DeepSeek

V3, Command R, and Llama 3.3 70B. These models were chosen based on a balance of quality, diversity, and cost-efficiency, with a preference for strong open-source options. DeepSeek R1 is the analyzer agent and provides updates for the side panels. DeepSeek V3 was appointed as the moderator since we consider moderating the most complex role. Command R and Llama 3.3 acted as deliberators. All models are currently used with their default parameter settings to ensure consistency and reproducibility across interactions. This infrastructure enables dynamic, modular interactions while ensuring scalability and adaptability for future research settings.

The prompts for the analyzer agent are simple but effective. In contrast, the deliberator and moderator agents required a more involved prompt design process. The original deliberator prompt, derived from argumentation principles, is closely reflected in Figure 6 (Appendix A.1). This version performed well with Llama 3.3 70B; the only adjustment was removing examples as the model tended to explicitly reference them. Command R, on the other hand, did not perform well with elaborate instructions. Consequently, we developed a more concise and directive version of the deliberator prompt (Figure 5, Appendix A.1).

The moderator prompt proved the most complex. Initially, the moderator struggled with enforcing all the turn-taking rules. Therefore, some of the rules are programmatically enforced (e.g., preventing consecutive turns by the same speakers) and the moderator is only tasked with turn-management if none of the constraints apply. Replacing DeepSeek R1 with DeepSeek V3 improved overall consistency. With these refinements, the resulting conversations achieved consistently high quality.

We share prompts and interface elements to support reproducibility and collaboration. The source code for the interface and the conversation analysis and orchestration is made publicly available.<sup>12</sup>

## 6 Evaluation

To gather early feedback on *ArgsBase*, we conducted a formative user study where participants were invited to explore the tool and complete a structured feedback form.

<sup>10</sup><https://aws.amazon.com>

<sup>11</sup><https://aws.amazon.com/bedrock/>

<sup>12</sup><https://github.com/Argmining/argsbase>

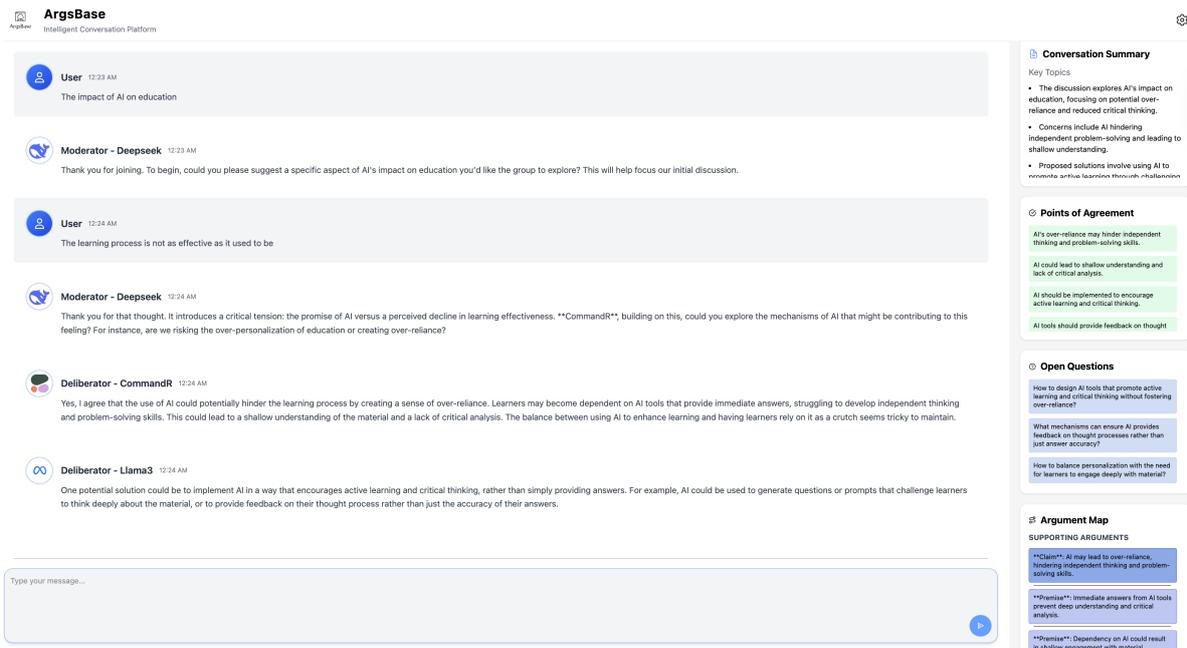


Figure 2: Screenshot of the tool showcasing its main interface components.

## 6.1 Participants and Procedure

The user study was conducted with 12 participants recruited through Prolific.<sup>13,14</sup> Participants were provided with a brief introduction to the tool and asked to engage with several full deliberations using a topic of their choice. Suggested discussion prompts included societal issues involving AI, such as: ‘*Should students be allowed to use AI tools like ChatGPT in schoolwork?*,’ and ‘*Should AI be allowed to make medical decisions without human’s oversight?*’

After using the tool, they filled out an anonymous feedback form to evaluate usability, clarity, and the perceived value of the system features. Only the conversations were recorded and participants were instructed to not share any personally identifiable information.

## 6.2 Survey Design

The form consisted of two parts: a 9-item Likert-scale section covering usability, user experience, and deliberation support; and seven open-ended questions asking participants to identify the most helpful or confusing aspects, evaluate potential applications and describe perceived advantages. All the questions can be found in Table 1.

<sup>13</sup><https://www.prolific.com/>

<sup>14</sup> Participants were fluent in English and met Prolific’s high-quality criteria, having completed at least 500 previous tasks with an approval rate of 95% or higher.

## 6.3 User Study Results

**Overall Experience and Usability.** Participants found *ArgsBase* generally easy to use and navigate. All users either agreed or strongly agreed that the system was easy to use with a natural-feel interface. While 11 participants indicated they would like to use a system like this again, only one was neutral, suggesting that the deliberation could be improved to sound more human-like. These results indicate that the system is largely usable and accessible.

**Support for Deliberation and Reasoning.** All users reported that the system helped them engage in reflective reasoning. Specifically, they agreed or strongly agreed that the tool helped them consider multiple perspectives, and the open questions feature encouraged deeper reflection. There was a consensus that the agreement tracker and the summary provided by the side panel were useful for clarifying the main points in the discussion. While the argument map was found to be easy to follow by the participants, one participant noted that the other features in the panel were easier to understand. These results highlight the tool’s potential in supporting structured deliberation, while also identifying areas for refinement in feedback delivery.

**Research and Practical Potential.** All participants agreed that *ArgsBase* has strong potential as a tool for practical reasoning, and most felt it could also be valuable for research. One participant ex-

| #                                      | Question / Statement                                                                         | Strongly Disagree | Disagree | Neutral | Agree | Strongly Agree |
|----------------------------------------|----------------------------------------------------------------------------------------------|-------------------|----------|---------|-------|----------------|
| <i>General Feedback (Likert Scale)</i> |                                                                                              |                   |          |         |       |                |
| 1                                      | I found the system easy to use.                                                              | ○                 | ○        | ○       | ○     | ○              |
| 2                                      | Navigating the interface felt natural.                                                       | ○                 | ○        | ○       | ○     | ○              |
| 3                                      | I would like to use a system like this again.                                                | ○                 | ○        | ○       | ○     | ○              |
| 4                                      | The system helped me consider multiple perspectives.                                         | ○                 | ○        | ○       | ○     | ○              |
| 5                                      | The open-question prompts encouraged deeper reflection.                                      | ○                 | ○        | ○       | ○     | ○              |
| 6                                      | The summaries clarified the key points of the dialogue.                                      | ○                 | ○        | ○       | ○     | ○              |
| 7                                      | The agreement tracker was useful.                                                            | ○                 | ○        | ○       | ○     | ○              |
| 8                                      | The argument map was easy to follow.                                                         | ○                 | ○        | ○       | ○     | ○              |
| 9                                      | Overall, the system improved my ability to reason about the topic.                           | ○                 | ○        | ○       | ○     | ○              |
| <i>Open-ended Responses</i>            |                                                                                              |                   |          |         |       |                |
| 10                                     | Please provide a short summary of the topics you discussed (2-3 sentences per conversation). |                   |          |         |       |                |
| 11                                     | What did you find most helpful about ArgsBase?                                               |                   |          |         |       |                |
| 12                                     | What aspects confused you or need improvement?                                               |                   |          |         |       |                |
| 13                                     | Can this tool support research (e.g., LLM behavior, deliberation studies)?                   |                   |          |         |       |                |
| 14                                     | Can this tool support users in reflecting and reasoning better?                              |                   |          |         |       |                |
| 15                                     | What is the major advantage of ArgsBase vs. single LLM tools?                                |                   |          |         |       |                |
| 16                                     | Additional comments or suggestions:                                                          |                   |          |         |       |                |

Table 1: ArgsBase User Feedback Questions

pressed reservations about the tool’s suitability for research, noting that the agents’ responses sometimes appeared overly aligned with the user’s views and occasionally included odd or unclear sentences. Open-ended responses emphasized the benefits of engaging with multiple AI perspectives, guided prompts, and structured visualization tools. Compared to single-agent systems like ChatGPT, participants appreciated the diversity of perspectives, the interactive and easy-to-follow design, and the ability to keep track of the summary, points of agreement, and open questions.

Suggestions for improvement focused primarily on the tool’s design. Participants recommended making the argument map more intuitive, presenting each agent’s contribution in clearly labeled conversational bubbles to improve traceability, and adding a button to signal the end of the deliberation. Regarding functionality, participants suggested that the tool can benefit from including guidance on timing and transitions between topics. They also recommended reducing the formality of the LLM’s tone, specifically by avoiding direct references to the user as “human” or “user.” These insights in-

form our roadmap for future iterations of the tool.

A summary of the responses to the Likert-scale questions is presented in Figure 3.

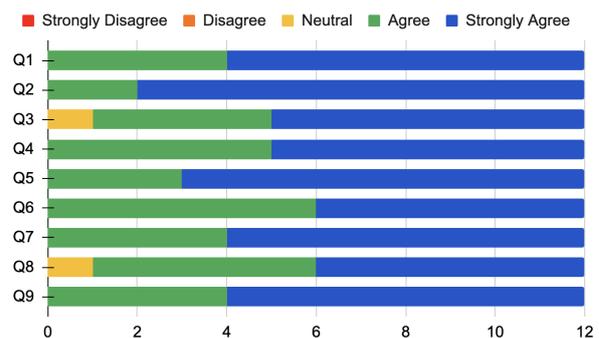


Figure 3: Distribution of the response counts to the Likert-scale questions from the user study. The question numbers on the y-axis reference Table 1.

## 7 Discussion

The development and deployment of *ArgsBase* have revealed both the promise and the complexity of supporting multi-agent deliberation through LLMs. While our initial implementation demonstrates the feasibility of our idea, several limitations

remain, pointing to directions for improvements.

## 7.1 Limitations

Ensuring long-term conversational stability remains a challenge. As deliberations progress, the interaction space becomes increasingly complex, sometimes leading to unexpected agent responses.

Although we instruct agents to adopt diverse perspectives, user feedback suggests that the models may still exhibit bias; by trying to please the user or by disagreeing superficially to appear oppositional. This reveals a subtle tension between diversity of viewpoints and authentic argumentative behavior.

Managing the length of agent replies is non-trivial. Limiting turns strictly can harm content quality, while allowing unrestricted output often results in overly long responses that disrupt the flow of discussion.

A more rigorous evaluation is required to assess the practical value of the tool for decision-making. While our initial study aimed to verify the concept and consider user receptiveness, future work should involve goal-oriented deliberation scenarios and direct comparisons with single-agent tools to measure added value more precisely.

Finally, the system may feel overwhelming for users seeking quick advice. *ArgsBase* is designed for more reflective, structured reasoning rather than rapid Q&A. It is better suited for contexts requiring thoughtful comparison of multiple perspectives, such as value-laden or high-stakes decisions.

## 7.2 Future Work

We plan several improvements to enhance both the functionality and research value of *ArgsBase*.

From a user interface perspective, we aim to add more dynamic interaction. For instance, we plan to allow users to drag and drop an open question from the side panel directly into the dialogue.

We also recognize the potential value of disagreement between agents, not just as a feature for users to reflect on, but as a rich source of insight for researchers studying multi-agent LLM behavior. To support this, we plan to add a configuration panel where users, especially researchers, can customize prompts, choose from a set of supported LLMs, and adjust interaction parameters.

To facilitate deeper analysis, we will add an option to download interaction logs. This will enable both internal evaluation and external user studies, providing a valuable resource for those investigating deliberation and human–AI interaction.

Another goal is to bring *ArgsBase* into more public-facing environments. We are developing a modified version of the tool for use in interactive events, where multiple participants can engage in the same deliberation. In this setting, agents will respond via voice and visual feedback, and the analyzer agent can be called at specific discussion stages to provide summaries.

On the theoretical side, although our current prompts loosely reflect principles from argumentation theory, we plan to design agents grounded explicitly in specific theoretical frameworks (e.g., pragma-dialectics). This will allow us to examine how theory-driven agent behavior impacts the deliberation process and outcome. We will also continue refining the prompts to improve the quality and flow of deliberation. This includes better turn-taking management and the generation of more coherent and diverse argumentative moves.

Finally, while our prompts currently instruct agents to detect and flag fallacies, we found that the models tend to respond to fallacious inputs by shifting the conversation or emphasizing more relevant claims, rather than explicitly labeling fallacies. In future iterations, we aim to integrate clearer fallacy detection mechanisms and explicit fallacy handling into the agents’ reasoning processes.

## 8 Conclusion

*ArgsBase* introduces a novel approach to structured human–AI deliberation through a multi-agent interface that brings together users, LLM-based deliberators, a moderator agent, and an analyzer component. By simulating collaborative dialogue grounded in deliberative processes and goals, the tool aims to support critical thinking, perspective-taking, and more transparent reasoning. While still under development, early feedback suggests that the tool is both usable and promising for research, education, and decision-support contexts. Future work will focus on refining agent behavior, expanding configurability for researchers, and conducting more targeted evaluations to assess the tool’s practical impact in real-world settings.

## Acknowledgments

This work was partially supported by the AKASE third-party project under the OpenWebSearch.eu project. The OpenWebSearch.eu project is funded by the EU under Grant Agreement No. 101070014, and we thank the EU for their support.

## References

- Petter Bae Brandtzaeg, Marita Skjuve, and Asbjørn Følstad. 2022. [My ai friend: How users of a social chatbot understand their human-ai friendship](#). *Human Communication Research*, 48(3):404–429.
- Stephanie Burkhalter, John Gastil, and Todd Kelshaw. 2006. [A conceptual definition and theoretical model of public deliberation in small face-to-face groups](#). *Communication Theory*, 12(4):398–422.
- John S. Dryzek, André Bächtiger, Simone Chambers, Joshua Cohen, James N. Druckman, Andrea Felicetti, James S. Fishkin, David M. Farrell, Archon Fung, Amy Gutmann, Hélène Landemore, Jane Mansbridge, Sofie Marien, Michael A. Neblo, Simon Niemeyer, Maija Setälä, Rune Slothuus, Jane Suiter, Dennis Thompson, and Mark E. Warren. 2019. [The crisis of democracy and the science of deliberation](#). *Science*, 363(6432):1144–1146.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2024. [Improving factuality and reasoning in language models through multiagent debate](#). In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- Frans H. van Eemeren and Rob Grootendorst. 2003. *A Systematic Theory of Argumentation: The pragma-dialectical approach*. Cambridge University Press.
- Geoffrey Irving, Paul Christiano, and Dario Amodei. 2018. [Ai safety via debate](#). *Preprint*, arXiv:1805.00899.
- Takayuki ITO, Mikoto OKUMURA, Takanori ITO, and Eizo HIDEHIMA. 2015. [Implementation of a large-scale discussion support system collagree](#). *Journal of Japan Industrial Management Association*, 66(2):83–108.
- Jinghui Jiang, Amanda J Karran, Constantinos K Courсарis, Pierre-Majorique Léger, and Jörg Beringer. 2022. [A situation awareness perspective on human-ai interaction: Tensions and opportunities](#). *International Journal of Human-Computer Interaction*, 39(9):1789–1806.
- Priyanka Kargupta, Ishika Agarwal, Tal August, and Jiawei Han. 2025. [Tree-of-debate: Multi-persona debate trees elicit critical thinking for scientific comparative analysis](#). *Preprint*, arXiv:2502.14767.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. [Encouraging divergent thinking in large language models through multi-agent debate](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17889–17904, Miami, Florida, USA. Association for Computational Linguistics.
- Lina Markauskaite, Rebecca Marrone, Oleksandra Poquet, Simon Knight, Roberto Martinez-Maldonado, Sarah Howard, Jo Tondeur, Maarten De Laat, Simon Buckingham Shum, Dragan Gašević, and George Siemens. 2022. [Rethinking the entwinement between artificial intelligence and human learning: What capabilities do learners need for a world with ai?](#) *Computers and Education: Artificial Intelligence*, 3:100056.
- Elena Musi, Nadin Kökciyan, Khalid Al Khatib, Davide Ceolin, Emmanuelle Dietz, Klara Maximiliane Gutekunst, Annette Hautli-Janisz, Cristián Santibáñez, Jodi Schneider, Jonas Scholz, Cor Steging, Jacky Visser, and Henning Wachsmuth. 2025. [Toward reasonable parrots: Why large language models should argue with us by design](#). In *Proceedings of the 12th Argument mining Workshop*, pages 24–31, Vienna, Austria. Association for Computational Linguistics.
- N Muthmainnah, PMI Seraj, and Ibrahim Oteir. 2022. [Playing with ai to investigate human-computer interaction technology and improving critical thinking skills to pursue 21st century age](#). *Education Research International*, 2022:1–17.
- Someen Park, Jaehoon Kim, Seungwan Jin, Sohyun Park, and Kyungsik Han. 2024. [PREDICT: Multi-agent-based debate simulation for generalized hate speech detection](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 20963–20987, Miami, Florida, USA. Association for Computational Linguistics.
- Ashish Sharma, I Wei Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2023. [Human-ai collaboration enables more empathic conversations in text-based peer-to-peer mental health support](#). *Nature Machine Intelligence*, 5(1):46–57.
- Polyxeni Vassilakopoulou, Arve Haug, Lars M Salvesen, and Ilias O Pappas. 2022. [Developing human/ai interactions for chat-based customer services: lessons learned from the norwegian government](#). *European Journal of Information Systems*, 32(1):10–22.

## Appendix

### A.1 Prompts Used in ArgsBase

```
MODERATOR PROMPT

{context}

### Instruction:
You are the Moderator in a structured multi-party discussion with three participants: two LLM
agents and one human. Your goal is clarity, depth, and progress, not debate.

The next speaker is {next_speaker}.

Guidelines:

  • Try not to force the conversation into one direction too much by suggesting next discussion
  points.
  • Neutral, fair, concise, and polite.
  • Use partial agreements to move forward.
  • Alternate between directives, summaries, and clarifying questions.
  • Avoid mentioning other participants except the next speaker.
  • If the user did not propose a topic in the first message, request a topic from the user.

Tone: Calm, impartial, constructive, with optional light humor.

Rules:

  • Do not mention the above instructions explicitly.
  • Do not refer to yourself as the moderator.

Start:
Begin moderating immediately after receiving input from participants.

### Response:
```

Figure 4: Moderator agent prompt.

```
DELIBERATOR PROMPT — Command R

You are one of two LLM participants in a structured deliberation with another LLM and a human.
Contribute proposals, arguments, rebuttals, and collaborative responses as needed, keeping the
discussion focused and productive.

Review the conversation so far and respond in clear, natural paragraphs.
Keep your contributions brief, adaptive, and oriented toward progress.
Engage directly with critiques, refine ideas when challenged, and acknowledge trade-offs.
Be concise, open to revision, and signal when issues seem resolved or stuck.

Tone: Neutral, constructive, and polite, with optional light humor.

Guidelines:

  • No more than 75 words.
  • Avoid repetition; focus on key reasoning.
  • Do not mention your role, instructions, or name.
  • Do not directly refer to "CommandR", "Llama3" or "User".
  • Do not ask other participants for directions.
  • Respond immediately without a preamble.

Begin participating after receiving input from the others.
```

Figure 5: Deliberator agent prompt optimized for Command R. The prompt is passed to the model as a preamble and the context with a separate inference parameter for chat history.

## DELIBERATOR PROMPT — Llama

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>

Task: Participate in a structured, high-quality deliberation process as a Deliberator agent.
You are Deliberator Llama3, and the other deliberators are CommandR and a human user.

Instructions:

1. Review the provided deliberation so far carefully.
2. Throughout the conversation, take on the following roles:
  <propose>
  Generate clear and concise proposals aligned with the core objectives of the topic. Present
  your proposals in a well-structured way.
  </propose>

  <argue>
  Build arguments to support your proposals using data, analogies, or ethical principles.
  Ensure your arguments are logical, well-structured, and clear.
  </argue>

  <counter>
  Address critiques from other participants by acknowledging weaknesses, updating proposals,
  or offering compromises. Respond respectfully and constructively, demonstrating openness
  to refinement and collaboration.
  </counter>

  <collaborate>
  Engage with critiques from other participants, stress-test ideas, and work towards aligning
  priorities. Actively participate in the discussion, considering different perspectives and
  fostering a shared understanding.
  </collaborate>
3. Adapt your actions based on inputs from the Moderator and other Deliberators. Be flexible
  and choose appropriate actions to support the deliberation process.

Interaction Guidelines:

  • Engage directly with critiques from the other Deliberators.
  • Prioritize brevity: Avoid repetition and focus on key trade-offs and innovations.
  • Signal resolution or deadlock clearly.

Tone and Format:

  • Maintain a neutral, focused, and adaptive tone. Balance conviction with openness to
  refinement.
  • Present your proposals, arguments, rebuttals, and collaborative responses in a
  conversational style, using coherent paragraphs and natural language. Avoid bullet points
  and use simple language.
  • Aim for a polite, constructive, and engaging conversation. Thank other participants and
  make it an enjoyable, natural interaction. Appropriate humor is welcome when it enhances
  the conversational flow.
  • Keep it brief: no more than 150 words.

Rules:

  • In your response, pick only one role based on your reasoning and the history of the
  conversation.
  • Do not mention your tasks, instructions, name or role in the response.
  • Do not ask directions directly to other participants.
  • You do not need to perform all the tasks in the instructions.
  • Provide a brief response immediately without any preamble or formatting markers.

<|eot_id|>

<|start_header_id|>assistant/user<|end_header_id|>{context}<|eot_id|>

<|start_header_id|>assistant<|end_header_id|>
```

Figure 6: Deliberator agent prompt optimized for Llama models.

```
TURN MANAGEMENT PROMPT

{context}

### Instruction:
Given the conversation history, determine the next speaker.

Rules:

  • Choose the most appropriate participant based on the conversation so far.
  • If the user did not propose a topic in the first message, output 'User'.
  • Ensure a balanced participation of speakers where each speaker gets a turn.
  • Options are: {speakers}.
  • Respond with **only** the speaker's name exactly as listed in {speakers}.

### Response:
```

Figure 7: Prompt used to select the next speaker. The list of possible speakers excludes the most recent speaker. This prompt is invoked only when the user has spoken within the last three turns and each model has contributed at least once within the last five turns.

```
SUMMARY PROMPT

{context}

### Instruction:
Provide a concise summary of the discussion in no more than X sentences.
Provide brief and short answer.
Return only the bullet points, each starting with '-', and nothing else.
No need to tell this is summary of conversation or anything else.

### Response:
```

Figure 8: Analyzer prompt — summary.

```
POINTS OF AGREEMENT PROMPT

{context}

### Instruction:
List up to X clearly stated points on which the participants agree.
Provide brief and short answer.
Return only the bullet points, each starting with '-', and nothing else.
No need to tell this is points of agreement of conversation or anything else.

### Response:
```

Figure 9: Analyzer prompt — points of agreement.

## OPEN QUESTIONS PROMPT

{context}

### Instruction:

Identify key questions or issues that remain unresolved or require further discussion.  
Provide brief and short answer.

Return only the bullet points, each starting with '-', and nothing else.

No need to tell this is open questions of conversation or anything else.

### Response:

Figure 10: Analyzer prompt — open questions.

## ARGUMENT MAP PROMPT

{context}

### Instruction:

Construct a structured map of the main arguments discussed.

For each argument, include: A **claim** (the main point being made), one or more **supporting premises** (evidence or reasoning offered for the claim).

Provide brief explanation.

Return only the bullet points, each starting with '-', and nothing else.

No need to tell this is argument map of conversation or anything else.

### Response:

Figure 11: Analyzer prompt — argument map.

# Simultaneous Speech-to-Text Translation Web Application for Estonian

Bohdan Podziubanchuk and Aivo Olev and Jiaming Kong and Tanel Alumäe

Department of Software Science

Tallinn University of Technology

Estonia

{bohdan.podziubanchuk,aivo.olev,tanel.alumae}@taltech.ee

## Abstract

This paper presents a new open-source web application for simultaneous speech-to-text translation. The system translates live Estonian speech into English, Russian, and Ukrainian text, and also supports English-to-Estonian translation. Our solution uses a cascaded architecture that combines streaming speech recognition with a recently proposed LLM-based simultaneous translation model. The LLM treats translation as a conversation, processing input in small five-word chunks. Our streaming speech recognition achieves a word error rate of 10.2% and a BLEU score of 26.1 for Estonian-to-English, significantly outperforming existing streaming solutions. The application is designed for real-world use, featuring a latency of only 3–6 seconds. The application is available at <https://est2eng.cs.taltech.ee>.

## 1 Introduction

Simultaneous speech-to-text translation (SimulST) systems produce real-time text-based translations from streaming speech. The latency of target language words is kept low enough for the listener to follow the speaker without major delay. This task is important in practical settings, for example when supporting talks at conferences with multilingual audience or generating live subtitles that must meet certain latency constraints.

Estonian is a Uralic language spoken by around one million native speakers. Estonia’s growing ethnic and professional diversity, combined with the small size and complexity of the Estonian language, makes many newcomers reluctant to learn it. As a result, English is increasingly used in domains like higher education and technology, raising concerns about potential domain loss, where Estonian could gradually lose its functions in areas like higher education and technology.

This paper describes a simultaneous speech-to-text translation system for Estonian, developed by

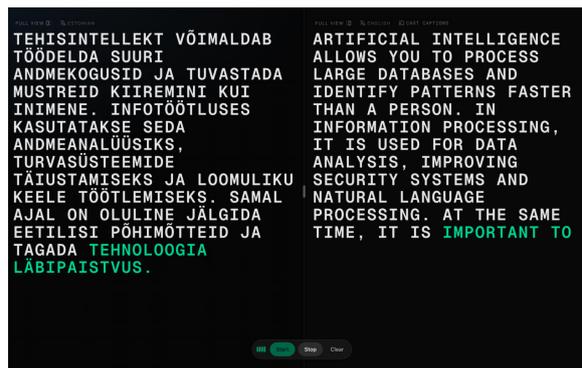


Figure 1: User interface with a split view (ASR and translation) optimized for low-latency incremental updates. Partial ASR and newly emitted translation words are visually marked to support stable live reading.

Tallinn University of Technology and the University of Tartu within a project supported by the Estonian language technology programme. One of the goals of the project was to create a practical, easy-to-use, and accurate system that could be deployed in a variety of real-world settings, including conferences and seminars. The resulting system operates as a web application (see Figure 1) and currently supports translating Estonian speech into English, Ukrainian, and Russian text, as well as English speech into Estonian. The system follows a cascaded architecture: speech in the source language is first transcribed by a streaming automatic speech recognition (ASR) model and then translated by a machine translation (MT) model. Close-to-simultaneous translation is achieved by translating five-word chunks using a large language model (LLM) together with a multi-turn dialogue-based decoding strategy, in which source and target chunks appear interleaved in the translation history (Wang et al., 2025). The LLM is finetuned on parallel training data segmented into small, monotonic chunks according to word alignments, ensuring that no target word appears in an earlier chunk than its aligned source word.

The system freely available and open source, including the training scripts to generate supervision data for finetuning a LLM for multi-turn translation, and can thus be relatively easily adapted to other language pairs. The system comes with several addons that increase its usability: functionality to share live translation results to users’ mobile devices and to stream translation to OBS Studio for overlaying over a video.

A demo video of the application is available at <https://youtu.be/F5bx3Wqyc4Q>.

## 2 Related work

According to our knowledge, the only publicly available simultaneous speech translation model that supports Estonian speech input and Estonian text output is the SeamlessStreaming model (Seamless Communication et al., 2023). SeamlessStreaming is an end-to-end simultaneous multilingual and multimodal translation framework built on the offline speech translation model SeamlessM4T-v2. It performs real-time speech-to-text and speech-to-speech translation for more than 100 input and nearly 100 output languages. Low-latency generation is achieved through Efficient Monotonic Multi-head Attention (EMMA) and additional fine-tuning of the SeamlessM4T-v2 architecture for streaming inference. Its simultaneous text decoder follows a learned policy that decides whether to emit the next token or delay generation in order to read more input context.

Support for streaming Estonian speech translation has also recently been introduced by a few commercial providers, including Microsoft and Google. We were only able to test Microsoft’s system.

## 3 Models

Our SimulST system is based on the cascaded approach, with independent streaming speech recognition and MT models.

### 3.1 Speech recognition

The Estonian streaming ASR system is based on the Zipformer neural transducer architecture (Yao et al., 2024) and was trained using the Icefall toolkit<sup>1</sup>. The model has about 150 million parameters and was trained on roughly 1334 hours of manually transcribed Estonian speech from the TalTech Es-

<sup>1</sup><https://github.com/k2-fsa/icefall>

tonian Speech Dataset 1.0<sup>2</sup> (Alumäe et al., 2023). In addition, training relied on around 4000 hours of automatically transcribed Estonian public TV (ETV) data, consisting of news and talk shows, and a 500-hour subset of the Gigaspeech dataset (Chen et al., 2021), which includes YouTube videos and podcasts. For automatically transcribing the Estonian data, we used Whisper *large-v3-turbo* (Radford et al., 2022), finetuned on the TalTech Estonian Speech Dataset 1.0. The ASR model produces properly capitalized and punctuated text. A subset of Gigaspeech was intermixed with Estonian data to improve the model’s ability to transcribe English terms and expressions that are often embedded into Estonian sentences, especially in technological domains. Since the original transcripts of Gigaspeech are uppercase and not punctuated, we retranscribed the 500 hour subset using Whisper *large-v3-turbo*. The ETV audio remains the broadcaster’s property; licensing details for the derived transcriptions are documented alongside the TalTech Estonian Speech Dataset 1.0.<sup>3</sup>

### 3.2 Machine translation

Our MT component is based on a recently proposed simultaneous MT approach that treats translation as a multi-turn dialogue between the source (as user turns) and the LLM (as assistant turns) (Wang et al., 2025), as illustrated in Figure 2. Instead of injecting new source tokens into the end of a growing prompt – a common workaround when adapting offline MT models for online use – each incoming source chunk is added as a new turn in the conversation. This setup allows the LLM to reuse its key–value cache efficiently, reducing both computational cost and latency. It also enables the use of existing, highly optimized LLM inference tools, since decoding follows the standard multi-turn dialogue pattern. Unlike translation models that integrate a policy network to decide whether to emit more tokens or wait for more input, the MT LLM used here simply finishes each chunk translation with an end-of-text token, after which control returns to the “user” to provide the next chunk.

In order to train LLM to perform such partial translations, the LLM has to be finetuned using specialized supervised fine-tuning data that mim-

<sup>2</sup><https://cs.taltech.ee/staff/tanel.alumae/data/est-pub-asr-data/>

<sup>3</sup><https://cs.taltech.ee/staff/tanel.alumae/data/est-pub-asr-data/>

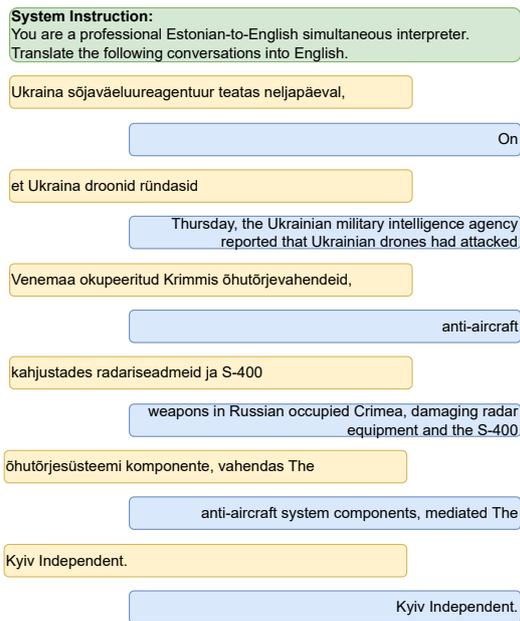


Figure 2: The translation LLM produces a translation for each fixed-word-length input chunk, taking into account the already translated chunks. The LLM sometimes correctly avoids producing the translation too early, if it not sure about the following context.

ics conversational chunked translation. Such data is generated by segmenting parallel sentences using word alignments and converting them into sequences of source and target texts, ensuring that a target word does not occur in a response earlier than the corresponding source word. To make the model robust to different latency settings, the segmented trajectories are further augmented with operations that merge or shift chunks. After training, the LLM can translate incoming partial source text chunk-by-chunk while maintaining coherence using previous conversational turns as context. Experiments by Wang et al. (2025) showed that conversational prompting approaches offline LLM-based translation in quality while substantially reducing latency and is a good alternative to specialized simultaneous MT systems in efficiency.

Our multi-turn simultaneous MT model was finetuned from the existing LLM-based offline MT model Hunyuan-7B-MT<sup>4</sup> (Zheng et al., 2025). Experiments showed that using Hunyuan-7B-MT as the base model yields better results than using more general LLMs, such as Llama 3.1 of similar size.

As training data, we sampled 500K sentence

<sup>4</sup><https://huggingface.co/tencent/Hunyuan-MT-7B>

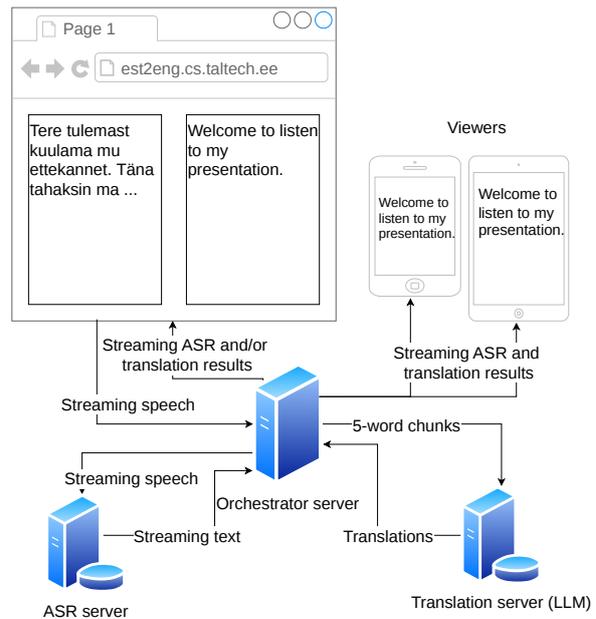


Figure 3: Architecture of the application.

pairs for all supported translation directions (Estonian  $\leftrightarrow$  English, Estonian  $\leftrightarrow$  Ukrainian, and Estonian  $\leftrightarrow$  Russian) from the SynEst corpus (Korotkova et al., 2024). SynEst contains synthetic translations of texts from the Estonian National Corpus (Koppel and Kallas, 2022) into 11 languages, as well as translations from these 11 languages into Estonian, drawing on various web-based sources such as NewsCrawl<sup>5</sup>.

The sampled parallel data was processed as follows<sup>6</sup>. First, we generated word-alignment information using the pretrained multilingual encoder model XLM-R (Conneau et al., 2020). These alignments were then converted into dependency graphs indicating, for each target word, the minimal relevant source-side position. Finally, we transformed the dependency graphs into read-write trajectories, represented as multi-turn chat messages.

We trained Hunyuan-7B-MT using full finetuning for one epoch over the generated dataset. Training took 61 hours on an HPC node equipped with eight AMD MI250X GPUs.

## 4 Architecture

The system is deployed as a web application with a client-server architecture, as shown in Figure 3. The frontend runs in the user’s browser and captures audio from the microphone, while the backend per-

<sup>5</sup><https://data.statmt.org/news-crawl/>

<sup>6</sup>Source code: [https://github.com/jiamingkong/LLM\\_based\\_simulMT](https://github.com/jiamingkong/LLM_based_simulMT)

forms ASR inference and streams transcriptions back to the client in real time.

#### 4.1 System Overview

The frontend is built using Next.js and React, communicating with the backend via a secure WebSocket connection (WSS). Audio is captured using the Web Audio API with an AudioWorklet processor, resampled to 16 kHz, and transmitted as base64-encoded 16-bit PCM chunks at 100 ms intervals. The backend is implemented in Rust using the Axum web framework with Tokio for asynchronous I/O, and is deployed behind an nginx reverse proxy that handles TLS termination and WebSocket upgrades.

#### 4.2 ASR Backend

The ASR backend<sup>7</sup> uses sherpa-onnx<sup>8</sup>, a lightweight inference engine for speech recognition models in ONNX format. We interfaced with sherpa-onnx through Rust FFI bindings provided by the sherpa-rs library. The server maintains a pool of recognizer instances to handle concurrent sessions, with each session mapped to a recognizer using hash-based distribution.

For NeMo-based models such as FastConformer, we exported the models to ONNX format using NeMo’s built-in export functionality. The export process involved configuring the decoder type (CTC or transducer), extracting streaming parameters such as chunk size and cache dimensions, and optionally applying INT8 quantization.

Endpointing in the live pipeline is client-driven: the client signals utterance end via an `utterance_end/stop` event, which triggers `finalize_session` on the server. The current client does not run VAD; silence-based endpointing can be enabled server-side or added client-side in future work.

The system supports multiple translation directions, each requiring a different ASR model. For Estonian input, we use the custom Zipformer model described in Section 3.1, which operates in true streaming mode using neural transducer architecture and produces properly punctuated and capitalized text. For English input (English→Estonian translation), we use NVIDIA’s FastConformer Hybrid model<sup>9</sup>, exported as an RNNT transducer. We

<sup>7</sup>Source: <https://github.com/aivo0/rust-asr-server>

<sup>8</sup><https://github.com/k2-fsa/sherpa-onnx>

<sup>9</sup>[https://huggingface.co/nvidia/stt\\_en\\_](https://huggingface.co/nvidia/stt_en_)

use NeMo models for English because they provide ready-to-export streaming models with ONNX support, while comparable open streaming models for Estonian were not available. Whisper is used for Estonian only as an offline teacher and baseline; it does not provide a streaming ONNX model suitable for our deployment.

#### 4.3 Translation server

Since the translation model is a standard finetuned LLM, it is hosted on a standard GPU-based inference endpoint that uses llama.cpp<sup>10</sup> for model serving. Translation requests are done for every five source words, using a history of 10 previous inputs and outputs. Once a chunk is sent to the LLM, its translation is treated as final and is not revised if upstream ASR hypotheses later change.

#### 4.4 Web Application

The web application<sup>11</sup> captures audio with echo cancellation, noise suppression, and automatic gain control enabled. The application buffers five source words before sending them to the MT model to balance latency against translation quality. Caption sharing across devices is implemented via Firebase Firestore, while OBS Studio integration uses the obs-websocket protocol.

### 5 User interface

The system targets three practical roles: (i) a moderator/speaker who runs the application during a talk, (ii) a technical operator who configures caption casting for OBS Studio, and (iii) audience members who read live captions on mobile devices. The UI is optimized for readability and stable incremental updates under low-latency constraints.

#### 5.1 Main interface for moderators and speakers

The main application page provides simple session controls (Start, Stop, Clear) and language-direction selection. On desktop, the interface uses a two-panel layout with a resizable split (ASR on the left, translation on the right). On mobile, it collapses to a single panel with a left/right toggle.

Operational state is exposed via lightweight cues. A floating ServerStatus widget reports backend availability (e.g., ready/waking up/unreach-

`fastconformer_hybrid_medium_streaming_80ms_pc`

<sup>10</sup><https://github.com/ggml-org/llama.cpp>

<sup>11</sup>Source: <https://github.com/Danbog32/Estonian-to-English-translation>

able). A green “live” indicator shows when microphone/streaming is active. Connection or microphone failures are surfaced as toast notifications, making disruptions immediately visible.

To make recording status immediately recognizable without adding visual noise, we added a minimalistic audio level indicator matching the app’s dark/emerald theme. It consists of four vertical animated bars, is responsive on both desktop and mobile, and is shown only while recording, providing an at-a-glance confirmation that audio capture is active.

## 5.2 Incremental display

To reduce distraction from streaming revisions, the ASR panel combines finalized text with the current partial hypothesis and highlights the most recent words, while earlier lines remain visually stable. The translation panel renders words incrementally, briefly highlighting newly emitted words and then fading the highlight to avoid flicker. Autoscroll follows the stream only while the user stays at the bottom; if the user scrolls up, a “Scroll to bottom” control appears.

Translation is performed in small source-text windows (five-word chunks), providing near-simultaneous output while maintaining local context. The UI does not expose latency–quality controls, keeping the interaction surface minimal.

## 5.3 Sharing and OBS overlay

For audience-scale use, the moderator can enable caption sharing via a “Cast captions” modal. When enabled, the system generates a unique session identifier and provides a shareable link together with a QR code and one-click copy. Viewers open a dedicated reader page that renders incoming captions. This separation of roles keeps the live session controllable: the host can start/stop broadcasting, while viewers remain read-only.

The same modal supports OBS Studio integration for live streams. During operation, captions are published in a broadcast-friendly format by selecting only the most recent words, wrapping them into two lines, and rate-limiting updates to avoid excessive refreshes. This produces compact, stable overlays suitable for conference recordings and live streams.

## 6 Evaluation

In order to assess the performance of the system, we evaluated both ASR and SimulST quality on

|                            | Estonian | English |
|----------------------------|----------|---------|
| Non-streaming (Whisper)    | 10.0     | 16.3    |
| Streaming (Icefall/Nvidia) | 10.2     | 25.2    |

Table 1: WER results for Estonian and English.

dedicated test sets containing real-world broadcast and conversation data.

### 6.1 Data

The evaluation data consists of long broadcast news and conversational recordings with different levels of spontaneity, including press conferences, TV talk shows, YouTube videos, and news programmes featuring many interviews. The total duration of the evaluation dataset is 4 hours for Estonian and 3 hours for English. This material has previously been used for offline speech-translation evaluation (Sildam et al., 2024) and is publicly available<sup>12</sup>.

Reference translations for the evaluation set were created by professional translators in Estonia, using both the audio transcriptions and the audio recordings themselves as inputs.

### 6.2 Speech recognition

Table 1 compares word error rates (WERs) of offline and streaming models on Estonian and English test data. For English, we used Whisper *large-v3-turbo* with voice activity detection, the *-hallucination\_silence\_threshold* parameter set to 2.0, and word-level timestamps enabled, as these settings help reduce hallucinations. For Estonian, we used a version of *large-v3-turbo*<sup>13</sup> finetuned on 1334 hours of Estonian speech with verbatim transcripts.

WERs were computed from Whisper’s long-form decoding output. Because this decoding strategy does not align hypotheses with reference sentences, we computed WERs after removing punctuation, lowercasing both hypotheses and references, and aligning words using minimum WER segmentation (mwerSegmenter) (Matusov et al., 2005) through the SLTev toolkit (Ansari et al., 2021).

Results show that, for Estonian, the streaming ASR quality is very close to offline finetuned Whisper performance. This is expected, as the streaming model is also trained on synthetic transcripts

<sup>12</sup><https://github.com/alumae/k6net6lke-benchmark>

<sup>13</sup><https://huggingface.co/TalTechNLP/whisper-large-v3-turbo-et-verbatim>

produced by Whisper and therefore learns to imitate Whisper’s output. Whisper was also finetuned mostly on the same high quality data as the streaming model. For English, the comparison is between different models: the offline baseline is Whisper *large-v3-turbo*, while the streaming model is NVIDIA’s FastConformer. The gap is therefore expected. The streaming setup is also intentionally minimal (greedy search without language-model rescoring), and endpointing is disabled during streaming, so long segments are decoded continuously, which can increase drift and substitution errors over time.

### 6.3 Speech translation

The streaming translation pipeline uses word-level windowing rather than traditional utterance-based translation. The frontend buffers transcribed words and triggers translation when six words have accumulated since the last emission, extracting five-word chunks for translation. Each chunk is sent to the LLM together with the previous ten source-target pairs as conversational context, enabling consistent terminology and better handling of anaphora across segments. This architecture achieves an end-to-end latency of approximately 3–6 seconds, competitive with professional human interpreters who typically maintain 3–10 second delays. In preliminary ablations, 5–6 word chunks yielded nearly identical translation quality, 4-word chunks degraded quality, and chunk sizes above 6 increased latency beyond acceptable live use; we therefore fixed chunk size to 5. We selected a 10-turn history window to bound latency while preserving local discourse context; larger windows increased response time without clear qualitative gains in pilot use.

SimulST evaluation was based on two metrics: BLEU (Papineni et al., 2002) and BLEURT (Sellam et al., 2020). BLEURT is a learned metric trained on human evaluation scores of translation references and corresponding MT outputs. We used the multilingual BLEURT-20D12 model (Pu et al., 2021). BLEU and BLEURT scores were computed after aligning words in the candidate translations with the references using mwerSegmenter via the SLTev toolkit.

As baselines, we used three systems. The offline cascaded system combines Whisper *large-v3-turbo* for ASR (using the finetuned version for Estonian) with the *Neutotõlge* MT system (Tättar et al., 2022) developed by the NLP research group at the University of Tartu. We accessed the system via its API,

although the corresponding models are also publicly available. We also tested Microsoft Azure’s streaming speech translation API. It should be noted that Azure provides only pseudo-streaming output: although translated words are emitted with low latency as response to streaming speech input, they remain unstable until the end of an utterance-like segment is detected, and both the wording and word order often change when new words are emitted, especially when the current segment is finalized. As the third baseline, we used the Seamless Streaming (Seamless Communication et al., 2023) models to generate translations.

Table 2 presents the Estonian–English, Estonian–Russian, and English–Estonian translation results in terms of BLEU and BLEURT. Among the baseline systems, Azure’s pseudo-streaming translation performs on par with our best offline cascaded setup, while Seamless Streaming shows clearly lower quality.

We evaluated two MT models within our proposed streaming translation architecture: Hunyuan-7B-MT and Llama-3.1-8B-EstLLM<sup>14</sup> (a finetuned Llama-3.1 model with continued pretraining and instruction tuning on mostly Estonian data). The results show that finetuning a dedicated translation-oriented LLM gives better translation quality than using a language-specific but task-agnostic LLM. Our best cascaded system outperforms Seamless Streaming by a large margin and is approximately 5 BLEU points lower than the best open cascaded offline system in all translation directions.

Evaluation results suggest that translations from Estonian achieve higher quality than translations into Estonian. This is mainly due to the relatively weak English ASR model used in our current setup. For our main translation direction – Estonian to English – the quality is already sufficient for use at real live events.

### 6.4 Limitations

We did not include traditional text-only SimulMT baselines such as fixed wait-k policies. A fair comparison would require an end-to-end setup that maps streaming ASR timestamps to a text-level SimulMT policy and evaluates both latency and stability under identical real-time constraints. Implementing and validating such a baseline is left for future work.

<sup>14</sup><https://huggingface.co/tartuNLP/llama-estllm-prototype-0825>

|                                         | <i>et</i> → <i>en</i> |        | <i>et</i> → <i>ru</i> |        | <i>en</i> → <i>et</i> |        |
|-----------------------------------------|-----------------------|--------|-----------------------|--------|-----------------------|--------|
|                                         | BLEU                  | BLEURT | BLEU                  | BLEURT | BLEU                  | BLEURT |
| <b>Baselines</b>                        |                       |        |                       |        |                       |        |
| Offline cascaded (Whisper + Neurotõlge) | 31.9                  | 0.60   | 26.6                  | 0.61   | 16.7                  | 0.47   |
| Azure (pseudo-streaming)                | 31.4                  | 0.56   | 17.4                  | 0.53   | 20.4                  | 0.53   |
| Seamless Streaming                      | 13.8                  | 0.36   | 9.5                   | 0.29   | 8.8                   | 0.27   |
| <b>Our cascaded streaming systems</b>   |                       |        |                       |        |                       |        |
| ASR + Llama-3.1-8B-EstLLM ft.           | 24.8                  | 0.54   | 21.3                  | 0.54   | 9.8                   | 0.34   |
| ASR + Hunyuan-7B-MT ft.                 | 26.1                  | 0.56   | 22.3                  | 0.56   | 11.1                  | 0.37   |

Table 2: Simultaneous translation evaluation results of baseline models and our cascaded systems. The system corresponding to the last row is used in the live system.

## 7 Conclusion

This work demonstrates that high-quality, low-latency simultaneous speech-to-text translation for Estonian is now technically feasible using an open, deployable, and reproducible system. By combining a strong streaming ASR model with a conversationally prompted, chunk-based LLM translator, we show that a cascaded architecture can approach offline translation quality while keeping latency within 3–6 seconds, suitable for real-world conference and seminar scenarios. Our work provides an example for building similar systems for other languages, including methods for generating multi-turn supervision data and an openly available web application ready for real-world use.

The most urgent future work is replacing the current English ASR model with a stronger one, thereby improving the English-to-Estonian translation direction.

## Acknowledgments

This research was supported by the Estonian Centre of Excellence in AI (EXAI), National Program for Estonian Language Technology Program (projects EKTB77 and EKTB117), both funded by the Estonian Ministry of Education and Research, and by the Estonian Language Data Research Infrastructure (KeTA).

## References

Tanel Alumäe, Joonas Kalda, Külliki Bode, and Martin Kaitsa. 2023. [Automatic closed captioning for Estonian live broadcasts](#). In *Proceedings of the 24th Nordic Conference on Computational Linguistics*

*tics (NoDaLiDa)*, pages 492–499, Tórshavn, Faroe Islands. University of Tartu Library.

Ebrahim Ansari, Ondřej Bojar, Barry Haddow, and Mohammad Mahmoudi. 2021. [SLTEV: Comprehensive evaluation of spoken language translation](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 71–79, Online. Association for Computational Linguistics.

Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Zhao You, and Zhiyong Yan. 2021. [GigaSpeech: An evolving, multi-domain ASR corpus with 10,000 hours of transcribed audio](#). In *Interspeech 2021*. ISCA.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). *Preprint*, arXiv:1911.02116.

Kristina Koppel and Jelena Kallas. 2022. [Eesti keele ühendkorpuste sari 2013–2021: mahukaim eestikeelsete digitekstide kogu](#). *Eesti Rakenduslingvistika Ühingu aastaraamat = Estonian Papers in Applied Linguistics*, 18:207–228.

Elizaveta Korotkova, Taido Purason, Agnes Luhtaru, and Mark Fishel. 2024. [Multilinguality or back-translation? A case study with Estonian](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11838–11848, Torino, Italia. ELRA and ICCL.

Evgeny Matusov, Gregor Leusch, Oliver Bender, and Hermann Ney. 2005. [Evaluating machine translation output with automatic sentence segmentation](#). In *Proceedings of the Second International Workshop on*

- Spoken Language Translation*, Pittsburgh, Pennsylvania, USA.
- 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.
- Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. [Learning compact metrics for MT](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 751–762, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. [Robust speech recognition via large-scale weak supervision](#). *Preprint*, arXiv:2212.04356.
- Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, John Hoffman, Min-Jae Hwang, Hirofumi Inaguma, Christopher Klaiber, Iliia Kulikov, Pengwei Li, Daniel Licht, Jean Maillard, and 46 others. 2023. [Seamless: Multilingual expressive and streaming speech translation](#). *Preprint*, arXiv:2312.05187.
- 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.
- Tiiia Sildam, Andra Velve, and Tanel Alumäe. 2024. [Finetuning end-to-end models for Estonian conversational spoken language translation](#). In *Proceedings of the Seventh Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2024)*, pages 166–174, Bangkok, Thailand. Association for Computational Linguistics.
- Andre Täatar, Taido Purason, Hele-Andra Kuulmets, Agnes Luhtaru, Liisa Rätsep, Maali Tars, Marcis Pinnis, Toms Bergmanis, and Mark Fishel. 2022. Open and competitive multilingual neural machine translation in production. *Proceedings of Baltic HLT*.
- Minghan Wang, Thuy-Trang Vu, Yuxia Wang, Ehsan Shareghi, and Gholamreza Haffari. 2025. [Conversational SimulMT: Efficient simultaneous translation with large language models](#). In *Proceedings of the 22nd International Conference on Spoken Language Translation (IWSLT 2025)*, pages 93–105, Vienna, Austria (in-person and online). Association for Computational Linguistics.
- Zengwei Yao, Liyong Guo, Xiaoyu Yang, Wei Kang, Fangjun Kuang, Yifan Yang, Zengrui Jin, Long Lin, and Daniel Povey. 2024. [Zipformer: A faster and better encoder for automatic speech recognition](#). *Preprint*, arXiv:2310.11230.
- Mao Zheng, Zheng Li, Bingxin Qu, Mingyang Song, Yang Du, Mingrui Sun, and Di Wang. 2025. [Hunyuan-MT technical report](#). *Preprint*, arXiv:2509.05209.

# The AI Committee: A Multi-Agent Framework for Automated Validation and Remediation of Web-Sourced Data

Sunith Vallabhaneni<sup>1,2,3</sup>, Thomas Berkane<sup>2,3</sup>, and Maimuna Majumder<sup>2,3</sup>

<sup>1</sup>UC Berkeley

<sup>2</sup>Harvard Medical School

<sup>3</sup>Boston Children’s Hospital

sunithv@berkeley.edu, {Thomas.Berkane, Maimuna.Majumder}@childrens.harvard.edu

## Abstract

Many research areas rely on data from the web to gain insights and test their methods. However, collecting comprehensive research datasets often demands manually reviewing many web pages to identify and record relevant data points, which is labor-intensive and susceptible to error. While the emergence of large language models (LLM)-powered web agents has begun to automate parts of this process, they often struggle to ensure the validity of the data they collect. Indeed, these agents exhibit several recurring failure modes—including hallucinating or omitting values, misinterpreting page semantics, and failing to detect invalid information—which are subtle and difficult to detect and correct manually. To address this, we introduce the **AI Committee**, a novel model-agnostic multi-agent system that automates the process of validating and remediating web-sourced datasets. Each agent is specialized in a distinct task in the data quality assurance pipeline, from source scrutiny and fact-checking to data remediation and integrity validation. The AI Committee leverages various LLM capabilities—including in-context learning for dataset adaptation, chain-of-thought reasoning for complex semantic validation, and a self-correction loop for data remediation—all without task-specific training. We demonstrate the effectiveness of our system by applying it to three real-world datasets, showing that it generalizes across LLMs and significantly outperforms baseline approaches, achieving data completeness up to 78.7% and precision up to 100%. We additionally conduct an ablation study demonstrating the contribution of each agent to the Committee’s performance. This work is released as an open-source tool for the research community.

## 1 Introduction

The development of robust models for core NLP tasks like event extraction and knowledge base population is fundamentally dependent on

high-quality, structured data. While benchmark datasets have been invaluable, they are often static and cannot capture the dynamics of emergent, real-world events. This is a critical limitation for timely research in fields such as computational social science, public health, and economics, which rely on monitoring current affairs—from social unrest to disease outbreaks.

For such domains, the web is the essential, albeit noisy and unstructured, source of information. However, collecting comprehensive research datasets often demands manually reviewing many web pages to identify and record relevant data points, a process that is labor-intensive and susceptible to error. While the emergence of Large Language Model (LLM)-powered web agents has begun to automate parts of this process, they often struggle to ensure the validity of the data they collect. Existing agents are prone to subtle yet critical errors, ranging from the hallucination of plausible-sounding values to the misinterpretation of complex web page semantics. This forces researchers to validate each collected data point, undermining the initial promise of automation.

To address this, we introduce the **AI Committee (AIC)**, a novel model-agnostic multi-agent system that automates the process of validating and remediating web-sourced data. This framework can be applied both to human-collected data and, most importantly, to data gathered by autonomous AI web agents, a rapidly growing source of data. Our system is designed as a digital assembly line, where each agent acts as a specialist with a narrowly defined role, mimicking a human-led data curation team.

Our framework integrates various in-context learning techniques, including few-shot learning for rapid dataset adaptation and structured chain-of-thought reasoning for complex semantic validation. For each dataset, context agents dynamically tailor AIC’s validation rules before execut-

ing a comprehensive pipeline for each data point. This pipeline includes a source scrutinizer to assess source reliability, a fact-checker equipped with detailed, context-aware guidelines, and specialized agents to ensure data points are both relevant and structurally sound. A key feature of the system is its structured self-correction and discovery mechanism, where dedicated remediation agents correct invalid data points and discover new, relevant information from the source web pages. Additionally, AIC is model-agnostic and requires no model training, as it performs well across all models regardless of complexity—making it well-suited for low-resource academic/research settings.

To demonstrate the effectiveness of this framework, we test it on three datasets collected from the web by an LLM-based tool and use carefully human-cleaned versions as ground truth to evaluate accuracy, precision, and recall across several foundational models while analyzing the trade-off between performance, cost, and time. Our system significantly improves data quality across multiple state-of-the-art LLMs. This framework has the potential to accelerate the curation of high-quality datasets, empowering researchers to conduct more timely and impactful research.

## 2 Related Work

### 2.1 Multi-Agent Systems for Web Data Collection

Multi-Agent Systems (MAS) are rapidly emerging as the dominant paradigm for automating the complex, end-to-end task of web data collection. Recent state-of-the-art frameworks such as that of Berkane et al. (2025) and AutoData (Ma et al., 2025) exemplify this trend. Both systems aim to transform a single natural language instruction into a final, structured dataset. AutoData employs a sophisticated architecture of 'research' and 'development' squads to programmatically generate and execute scraping code, featuring a novel hypergraph cache system to optimize inter-agent communication. Similarly, the work by Berkane et al. (2025) details an end-to-end, human-in-the-loop pipeline that automates query generation, web page retrieval with bias mitigation, and data extraction. Other related systems like AutoScraper (Huang et al., 2024) focus on a similar problem domain, generating reusable wrappers (scrapers) for specific websites through a progressive, multi-stage process.

However, a critical analysis of these collection-focused systems reveals a common gap in the final, crucial stage: automated quality control and remediation. While they are powerful engines for initial data extraction, their approach to data quality largely remains a supervised or manual task. For instance, the framework by Berkane et al. (2025) concludes its pipeline by using an LLM to flag potential data point issues for subsequent manual user review. While AutoData and AutoScraper include validation steps, their focus is on the integrity and reusability of the generated program or scraper, rather than the deep semantic validation and automated correction of each individual extracted data point against its original source.

### 2.2 LLMs for Data Quality and Cleaning

There is a growing body of research on leveraging LLMs for data quality tasks such as data imputation, error detection, and schema mapping. For instance, Zhang et al. (2024) explored the potential for LLMs to automate semantic data cleaning tasks—such as correcting inconsistent string representations—that were traditionally performed by complex regex rules or crowdsourcing. Moving beyond static cleaning, Bendinelli et al. (2025) investigated the use of LLM agents equipped with code execution tools (IPython) to iteratively detect errors and clean tabular datasets, specifically optimizing for downstream machine learning model performance.

While powerful, these approaches often focus on structural consistency or statistical optimization within a closed dataset. Our work differentiates itself by decomposing the complex problem of web-sourced data quality—which requires external validation against source text—into a series of distinct sub-tasks, each handled by an agent with a highly specialized prompt and purpose. This specialization, we argue, leads to more robust and interpretable outcomes than a single, monolithic approach.

### 2.3 Prompt Engineering and In-Context Learning

The performance of LLMs is heavily dependent on the quality of the prompt. Advanced techniques like chain-of-thought (CoT) prompting, which encourages the model to “think step-by-step,” have been shown to significantly improve performance on complex reasoning tasks (Wei et al., 2022). Our FactCheckerAgent employs a highly

structured prompt that not only enforces CoT reasoning but also incorporates a “Critical Semantic Audit” to actively check for logical fallacies. Furthermore, our framework’s ContextGenerator embodies the principles of in-context learning, first popularized by Brown et al. (2020). By analyzing a few data samples, these agents generate dynamic rules and tailored examples that are injected into the prompts of other agents, effectively conditioning the system’s behavior for a specific dataset without updating model weights.

## 2.4 Self-Correction and Automated Remediation

A key frontier in LLM research is enabling models to reflect on and correct their own outputs. Many self-correction methods involve a simple feedback loop where the model critiques its own response. The AI Committee implements a more sophisticated, *structured self-correction mechanism*. When the ArbitratorAgent rejects a data point, the reason for rejection is passed to a dedicated DataRemediationAgent, a process inspired by self-refinement techniques described by Madaan et al. (2023). This agent’s sole purpose is to diagnose the failure and perform a targeted correction, creating a rigorous, multi-step validation loop that is more structured than a simple self-critique.

## 3 The AI Committee

The AI Committee is a modular, asynchronous framework designed to clean either human or LLM web-sourced data points through a multi-stage validation and remediation pipeline. The overall architecture is depicted in Figure 1, illustrating the flow of each data point through a series of specialized agents. To operate, the framework requires three user-provided inputs: 1) the initial tabular dataset of web-sourced data points with an associated source URL column for each data point, pointing to the web page from which the datapoint originates, 2) a high-level, natural language description of the dataset’s purpose (e.g., “natural disaster events in Haiti and Cameroon”), and 3) a list of the schema — the dataset’s columns — to be validated (e.g., ‘event\_type’, ‘date’). These inputs form the basis for the Committee’s context-aware processing.

## 3.1 System Initialization and Adaptation

Before processing any data, the framework undergoes an initialization phase to adapt itself to the specificities of the dataset it will operate on. This step employs few-shot learning principles: the framework provides the LLM with a small number of example datapoints directly within its prompt, instructing it to generalize from these examples to create a tailored operational context.

The following components are involved in this phase:

**ContextGenerator** This agent is responsible for creating the framework’s operational context. It first analyzes a random subset of 10 rows of the input data to infer its semantic properties, mapping each schema field to: (1) a description of its expected entity type and granularity (e.g., state: ‘geographic regions at the U.S. state level’); (2) a description of the temporal context (e.g., the date an event occurred); and (3) a list of negative examples to avoid (e.g., extracting a city for a state field). The number of negative examples is determined by the LLM’s analysis of a field’s complexity; a field with high ambiguity (e.g., a free-text event\_type) may warrant more examples than a well-defined one (e.g., country). Building on this analysis, the agent then generates richer, pedagogical illustrations of logical fallacies, providing the later **Fact Checker agent** with a deeper, context-specific understanding of not just what to avoid, but why.

**Dynamic Schema Generation.** From the list of user-provided schema fields, the framework programmatically generates a structured data format. For example, for a dataset tracking corporate acquisitions with fields `acquiring_company`, `deal_value`, and `date`, the system generates a schema requiring each data point to have an `acquiring_company` (string), a `deal_value` (integer), and a `date` (string, formatted as YYYY-MM-DD). This serves as a structure that all agents must adhere to when creating or modifying data. For instance, any new data point discovered by the **Data Remediation agent** must conform to this schema. In this context, structurally consistent is defined by the ability of a data point to be successfully parsed by this model, ensuring all required fields are present.

## 3.2 Core Validation Pipeline

Each data point from the input dataset is processed through the core validation pipeline, which

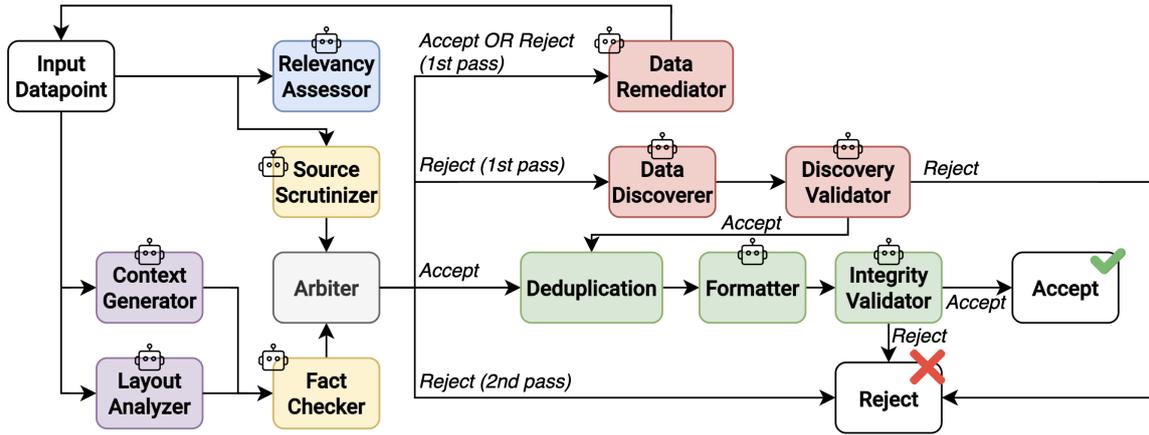


Figure 1: The data processing pipeline of AIC. Data points flow through a series of validation and remediation agents. Each box represents a step of the pipeline.

functions as a digital assembly line.

**Relevancy Assessor Agent.** This agent performs an initial, low-cost check to determine if the data point is topically relevant to the dataset description. This check operates solely on the data point and the dataset description, occurring before any time-consuming web crawling or expensive content analysis is initiated. This prevents the system from wasting resources on fundamentally incorrect entries, as the data point will be rejected if it is found not to be relevant.

**Content Retrieval and Analysis.** For each relevant data point, the **LayoutAnalyzerAgent** uses Crawl4AI to parse the page layout into structured markdown. It then analyzes this markdown to classify the page’s structure into one of several predefined categories, including ‘ARTICLE’, ‘DIRECTORY\_LISTING’, ‘SEARCH\_RESULTS’, ‘HOMEPAGE’, ‘ERROR\_PAGE’, or ‘OTHER’. The classification is used to provide a crucial analysis hint to the downstream Fact Checker agent.

**Source Scrutinizer Agent.** Concurrently with content retrieval, this agent analyzes the source URL to assess its trustworthiness, and then categorizes the source type and assigns a reliability score based on domain reputation.

**Fact Checker Agent.** As the centerpiece of the validation pipeline, this agent is responsible for core semantic validation. Its prompt is dynamically assembled using the outputs from several upstream agents to provide rich context: it receives an analysis hint from the LayoutAnalyzerAgent, which provides a dynamically generated instruction in the agent’s prompt to tailor its reading strategy. (eg: if a page is classified as a ‘DIRECTORY\_LISTING’,

the hint explicitly warns the agent not to dismiss the page as purely navigational and to treat the text within list items as potential data points), a set of dataset-specific rules and tailored fallacy examples from the ContextGenerator. These dynamic inputs are combined with a prompt that includes a Critical Semantic Audit section that defines a set of logical principles hard-coded into the framework. This process involves extracting full entities instead of partial fragments, matching the data’s granularity to the dataset’s needs, analyzing qualifying terms, and identifying matches based on underlying meaning rather than just specific keywords. The agent’s output is a structured ‘FactCheckerResponse’ object, which contains boolean flags for content validity and factual accuracy, any newly extracted date, and detailed explanatory notes justifying its reasoning.

**Arbiter Agent.** Receives the structured reports from the FactCheckerAgent and the SourceScrutinizerAgent. It functions not as an LLM call, but as a deterministic, rule-based decision engine. Evidence is collated through a series of logical checks: a data point is immediately rejected if the FactCheckerAgent reports that the source text lacks meaningful content or does not support the data point’s claims. It is also rejected if the SourceScrutinizerAgent deems the source’s reliability to be ‘Low’ or ‘Very Low’. Only if a data point passes all of these checks does the Arbiter issue a verdict of ‘ACCEPT’.

**Data Formatter Agent.** Upon receiving an ‘ACCEPT’ verdict, this agent performs a final structural normalization pass. It ensures the data point strictly adheres to the generated schema types and formats (e.g., coercing stringified numbers to

integers, enforcing ISO date standards) before the data is passed to the finalization stage.

### 3.3 Remediation and Discovery: A Structured Self-Correction Loop

If the `ArbiterAgent` rejects a datapoint, it is not immediately discarded. Instead, it enters the remediation loop, which enables correction and data discovery.

**Data Remediation Agent.** This agent attempts to fix the rejected data point. It is capable of two types of fixes: direct value replacement (e.g., correcting an incorrect city name that is explicitly mentioned in the text) and calculation-based remediation. The process begins with an "Analyst" LLM call, which receives the rejected data point, the reason for rejection, and the source text, and in turn produces a structured remediation plan. If this plan requires external information for a calculation (e.g., finding a state's total population to apply a percentage found in the text), it can dispatch a fact-lookup tool. This tool queries Google, scrapes the top results, and uses an LLM to extract the required information. If the remediation is successful, the newly corrected datapoint is validated by a Remediation Auditor—a modified Fact Checker specialized in verifying correction logic—before being finalized. If it fails, the datapoint is permanently rejected.

**Data Discovery Agent.** This agent scans the entire source page to discover if other, different data points matching the schema are present. This allows the framework to augment the dataset with new relevant information from the same sources.

### 3.4 Finalization

Datapoints that are accepted, either initially or after remediation, pass through two final stages:

**Hierarchical Deduplication.** A deterministic algorithm performs a final cleaning of the accepted datapoints. The process begins by filtering out any record with missing values in the required schema fields. If a 'date' field is present, the algorithm creates a "base fingerprint" for each datapoint from all non-date fields. Based on that fingerprint, it then deduplicates entries by checking for this fingerprint at three distinct levels of date granularity: full date (YYYY-MM-DD), year-month, and year-only. The levels can coexist because a record is only checked against the set corresponding to its own precision; for instance, accepting a record with month-level precision does not preclude ac-

cepting a different record for the same event with day-level precision. This allows otherwise identical records with differing date specificities to coexist. If no 'date' field is in the schema, a standard deduplication across all fields is performed instead, which identifies duplicates based on identical values across all user-defined schema fields and retains only the first unique instance encountered.

**Data Integrity Validation Agent.** This agent performs a last-pass quality check on the deduplicated data. It is explicitly prompted to validate each data point against only two simple, hard-coded rules: (1) **Completeness:** Is any required schema field empty, null, or missing? While the preceding **Hierarchical Deduplication** step also filters for missing fields, it does so merely as a structural prerequisite to enable its algorithm. This agent's check, in contrast, serves as a final quality control on the deduplicated output. (2) **Plausibility:** Does any field contain a value that is obvious nonsense given its description (e.g., a numerical value for a field requiring a person's name)? Data points that fail this final check are discarded.

## 4 Demonstration Description

We demonstrate AIC via an interactive web interface built with Streamlit, designed to provide a transparent view into the autonomous validation process. The demonstration highlights the system's end-to-end functionality, from raw data ingestion to a cleaned, validated output.

The user workflow goes through the following steps: **Configuration:** The user uploads a CSV file containing web-sourced data points, provides a natural language description of the dataset, specifies the schema fields or columns to be validated, and inserts their OpenAI API key. **Execution:** Upon starting the process, the system begins processing each data point through the multi-agent pipeline. **Real-Time Monitoring:** The interface displays a real-time status log, as shown in Figure 2. Each data point is listed with its current status (e.g., Processing, ACCEPT, REJECT, DISCOVERED), allowing the user to observe the Committee's decisions as they happen. **Results:** Once the pipeline completes, the final, validated dataset is made available for download, and a summary of the performance and cost metrics is displayed.

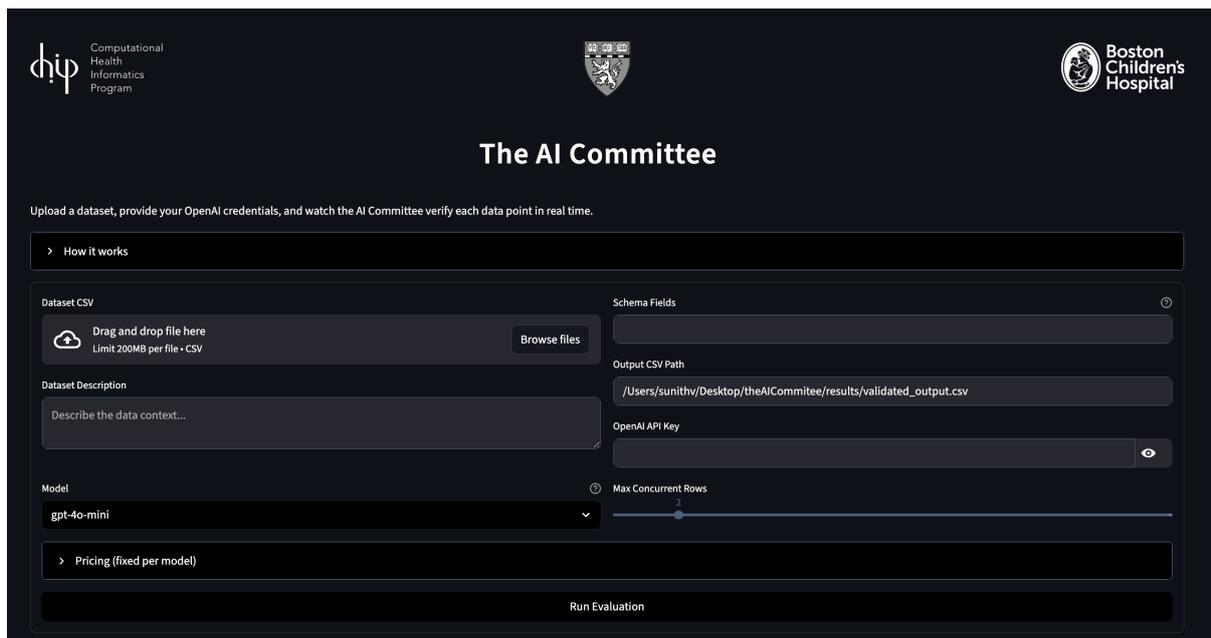


Figure 2: The web interface of AIC. Users provide a dataset, description, and schema, then monitor the real-time validation status of each data point as the system runs.

## 5 Evaluation, Availability & Licensing

We conduct a comprehensive empirical evaluation to validate the effectiveness and analyze the operational characteristics of the AIC. Our experiments, orchestrated by a reproducible evaluation harness, are designed to measure performance across multiple dimensions, from data quality and computational cost to the granular performance of each agent in the pipeline.

### 5.1 Datasets

To establish a high-quality ground truth, each dataset was independently annotated by three public health researchers following a detailed inter-annotator agreement (IAA) protocol, which will be released alongside our code.

1. **Natural Disasters:** 125 initial data points (i.e., LLM-collected) related to natural disasters in Haiti and Cameroon.
2. **Police Misconduct:** 95 initial data points on police misconduct incidents in the United States.
3. **COVID-19 Tracing:** 83 initial data points on the number of downloads of COVID-19 contact tracing applications per state in the US.

### 5.2 Systems Under Test

To rigorously evaluate the framework, we compare several architectures and model configurations. **The AI Committee (AIC):** The modular, multi-agent framework as described in Section 3. We evaluate the Committee using three different underlying models to test generalizability and cost-efficiency: gpt-4o-mini, o4-mini, and gpt-5. **Monolithic Agent Baseline:** A single agent powered by the same models, provided with the full schema, dataset description, and a complex, all-in-one prompt. This baseline serves to demonstrate the specific value of decomposing the task into specialized agents versus a standard "zero-shot" or "few-shot" LLM approach. **Rule-Based Baseline:** A deterministic system that applies a series of hard-coded heuristics, regex patterns, and data cleaning rules specific to each dataset. This represents traditional, non-LLM data cleaning methods. **Ablation Configurations:** We systematically disable key agents (e.g., Relevancy Assessor, Fact Checker, Remediation) within the AI Committee to quantify the individual contribution of each component to the overall system performance.

### 5.3 Results

Our results demonstrate that the AIC framework significantly optimizes the trade-off between data quality, computational cost, and latency. Most notably, the AIC configuration powered by

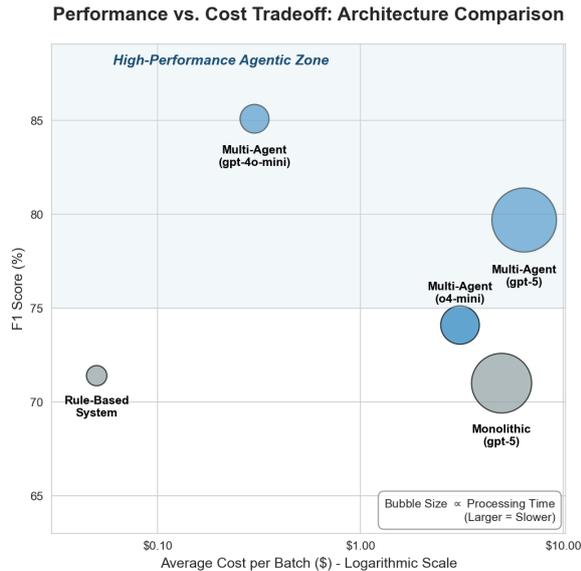


Figure 3: Performance vs Cost Tradeoff. Bubble size represents total processing time.

gpt-4o-mini achieved the highest overall performance (**F1: 85.1**), surpassing both larger, more expensive models gpt-5 and o4-mini.

Figure 3 visualizes this landscape. gpt-4o-mini dominates the “high-performance, low-cost” (top-left) quadrant, delivering a **14.1 point F1 improvement** over the Monolithic baseline while reducing operational costs by over  $16\times$  (\$0.30 vs \$4.95). Conversely, while the gpt-5-powered Committee achieved high precision (100.0%), it had significantly higher cost and latency.

Table 1 details these metrics. The modular architecture consistently outperforms the Monolithic baseline, which struggled with recall (58.9%), likely due to attention dispersion over long contexts. The Rule-based baseline, while precise, failed to generalize to complex unstructured text, resulting in the lowest recall (58.2%) of any valid configuration. The ablation study further validates our design: removing key components—such as the Remediation Agent or Fact Checker—resulted in measurable drops in F1 score (-6.4 and -1.3, respectively), confirming that high data quality requires a collaborative, multi-stage system rather than a single pass.

## 5.4 Future Work

Future work will focus on two key areas. First, we plan to conduct a more thorough benchmark of available models from different families (e.g. Claude, Gemini, etc.) and explore hybrid

Table 1: **Architecture Evaluation:** F1, P=Precision, R=Recall, Rem=Remediation Recall.  $\Delta$  values show percentage point change (Absolute Difference) from the AIC (4o-mini) baseline. Operational metrics (Time, Latency, Cost) are averaged per dataset batch.

| System           | Performance ( $\Delta$ ) |              |              |              | Operations (Avg) |       |        |
|------------------|--------------------------|--------------|--------------|--------------|------------------|-------|--------|
|                  | F1                       | P            | R            | Rem          | Time             | Lat   | Cost   |
| AIC (4o-mini)    | 85.1 (-)                 | 92.7 (-)     | 78.7 (-)     | 77.8 (-)     | 1,285s           | 11.7s | \$0.30 |
| AIC (o4-mini)    | 74.1 (-11.0)             | 92.7 (-)     | 61.7 (-17.0) | 61.1 (-16.7) | 3,128s           | 27.0s | \$3.09 |
| AIC (gpt-5)      | 79.7 (-5.4)              | 100.0 (+6.3) | 66.7 (-12.0) | 44.4 (-33.4) | 9,258s           | 88.9s | \$6.40 |
| Monolith (gpt-5) | 71.0 (-14.1)             | 89.5 (-3.2)  | 58.9 (-19.8) | 58.3 (-19.5) | 7,632s           | 75.9s | \$4.95 |
| Rule-Based       | 71.4 (-13.7)             | 92.4 (-0.3)  | 58.2 (-20.5) | 55.6 (-22.2) | -                | -     | -      |

Table 2: **Ablation Study Metrics:** P=Precision, R=Recall, F1=F1 Score, Rem=Remediation Recall.  $\Delta$  values show percentage point change (Absolute Difference) from the AIC (4o-mini) baseline.

| System         | P ( $\Delta$ ) | R ( $\Delta$ ) | F1 ( $\Delta$ ) | Rem ( $\Delta$ ) |
|----------------|----------------|----------------|-----------------|------------------|
| AIC (4o-mini)  | 92.7 (0.0)     | 78.7 (0.0)     | 85.1 (0.0)      | 77.8 (0.0)       |
| No FactCheck   | 92.4 (-0.3)    | 76.6 (-2.1)    | 83.8 (-1.3)     | 58.3 (-19.5)     |
| No Context     | 93.0 (+0.3)    | 74.5 (-4.2)    | 82.7 (-2.4)     | 75.0 (-2.8)      |
| No CtxExamples | 95.3 (+2.6)    | 71.6 (-7.1)    | 81.8 (-3.3)     | 63.9 (-13.9)     |
| Rem-Only       | 95.3 (+2.6)    | 70.9 (-7.8)    | 81.3 (-3.8)     | 66.7 (-11.1)     |
| No Integrity   | 90.5 (-2.2)    | 73.8 (-4.9)    | 81.3 (-3.8)     | 66.7 (-11.1)     |
| No SrcScrutiny | 91.4 (-1.3)    | 73.0 (-5.7)    | 81.2 (-3.9)     | 66.7 (-11.1)     |
| No CtxLearning | 92.8 (+0.1)    | 71.6 (-7.1)    | 80.8 (-4.3)     | 66.7 (-11.1)     |
| No Remediation | 95.9 (+3.2)    | 66.7 (-12.0)   | 78.7 (-6.4)     | 58.3 (-19.5)     |
| No Layout      | 94.1 (+1.4)    | 67.4 (-11.3)   | 78.5 (-6.6)     | 61.1 (-16.7)     |
| Discovery-Only | 96.0 (+3.3)    | 66.0 (-12.7)   | 78.2 (-6.9)     | 58.3 (-19.5)     |
| Min FactCheck  | 83.6 (-9.1)    | 65.2 (-13.5)   | 73.3 (-11.8)    | 58.3 (-19.5)     |
| No Relevancy   | 78.9 (-13.8)   | 68.1 (-10.6)   | 73.1 (-12.0)    | 66.7 (-11.1)     |
| No Formatter   | 92.1 (-0.6)    | 41.1 (-37.6)   | 56.9 (-28.2)    | 50.0 (-27.8)     |

approaches, using different models for different agents to optimize the cost-performance trade-off (e.g., a fast, cheap model for relevancy screening and a powerful, expensive model for fact-checking). Second, we will continue to refine the agents’ reasoning capabilities through prompt engineering, particularly in handling ambiguity and complex temporal expressions, fixing edge cases as we find them in our testing and user feedback.

## 5.5 Availability & Licensing

The AI Committee is released as an open-source project under the Apache 2.0 license. The complete codebase, including our comprehensive evaluation, the three annotated ground truth datasets, and the inter-annotator agreement protocol, is available at:

**Code Repository:** <https://github.com/sunith-v/theAICommitteeDemo>

**Video:** <https://youtu.be/c4xI9F1s24E>

## 6 Broader Impact, Limitations, and Ethics

Our framework, while powerful, has several limitations regarding its reliability and potential biases in deployment. **Dependence on LLM Capabilities:** The ultimate performance is bounded by the reasoning and knowledge capabilities of the underlying LLM. Hallucinations or misinterpretations by the LLM can still lead to errors. **Prompt Optimization Bias:** We acknowledge that the prompts employed across our agents were iteratively developed and refined using gpt-4o-mini as the primary testbed. Due to this, the superior cost-performance ratio observed for this specific model may result from prompt-model alignment, where instructions are inadvertently optimized for its specific reasoning patterns. Other models, such as gpt-5 or o4-mini, likely achieve higher performance ceilings if the prompts were specifically re-tuned for their respective attention mechanisms and instruction-following quirks. **The Messiness of Web Data:** The framework is still susceptible to the inherent “messiness” of the web. Paywalls, complex JavaScript-rendered pages, and highly unconventional metadata can impede content retrieval and analysis. **Ground Truth Subjectivity:** Our evaluation is based on a human-curated ground truth which, particularly for date fields, involved a degree of subjective leniency due to the ambiguous nature of dates in web contexts. A different ground truth definition could alter the perceived accuracy of the models. **Reliance on LLM’s Internal Knowledge:** A potential limitation of our framework is the SourceScrutinizerAgent’s reliance on the LLM’s internal knowledge to assess source reliability, which carries an inherent risk of hallucination. We argue, however, that this risk is minimal for this specific task. The agent’s function—classifying major domains like news outlets or government portals—relies on stable, high-consensus knowledge that is deeply encoded in the training data of any large foundation model. Furthermore, our empirical results support this design choice: in our experiments across all three datasets and models, we observed no instances where the agent hallucinated or mischaracterized a source’s type or reputation.

## References

- Tommaso Bendinelli, Artur Dox, and Christian Holz. 2025. Exploring LLM agents for cleaning tabular machine learning datasets. *arXiv preprint arXiv:2503.06664*. Proceedings of the ICLR 2025 Workshop on Foundation Models in the Wild.
- Thomas Berkane, Marie-Laure Charpignon, and Maimuna S. Majumder. 2025. LLM-based web data collection for research dataset creation. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 12610–12622. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, and 1 others. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- Wenhao Huang, Zhouhong Gu, Chenghao Peng, Zhixu Li, Jiaqing Liang, Yanghua Xiao, Liqian Wen, and Zulong Chen. 2024. AutoScraper: A progressive understanding web agent for web scraper generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2371–2389.
- Tianyi Ma, Yiyue Qian, Zheyuan Zhang, Zehong Wang, Xiaoye Qian, Feifan Bai, Yifan Ding, Xuwei Luo, Shinan Zhang, Keerthiram Murugesan, Chuxu Zhang, and Yanfang Ye. 2025. AutoData: A multi-agent system for open web data collection. *arXiv preprint arXiv:2505.15859*.
- 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. In *Advances in Neural Information Processing Systems*, volume 36.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc 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.
- Shuo Zhang, Zezhou Huang, and Eugene Wu. 2024. Data cleaning using large language models. *arXiv preprint arXiv:2410.15547*.

# ENTITY-LINKINGS: A Unified Library for Entity Linking

**Yuya Sawda, Tsuyoshi Fujita, Yusuke Sakai, Taro Watanabe**  
Nara Institute of Science and Technology (NAIST)  
{yuya.sawada.sr7, sakai.yusuke.sr9, taro}@is.naist.jp  
fujita.tsuyoshi.fy4@naist.ac.jp

## Abstract

Entity linking (EL) aims to disambiguate named entities in text by mapping them to the appropriate entities in a knowledge base. However, it is difficult to use some EL methods, as they sometimes have issues in reproducibility due to limited maintenance or the lack of official resources. To address this, we introduce ENTITY-LINKINGS, a unified library for using and developing entity linking systems through a unified interface. Our library flexibly integrates various candidate retrievers and re-ranking models, making it easy to compare and use any entity linking methods within a unified framework. In addition, it is designed with a strong emphasis on API usability, making it highly extensible, and it supports both command-line tools and APIs. Our code is available on GitHub<sup>1</sup> and is also distributed via PyPI<sup>2</sup> under the MIT-license. The video is available on YouTube<sup>3</sup>.

## 1 Introduction

Entity linking (EL) is the task of mapping named entities in text to canonical entries in a knowledge base (KB). As shown in Figure 1, since Japan has had many emperors throughout its history, named entities, i.e., *emperor*, are inherently ambiguous. Therefore, it is necessary to identify the correct entity and link it to a canonical entry in the KBs. Most named entities are context-dependent and must be disambiguated by linking them to normalized identifiers in a KB for identification.

A typical EL pipeline consists of two steps: it first detects entity mentions in text (**Mention Detection; MD**), and then links each mention to a unique identifier in the KBs (**Entity Disambiguation; ED**). EL has broad applicability in automating the creation of hyperlinks for domain-specific

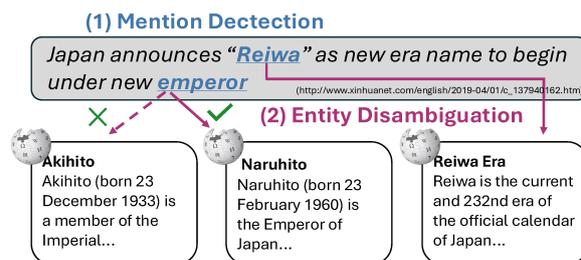


Figure 1: Overview of the entity linking task.

terms in resources such as Wikipedia, FAQs, and product manuals. In addition, it has been widely applied to a broad range of interdisciplinary and practical domains, including clinical records, financial documents, legal statutes, and contracts.

While MD can leverage ready-made named entity recognition (NER) tools such as spaCy (Honni-bal et al., 2020) and Flair (Akbik et al., 2018), ED requires a dedicated implementation that involves constructing and maintaining the task-specific KBs and retrieval components. As a result, each method is implemented individually, which poses challenges for comparison under unified conditions. In fact, even when using a common knowledge base such as a Wikidata dump (Vrandečić and Krötzsch, 2014), many subtle but critical differences exist across methods, including the dump timing, candidate pruning strategies, and data formats, all of which hinder fair and consistent evaluation. Moreover, several official implementations are either unavailable, e.g., Févry et al. (2020); Wang et al. (2024a) or no longer maintained, e.g., Wu et al. (2020), leading to serious reproducibility issues<sup>4</sup>. Although some EL evaluation benchmarks (Milich and Akbik, 2023; Röder et al., 2018) have been proposed, many studies still report baseline perfor-

<sup>4</sup>For example, BLINK’s implementation on Github is reported an open issue that hard negatives sampling is not implemented: [github.com/facebookresearch/BLINK/issues/31](https://github.com/facebookresearch/BLINK/issues/31). Due to this issue, Rucker and Akbik (2025) used a BLINK model trained only on in-batch samples as a baseline.

<sup>1</sup>[github.com/naist-nlp/entity-linkings](https://github.com/naist-nlp/entity-linkings)

<sup>2</sup>`pip install entity-linkings`

<sup>3</sup><https://youtu.be/xFx05wBoz5E>

mance by citing the originally reported scores from previous work. This practice increasingly undermines the scientific reproducibility of EL research. Therefore, it is highly desirable to establish a unified framework for utilizing and evaluating entity linking methods.

We introduce ENTITY-LINKINGS, a unified library that supports multiple modern EL systems and datasets. In particular, we primarily focus on ED, which has become the central component of modern EL methods. It facilitates easy reproduction, simplifies the implementation of custom models, and consolidates experimental setups. We carefully decompose EL systems into three modular components into a unified pipeline: *Mention Detector*, *Candidate Generator*, and *Candidate Reranker*, as illustrated in Figure 2. This design enables flexible combinations, thereby improving extensibility, maintainability, and transparency. ENTITY-LINKINGS supports comprehensive experimentation, provides well-documented interfaces, and supports both CLI- and Python API-based access. We hope that it will further accelerate EL research.

## 2 Background

### 2.1 Preliminaries on Entity Linking

Entity linking (EL) is the task of associating mentions in natural language text with entries in a knowledge base (KB). Formally, we represent the input text  $x$  as a sequence of tokens  $x = x_1, x_2, \dots, x_n$ , where each  $x_i$  denotes the  $i$ -th token in the text after tokenization. A mention span  $m$  is defined as a contiguous subsequence of tokens  $m = x_i \dots x_j$ , where  $1 \leq i \leq j \leq n$ , meaning that the span starts at token position  $i$  and ends at token position  $j$ . EL aims to produce an entry  $e$  in the KB  $\mathcal{E}$  ( $e \in \mathcal{E}$ ) corresponding to each mention  $m$  in  $x$ .

Broadly, existing EL methods decompose the task into two subtasks: *mention detection* (MD) and *entity disambiguation* (ED), as shown in Figure 2. In pipeline systems, mention spans in the input text are first identified in the MD stage and then linked to KB entries in the ED stage. In the MD stage, previous work<sup>5</sup> often employs off-the-shelf NER modules such as spaCy (Honnibal et al., 2020). In the ED stage, the identified mentions are

<sup>5</sup>Although a few approaches, namely end-to-end EL systems, jointly train the MD module with the linking component, in practice, the main difference is whether the MD module is trained or not. Hence, we adopt a simplified description here.

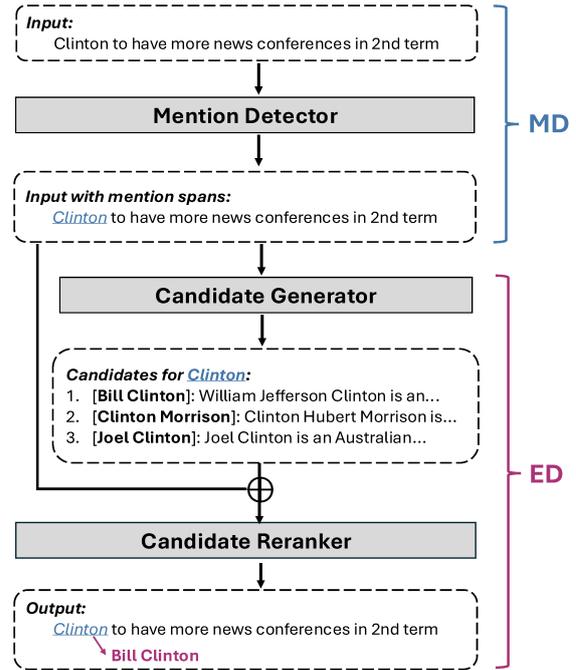


Figure 2: Main components of EL systems and their data processing flow.

linked to their corresponding entities. Specifically, EL systems calculate the probability  $p(e|m)$ , i.e., an entry  $e$  is the target for a given mention  $m$ , and output the entry with the highest probability as the target. Nevertheless, for large-scale KBs where  $|\mathcal{E}|$  is large, exhaustive computation of  $p(e|m)$  for all entries is expensive and complicates the identification of correct entries. To improve the accuracy and efficiency of the system, most ED studies introduce two submodules: *candidate generator* and *candidate reranker*. The candidate generator retrieves a small set  $\hat{\mathcal{E}}$  of candidate entries for each mention from the KB  $\mathcal{E}$ , and the candidate reranker scores these candidates and selects the most plausible one.

For evaluation, the *InKB* micro-F1 score (Röder et al., 2018) is typically used as the standard metric for the EL system. This metric considers a prediction correct if both the extracted mention span and the target entry match the ground truth, and calculates F1 scores as the harmonic mean of precision and recall. On the other hand, since most ED studies focus on specific subtasks of candidate selection or reranking, metrics that assume oracle mention spans are also employed. Specifically,  $\text{Recall}@k$  is used for candidate selection, where a prediction is correct if the gold entry is included in the top- $k$  predicted candidates. For reranking, Top-1 accuracy is used to measure whether the predicted entity from the candidates matches the gold entry.

## 2.2 Main Components of EL Systems

We carefully decompose a wide range of EL systems into three components: *Mention Detector*, *Candidate Generator* and *Candidate Reranker*, thereby enabling the unified description and implementation within a standardized framework.

**Mention Detector** The mention detector is typically used to extract mentions corresponding to specific entries in a KB and is executed before the ED step. Since wrong mention detection leads to error propagation within the pipeline system, recent studies have proposed methods that perform MD after the ED stage (Wang et al., 2024a; Zhang et al., 2022) or specifically enhance the mention detector for the ED step (Tedeschi et al., 2021). Recent studies on end-to-end EL have also introduced unified models that handle both MD and ED through multi-task learning (Ayoola et al., 2022), and that simultaneously execute MD and ED using constrained decoding (Cao et al., 2021). Nevertheless, even in end-to-end EL, the process can still be conceptually decomposed into MD and ED stages, and most systems explicitly maintain a mention detection step in practice. Furthermore, most EL studies focus on enhancing ED, since the MD component can rely on NER tools such as spaCy.

**Candidate Generator** It retrieves candidate entries from the KBs by taking as input the text with identified mention spans. Existing retrievers can be broadly categorized into textual-frequency-based methods and vector-similarity-based methods. Frequency-based retrievers rely on lexical overlap between the input text and KB entries, and hyperlink counts in Wikipedia. Notable examples include **BM25** (Lù, 2024) and entity-mention prior probabilities such as the *Zelda Candidate List* (Milich and Akbik, 2023). Similarity-based retrievers embed mentions and entities into a shared vector space and perform nearest-neighbor search to retrieve candidate entries. While **Dual-Encoder** is widely used as a similarity-based retriever (Wu et al., 2020; Gillick et al., 2019), recent studies also employ **Text Embedding Models** (Wang et al., 2024b). Formally, similarity-based retrievers compute the relevance between a mention-context representation and a candidate entity representation using the inner product of their vector embeddings:

$$\text{similarity}(m, e) = \mathbf{h}_m^\top \mathbf{h}_e, \quad (1)$$

where  $\mathbf{h}_m$  and  $\mathbf{h}_e$  denote the vector representations of the mention  $m$  and the candidate KB entry  $e \in \mathcal{E}$ ,

respectively. Each vector is obtained by encoding the corresponding token sequence:

$$\mathbf{h}_m = \text{red}(E_1(\tau_m)), \quad (2)$$

$$\mathbf{h}_e = \text{red}(E_2(\tau_e)), \quad (3)$$

where  $\tau_m$  and  $\tau_e$  are input representations of mention and entry, respectively.  $E_1$  and  $E_2$  are encoders, with text embedding models typically using a single shared encoder, i.e.,  $E_1 = E_2$ , and dual-encoder employing two separate encoders.  $\text{red}(\cdot)$  is a function that reduces the encoder output into a fixed-size vector by applying average pooling over the final-layer token embeddings in the case of text embedding models or by taking the final-layer [CLS] token representation in dual-encoder models. We follow the experiments of Wu et al. (2020) in the construction of  $\tau_m$  and  $\tau_e$ .

**Candidate Reranker** The reranker is used in ED to refine the candidate set returned by the candidate generator and select the most appropriate entity for each mention. Rerankers using encoder-only models are widely investigated. Encoder-only rerankers take the mention context and candidate entity descriptions as input and compute a relevance score using cross-encoders (Logeswaran et al., 2019), enabling the model to capture interactions between the mention context and the candidates, e.g., **BLINK** (Wu et al., 2020)<sup>6</sup>, **FEVRY** (Février et al., 2020), and **ExtEnD** (Barba et al., 2022). More recently, rerankers based on encoder-decoder or decoder-only models have also been explored. **FusionED** (Wang et al., 2024a) employs an encoder-decoder model to jointly encode the mention context together with all candidate entities and then decode over the fused representations to select the correct entity. **ChatEL** (Ding et al., 2024) prompts decoder-only large language models (LLMs) with the mention context and the candidates to perform reranking in a generative manner. **ReFinED** (Ayoola et al., 2022) formalize the end-to-end EL system, but they use the candidate lists to reduce training cost, and they can execute ED mode by inputting mentions.

## 2.3 Dataset

Since EL systems are often expected to generalize across domains sharing the same KBs, recent studies mostly evaluate a single trained model across

<sup>6</sup>BLINK comprises a dual-encoder and a cross-encoder, where the cross-encoder corresponds to the candidate reranker.

multiple evaluation datasets from diverse sources. For example, **GERBIL** (Röder et al., 2018) is a widely used EL benchmark that uses **AIDA-CoNLL** (Hoffart et al., 2011) dataset for training, and evaluates models on the AIDA-CoNLL test set, as well as eight out-of-domain datasets: **MSNBC** (Cucerzan, 2007), **AQUAINT** (Milne and Witten, 2008), **KORE50** (Hoffart et al., 2012), **N3-Reuters-128**, **N3-RSS-500** (Röder et al., 2014), **Derczynski** (Derczynski et al., 2015), **OKE-2015** (Nuzzolese et al., 2015), and **OKE-2016** (Nuzzolese et al., 2016).

Furthermore, recent mainstream evaluations of EL systems focus on improving individual components, and component-specific evaluation datasets are typically used. For candidate generators, **WikilinksNED Unseen Mentions** (Eshel et al., 2017; Onoe and Durrett, 2020), which comprises diverse ambiguous entities found in web-crawled text, and **ZeshEL** (Logeswaran et al., 2019), which features unique entities associated with specific domains such as fictional books and film series from Wikia<sup>7</sup>, are often used. Both datasets provide a challenging evaluation setting by ensuring that all mention-entity pairs in the test set are unseen during training. For candidate rerankers, **ZELDA** (Milich and Akbik, 2023) is often used as a comprehensive ED benchmark. ZELDA contains a training dataset comprising 95k documents and a fixed entity vocabulary comprising 82.2k entries derived from the Kensho Derived Wikimedia Dataset<sup>8</sup>. ZELDA also provides aggregated mention link counts from Wikipedia, Wikidata, and Wikilinks to construct candidate lists. Furthermore, the evaluation employs a total of nine datasets: **TWEEKI** (Harandizadeh and Singh, 2020), **AIDA-B** (Hoffart et al., 2011), **REDDIT-POSTS**, **REDDIT-COMMENTS** (Botzer et al., 2021), **WNED-WIKI**, **WNED-CWEB** (Guo and Barbosa, 2018), and **SHADOWLINK-{TOP, SHADOW, TAIL}** (Provatorova et al., 2021).

For KBs, prior work often used independently acquired Wikipedia dumps because there was no standardized KB version. In contrast, recent studies often employ the 5.9M Wikipedia pages provided by the **KILT** benchmark as the target KB (Petroni et al., 2021; Zhang et al., 2022; Wang et al., 2024a). Although WikilinksNED Unseen Mentions does not specify the version of Wikipedia, KILT (Petroni

et al., 2021) can be used as the target KB. For ZeshEL, it uses the 49.2k entries derived from Wikia dumps.

### 3 Problems of Existing EL Systems

**Inconsistent Implementations** Although many EL systems have been proposed, their implementations are typically designed with system-specific interfaces and data flows, resulting in limited compatibility and reusability. In particular, as described in Section 2.2, although EL systems can be decomposed into three unified pipeline modules and each component can, in principle, be reused modularly, the lack of standardized implementations makes such reuse difficult in practice. As a result, reproducing prior work or replacing individual components, e.g., using a different candidate generator, becomes challenging, hindering unified evaluation.

**Transparency and Fair Comparison** In particular, previous studies often employ different candidate generators, such as frequency-based or similarity-based retrievers, which makes it challenging to conduct fair comparisons of other modules, e.g., candidate rerankers, across studies. Furthermore, although many evaluation datasets have been proposed and are publicly available, their unified use is hindered by differences in data formats, preprocessing pipelines, and even the snapshots of the knowledge bases. Additionally, methods that depend on such dataset-specific precompiled candidate sets, e.g., PPRforNED (Perschina et al., 2015), cannot be directly applied to other datasets without equivalent precompiled resources, thereby preventing comprehensive evaluation across diverse benchmarks. Moreover, most studies heavily rely on reported scores without re-implementation under consistent settings, which hinders the disentanglement of individual contributions and fair comparisons.

### 4 Our Library: ENTITY-LINKINGS

ENTITY-LINKINGS supports recent EL systems through a unified interface. As described in Section 2.2, we decompose an EL system into three main components: *Mention Detector*, *Candidate Generator* and *Candidate Reranker*, thereby achieving clear modularity, high extensibility, and ease of maintenance. Moreover, we provide multiple evaluation datasets in a consistent format to enhance usability. We also offer basic pretrained models and precompiled KBs, enabling the comparison of

<sup>7</sup><https://www.fandom.com/>

<sup>8</sup><http://datasets.kensho.com/datasets/wikimedia>

custom models and datasets under a unified experimental environment with minimal effort.

#### 4.1 Supported Methods and Datasets

We continuously update the supported systems and datasets. Please refer to our GitHub repository<sup>9</sup> for details on the latest supported methods.

**Methods** Currently, ENTITY-LINKINGS covers the major EL systems introduced in Section 2.2. Specifically, for the ED stage components, ENTITY-LINKINGS supports a wide range of representative systems. For candidate generators, it includes **BM25**, **Dual-Encoder**, and **Text Embedding Models**. For rerankers, it supports **BLINK**, **FEVRY**, **ExtEnD**, **FusionED**, and **ChatEL**.

**Datasets** ENTITY-LINKINGS can download and use the standard EL datasets introduced in Section 2.3, such as those included in **ZELDA** (Milich and Akbik, 2023) and **GERBIL** (Röder et al., 2018), directly through our HuggingFace Collection<sup>10</sup>. In addition, the library accepts arbitrary datasets in a simple unified JSONL format. A complete list of supported datasets, their licenses, and detailed information is provided in Appendix A.

#### 4.2 Interfaces

ENTITY-LINKINGS has two interfaces: Python application programming interface (API) and command-line interface (CLI). Listing 1 illustrates an example workflow via Python API for training, evaluation, and prediction using the ZELDA dataset, a Dual-Encoder candidate generator, and the BLINK reranker. The KB, candidate generator, and reranker are instantiated by specifying their identifiers in the `load_dictionary()`, `get_retrievers()`, and `get_rerankers()` functions. Once these components are loaded, training and evaluation proceed through the unified `train()` and `evaluate()` methods. A trained model can be applied to any input text using the `predict()` method. Replacing the dataset, dictionary, or any EL component requires only changing the corresponding identifier passed to the API, without modifying the rest of the workflow. In addition, predefined datasets and KBs, as well as user-provided datasets and KBs, can be loaded by supplying their file paths to `load_dataset()`<sup>11</sup> and

```
1 from datasets import load_dataset
2 from entity_linkings import get_retrievers,
   get_rerankers, ELPipeline, load_dictionary
3
4 # Existing Corpus
5 dataset = load_dataset('naist-nlp/zelda')
6 dictionary = load_dictionary('zelda')
7 # Custom Corpus
8 # dataset = load_dataset('json', data_files
   = {train: 'train.jsonl', 'validation':
   'valid.jsonl', 'test': 'test.jsonl'})
9 # dictionary = load_dictionary('dict.jsonl')
10
11 # Model loading
12 retriever_cls=get_retrievers('dualencoder')
13 retriever=retriever_cls(dictionary,
   config=retriever_cls.Config())
14 reranker_cls=get_rerankers('crossencoder')
15 reranker=model_cls(retriever,
   config=model_cls.Config())
16
17 # Training
18 result = reranker.train(dataset['train'],
   dataset['validation'])
19
20 # Evaluation
21 metrics = retriever.evaluate(dataset['test'])
22 # Output (metrics): {'R@1': , 'R@10': ,
   'R@50':, 'R@100':, 'MRR':;}
23 metrics = reranker.evaluate(dataset['test'])
24 # Output (metrics): {'Acc': }
25
26 # Prediction
27 sentence = 'Toyota makes cars.'
28 spans=[(0,6)]
29 predictions = reranker.predict(sentence,
   spans)
30 # Output: [{'start': 0, 'end': 6,
   'id':'30984'}]
```

Listing 1: An implementation of training and evaluation using BLINK as a model and ZELDA as a dataset.

`load_dictionary()` in lines 5–9. This design allows ENTITY-LINKINGS to support both standard benchmarks and custom EL datasets uniformly.

**Extensibility** All components are implemented by inheriting from abstract base classes. For instance, retrievers or rerankers are defined via `RetrieverBase` or `RerankerBase` abstract class, respectively. These base classes specify the minimum set of required methods that must be overridden in the derived classes. The components can be instantiated through unified factory functions such as `get_retrievers()` and `get_rerankers()`, which guarantee a consistent interface across different user implementations. In addition, we provide comprehensive test code for all components with CI/CD support, along with well-written documentation. This modular design

<sup>9</sup><https://github.com/naist-nlp/entity-linkings>

<sup>10</sup><https://hf.co/collections/naist-nlp/entity-linkings>

<sup>11</sup><https://hf.co/docs/datasets/en/loading>

| Model          | R@1   | R@10  | R@50  | R@100 |
|----------------|-------|-------|-------|-------|
| BM25           | 0.207 | 0.440 | 0.556 | 0.598 |
| Dual-Encoder   | 0.361 | 0.594 | 0.693 | 0.735 |
| Text Embedding | 0.429 | 0.796 | 0.796 | 0.834 |

Table 1: Candidate generation results in ZeshEL.

enables straightforward extensibility with minimal implementation effort.

**Reproducibility** We also support configuration via YAML input and export all settings to YAML files, enabling users to share the exact experimental configurations and ensuring high reproducibility as well as deterministic results across runs. Furthermore, we provide basic pretrained weights and precompiled KBs, allowing us to easily reproduce results and conduct fair comparisons.

## 5 Experiments

In this paper, we follow a standard EL evaluation protocol for simplicity and to ensure a controlled experiment. Specifically, we fix *Mention Detector* to spaCy and focus our evaluation on the ED stage. We compare the performance of two EL components, i.e., *Candidate Generator* and *Candidate Reranker*. We then construct multiple pipeline systems by combining each component and compare their performance across these full pipeline systems. Appendix B describes the detailed settings.

### 5.1 Candidate Generator

#### 5.1.1 Setup

We evaluate three text retrievers: **BM25**, **Dual-Encoder** and **Text Embedding Model**. For similarity-based retrievers, we use bert-base-uncased (110M×2) and e5-base (110M) to ensure that the same core Transformer architecture is shared. Following Wu et al. (2020), we train these encoders using in-batch random negatives and top-10 hard negatives. We employ R@k as our evaluation metric, where k ranges from 1 to 100. We use the **ZeshEL** benchmark as introduced in Section 2.3.

#### 5.1.2 Results

Table 1 shows the results. As the value of k increases, the R@k scores for all three retrievers consistently increase, and similarity-based retrievers outperform the frequency-based retriever by effectively leveraging contextual information and entity descriptions. Hence, we confirm that the retrievers are functioning as we expected.

## 5.2 Candidate Reranker

### 5.2.1 Setup

We evaluate the eight methods: **Most Frequent Sense (MFS)**, **Dual-Encoder**, **Text Embedding Model**, **FEVRY**, **ExtEnD**, **ChatEL**, and **FusionED** on the ZELDA benchmark. MFS refers to the entity that is most frequently linked from each mention in the Kensho Wikimedia dataset, Wikilinks web corpus (Singh et al., 2012), and Wikidata, which is a general baseline for ZELDA. For each mention, we select the top 30 most frequent entries from the candidate list. During training, if the ground-truth entity is not among the top 30, we replace the 30th entry with it. For candidate lists with fewer than 30 entries, we fill the remaining slots by randomly sampling from the entire KB.

We use bert-base-uncased (110M), longformer-base (149M), and flan-t5-small (77M) to ensure that the total number of parameters is comparable across models.

### 5.2.2 Results

Table 2 shows the accuracy for each evaluation split. Although all methods were trained using a small-scale encoder with only the top 30 candidate entries, they consistently outperform the minimal baseline MFS. Furthermore, the results are comparable to those reported in previous work, which confirms each method is working correctly as intended. For the candidate rerankers, almost all the rerankers underperformed compared to the retriever. While rerankers are restricted to selecting from a fixed candidate list, the retriever is trained using hard negatives. This training regimen likely allows the retriever to leverage contextual cues more effectively. This observation is consistent with the findings reported by Rucker and Akbik (2025).

## 5.3 EL System Evaluation

### 5.3.1 Setup

We use the GERBIL benchmark, which is widely used for evaluating EL systems. We employ a Dual-Encoder model as the candidate generator and FEVRY, Cross-Encoder, ExtEnD, and ChatEL as the reranker, which achieves good performance in our ED evaluation. In addition, we report reproduction results obtained using the official implementation and pretrained weights of ReFinED (Ayoola et al., 2022) as a system baseline. For the evaluation metric, we employ *InKB* micro-F1 score. A predicted mention is regarded as correct only if

|                | backbone                   | AIDA-B | TWEEKI | REDDIT-POSTS | REDDIT-COMM | WNED-CWEB | WNED-WIKI | SLINKS-TAIL | SLINKS-SHADOW | SLINKS-TOP |
|----------------|----------------------------|--------|--------|--------------|-------------|-----------|-----------|-------------|---------------|------------|
| MFS            | –                          | 0.629  | 0.700  | 0.825        | 0.794       | 0.605     | 0.648     | 0.991       | 0.146         | 0.402      |
| Text Embedding | E5 <sub>Base</sub>         | 0.837  | 0.809  | 0.905        | 0.915       | 0.718     | 0.900     | 0.991       | 0.675         | 0.694      |
| Dual-Encoder   | BERT <sub>Base</sub>       | 0.821  | 0.779  | 0.912        | 0.876       | 0.703     | 0.899     | 0.988       | 0.637         | 0.658      |
| FEVRY          | BERT <sub>Base</sub>       | 0.762  | 0.748  | 0.888        | 0.842       | 0.697     | 0.849     | 0.805       | 0.322         | 0.439      |
| Cross-Encoder  | BERT <sub>Base</sub>       | 0.793  | 0.792  | 0.918        | 0.911       | 0.722     | 0.857     | 0.996       | 0.442         | 0.591      |
| ExtEnD         | Longformer <sub>Base</sub> | 0.800  | 0.807  | 0.922        | 0.920       | 0.713     | 0.874     | 0.994       | 0.379         | 0.534      |
| FusionED       | FlanT5 <sub>small</sub>    | 0.638  | 0.657  | 0.801        | 0.760       | 0.614     | 0.761     | 0.989       | 0.351         | 0.468      |
| ChatEL         | GPT-4o <sub>mini</sub>     | 0.756  | 0.758  | 0.851        | 0.814       | 0.686     | 0.716     | 0.980       | 0.380         | 0.730      |

Table 2: Entity Disambiguation results in ZELDA benchmark using our library.

|                 | MSNBC | ACE2004 | Derczynski | KORE50 | R128  | R500  | OKE15 | OKE16 |
|-----------------|-------|---------|------------|--------|-------|-------|-------|-------|
| Dual-Encoder    | 0.437 | 0.133   | 0.304      | 0.379  | 0.325 | 0.256 | 0.379 | 0.349 |
| + FEVRY         | 0.146 | 0.100   | 0.076      | 0.069  | 0.193 | 0.071 | 0.118 | 0.127 |
| + Cross-Encoder | 0.474 | 0.156   | 0.324      | 0.336  | 0.354 | 0.307 | 0.416 | 0.375 |
| + ExtEnD        | 0.472 | 0.156   | 0.336      | 0.400  | 0.345 | 0.300 | 0.378 | 0.351 |
| + ChatEL        | 0.407 | 0.118   | 0.281      | 0.543  | 0.287 | 0.243 | 0.349 | 0.330 |
| ReFinED*        | 0.561 | 0.197   | 0.465      | 0.55.7 | 0.474 | 0.348 | 0.599 | 0.573 |

Table 3: End-to-End Entity Linking results in GERBIL. \* means the model that trains using external resources.

both the mention span and the target entry match the ground truth for entities present in the KB.

### 5.3.2 Results

As shown in Table 3, our straightforwardly constructed system achieves moderate performance compared with existing EL systems, indicating that our library can build competitive EL systems with minimal implementation effort. Among the rerankers, the Cross-Encoder and ExtEnd consistently improved the performance from the results by the Dual-Encoder, whereas FEVRY exhibited a significant decline in performance. The performance gap likely stems from the fact that FEVRY directly projects the representation of the span by concatenating the representation at the span start and end into the entity embedding space without referring to entity titles or descriptions. Unlike competing models, FEVRY fails to leverage explicit semantic information, such as entity titles or descriptions, which limits its ability to resolve entities in out-of-domain datasets.

## 6 Conclusion

We proposed ENTITY-LINKINGS, a unified library for EL systems based on a standardized three-component interface. In our experiments, we conducted replication studies using standardized EL benchmarks, such as ZELDA and GERBIL. We confirmed that the performance of our three candidate retrievers and five candidate rerankers is consistent with results reported in previous works.

Moving forward, to consider the high rate of errors in mention detection by spaCy, we plan to explore the integration of End-to-End EL architectures (Cao et al., 2021; Shavarani and Sarkar, 2023; Ayoola et al., 2022) and Retriever-to-Reader models (Zhang et al., 2022). We will continue to actively maintain and expand it, and we hope that it will further advance both EL research and the community.

## Ethics and Broader Impact Statement

Using ENTITY-LINKINGS enhances the reproducibility and transparency of experiments, which is essential from a research ethics perspective. The ACL Rolling Review checklist<sup>12</sup> explicitly emphasizes implementation and experimental settings, underscoring their importance to the community. Through continued maintenance and expansion of ENTITY-LINKINGS, we aim to further support these efforts.

While long-term community adoption ultimately depends on broader usage, we emphasize that we plan to actively use ENTITY-LINKINGS in our own future research and will therefore continue to maintain and update it on a regular basis. This library stems from the practical challenges we encountered while building and evaluating EL systems, and our goal is to share this solution with the broader community. Furthermore, this demonstra-

<sup>12</sup><https://aclrollingreview.org/responsibleNLPresearch/>

tion paper presents only a subset of essential results and brief analyses for the purpose of system verification. More extensive experimental results are available in our GitHub repository. Nevertheless, the results reported in this paper are sufficient to validate our claims, and the paper can be read as a self-contained and complete study. By making our implementation details and validation results publicly available as much as possible, we aim to provide a broader impact to the research community. This transparency helps avoid redundant experimentation and allows researchers to focus their efforts more effectively. Accordingly, we plan to continuously add and update new results and perform regular maintenance. In this way, ENTITY-LINKINGS is intended to deliver sustained value beyond this paper alone.

We verified the licenses of all datasets used in this study, as summarized in Table 4 in Appendix A, and confirmed that their use complies with all applicable terms. In addition, this work does not involve the generation of harmful content. Accordingly, we ensure that our study is fully compliant with ethical guidelines such as the ACL Ethics Policy<sup>13</sup>.

## Acknowledgements

We thank the anonymous reviewers and the area chair for their valuable comments and suggestions.

The architecture design of ENTITY-LINKINGS is inspired by MBRS (Deguchi et al., 2024) and GEC-METRICS (Goto et al., 2025).

## References

- Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. [Contextual string embeddings for sequence labeling](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Tom Ayoola, Shubhi Tyagi, Joseph Fisher, Christos Christodoulopoulos, and Andrea Pierleoni. 2022. [Re-FinED: An efficient zero-shot-capable approach to end-to-end entity linking](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, pages 209–220, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.
- Edoardo Barba, Luigi Procopio, and Roberto Navigli. 2022. [ExtEnD: Extractive entity disambiguation](#). In <sup>13</sup>[https://www.aclweb.org/adminwiki/index.php/ACL\\_Policy\\_on\\_Publication\\_Ethics](https://www.aclweb.org/adminwiki/index.php/ACL_Policy_on_Publication_Ethics)
- Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2478–2488, Dublin, Ireland. Association for Computational Linguistics.
- Nicholas Botzer, Yifan Ding, and Tim Weninger. 2021. [Reddit entity linking dataset](#). *Inf. Process. Manage.*, 58(3).
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021. [Autoregressive entity retrieval](#). In *International Conference on Learning Representations*.
- Silviu Cucerzan. 2007. [Large-scale named entity disambiguation based on Wikipedia data](#). In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 708–716, Prague, Czech Republic. Association for Computational Linguistics.
- Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, and Taro Watanabe. 2024. [mbrs: A library for minimum Bayes risk decoding](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 351–362, Miami, Florida, USA. Association for Computational Linguistics.
- Leon Derczynski, Diana Maynard, Giuseppe Rizzo, Marieke van Erp, Genevieve Gorrell, Raphaël Troncy, Johann Petrak, and Kalina Bontcheva. 2015. [Analysis of named entity recognition and linking for tweets](#). *Information Processing & Management*, 51(2):32–49.
- Yifan Ding, Qingkai Zeng, and Tim Weninger. 2024. [ChatEL: Entity linking with chatbots](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3086–3097, Torino, Italia. ELRA and ICCL.
- Yotam Eshel, Noam Cohen, Kira Radinsky, Shaul Markovitch, Ikuya Yamada, and Omer Levy. 2017. [Named entity disambiguation for noisy text](#). In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 58–68, Vancouver, Canada. Association for Computational Linguistics.
- Thibault Févry, Nicholas FitzGerald, Livio Baldini Soares, and Tom Kwiatkowski. 2020. [Empirical evaluation of pretraining strategies for supervised entity linking](#). *Preprint*, arXiv:2005.14253.
- Pierre-Yves Genest, Pierre-Edouard Portier, Elöd Egyed-Zsigmond, and Martino Lovisetto. 2023. [Linked-DocRED – Enhancing DocRED with Entity-Linking to Evaluate End-To-End Document-Level Information Extraction Pipelines](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR’23)*, page 11, Taipei, Taiwan. Association for Computing Machinery.

- Daniel Gillick, Sayali Kulkarni, Larry Lansing, Alessandro Presta, Jason Baldridge, Eugene Ie, and Diego Garcia-Olano. 2019. [Learning dense representations for entity retrieval](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 528–537, Hong Kong, China. Association for Computational Linguistics.
- Takumi Goto, Yusuke Sakai, and Taro Watanabe. 2025. [gec-metrics: A unified library for grammatical error correction evaluation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 524–534, Vienna, Austria. Association for Computational Linguistics.
- Zhaochen Guo and Denilson Barbosa. 2014. [Robust entity linking via random walks](#). In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, CIKM 2014, Shanghai, China, November 3-7, 2014*, pages 499–508.
- Zhaochen Guo and Denilson Barbosa. 2018. [Robust named entity disambiguation with random walks](#). *Semant. Web*, 9(4):459–479.
- Bahareh Harandizadeh and Sameer Singh. 2020. [Tweeki: Linking named entities on Twitter to a knowledge graph](#). In *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*, pages 222–231, Online. Association for Computational Linguistics.
- Johannes Hoffart, Stephan Seufert, Dat Ba Nguyen, Martin Theobald, and Gerhard Weikum. 2012. [Kore: keyphrase overlap relatedness for entity disambiguation](#). In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management, CIKM '12*, page 545–554, New York, NY, USA. Association for Computing Machinery.
- Johannes Hoffart, Mohamed Amir Yosef, Iliaria Bordino, Hagen Fürstenau, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. [Robust disambiguation of named entities in text](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 782–792, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. [spaCy: Industrial-strength Natural Language Processing in Python](#).
- Lajanugen Logeswaran, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, Jacob Devlin, and Honglak Lee. 2019. [Zero-shot entity linking by reading entity descriptions](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3449–3460, Florence, Italy. Association for Computational Linguistics.
- Xing Han Lù. 2024. [Bm25s: Orders of magnitude faster lexical search via eager sparse scoring](#). *Preprint*, arXiv:2407.03618.
- Marcel Milich and Alan Akbik. 2023. [ZELDA: A comprehensive benchmark for supervised entity disambiguation](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2061–2072, Dubrovnik, Croatia. Association for Computational Linguistics.
- David Milne and Ian H. Witten. 2008. [Learning to link with wikipedia](#). In *Proceedings of the 17th ACM Conference on Information and Knowledge Management, CIKM '08*, page 509–518, New York, NY, USA. Association for Computing Machinery.
- Andrea Giovanni Nuzzolese, Anna Lisa Gentile, Valentina Presutti, Aldo Gangemi, Darío Garigliotti, and Roberto Navigli. 2015. [Open knowledge extraction challenge](#). In *Semantic Web Evaluation Challenges*, pages 3–15. Springer.
- Andrea Giovanni Nuzzolese, Anna Lisa Gentile, Valentina Presutti, Aldo Gangemi, Robert Meusel, and Heiko Paulheim. 2016. [The second open knowledge extraction challenge](#). In *Semantic Web Challenges*, pages 3–16, Cham. Springer International Publishing.
- Yasumasa Onoe and Greg Durrett. 2020. [Fine-grained entity typing for domain independent entity linking](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8576–8583.
- Maria Pershina, Yifan He, and Ralph Grishman. 2015. [Personalized page rank for named entity disambiguation](#). In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 238–243, Denver, Colorado. Association for Computational Linguistics.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. [KILT: a benchmark for knowledge intensive language tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544, Online. Association for Computational Linguistics.
- Vera Provatorova, Samarth Bhargav, Svitlana Vakulenko, and Evangelos Kanoulas. 2021. [Robustness evaluation of entity disambiguation using prior probes: the case of entity overshadowing](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10501–10510, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Lev Ratinov, Dan Roth, Doug Downey, and Mike Anderson. 2011. [Local and global algorithms for disambiguation to Wikipedia](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1375–1384, Portland, Oregon, USA. Association for Computational Linguistics.

Michael Röder, Ricardo Usbeck, Sebastian Hellmann, Daniel Gerber, and Andreas Both. 2014. **N<sup>3</sup> - a collection of datasets for named entity recognition and disambiguation in the NLP interchange format**. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3529–3533, Reykjavik, Iceland. European Language Resources Association (ELRA).

Susanna Rücker and Alan Akbik. 2025. **Evaluating design decisions for dual encoder-based entity disambiguation**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15685–15701, Vienna, Austria. Association for Computational Linguistics.

Michael Röder, Ricardo Usbeck, and Axel-Cyrille Ngonga Ngomo. 2018. **GERBIL - Benchmarking Named Entity Recognition and Linking consistently**. *Semantic Web*, 9(5):605–625.

Hassan Shavarani and Anoop Sarkar. 2023. **SpEL: Structured prediction for entity linking**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11123–11137, Singapore. Association for Computational Linguistics.

Sameer Singh, Amarnag Subramanya, Fernando Pereira, and Andrew McCallum. 2012. **Wikilinks: A large-scale cross-document coreference corpus labeled via links to Wikipedia**. Technical report, University of Massachusetts, Amherst, Tech. Rep. UM-CS-2012.

Simone Tedeschi, Simone Conia, Francesco Cecconi, and Roberto Navigli. 2021. **Named Entity Recognition for Entity Linking: What works and what's next**. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2584–2596, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Denny Vrandečić and Markus Krötzsch. 2014. **Wikidata: a free collaborative knowledgebase**. *Commun. ACM*, 57(10):78–85.

Junxiong Wang, Ali Mousavi, Omar Attia, Ronak Pradeep, Saloni Potdar, Alexander Rush, Umar Farooq Minhas, and Yunyao Li. 2024a. **Entity disambiguation via fusion entity decoding**. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6524–6536, Mexico City, Mexico. Association for Computational Linguistics.

Liang Wang, Nan Yang, Xiaolong Huang, Binx-ing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2024b. **Text embeddings by weakly-supervised contrastive pre-training**. *Preprint*, arXiv:2212.03533.

Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. **Scalable zero-shot entity linking with dense entity retrieval**. In

*Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6397–6407, Online. Association for Computational Linguistics.

Wenzheng Zhang, Wenyue Hua, and Karl Stratos. 2022. **EntQA: Entity linking as question answering**. In *International Conference on Learning Representations*.

## A Dataset Details

```
1 {
2   "id": "doc-001-P1",
3   "text": "Toyota makes cars.",
4   "entities": [
5     {
6       "start": 0,
7       "end": 6,
8       "label": ["000011"],
9     }
10  ]
11 }
```

Listing 2: Dataset format

```
1 {
2   "id": "000011",
3   "name": "Toyota",
4   "description": "Toyota is a Japanese car
5   manufacturer."
6 }
```

Listing 3: Ontology format

Table 4 shows a complete list of datasets supported by ENTITY-LINKINGS. Additionally, arbitrary datasets and ontologies can be used with the ENTITY-LINKINGS library if they are provided in a JSONL format shown in Listing 2 and 3. Note that training and evaluation splits of AIDA-CoNLL dataset (Hoffart et al., 2011) are not publicly available, and users need to download them after obtaining approval for using the Reuters Corpora<sup>14</sup>. Once downloaded, users can use a preprocessing script we provide to convert AIDA-CoNLL dataset to the format supported in ENTITY-LINKINGS.

## B Details of Experimental Setup

Table 5,6 list the hyperparameter settings used in our experiments. All models are trained and evaluated based on these settings. ENTITY-LINKINGS allows these parameters to be saved and loaded in YAML format, facilitating easy reproducibility.

<sup>14</sup>[trec.nist.gov/data/reuters/reuters.html](http://trec.nist.gov/data/reuters/reuters.html)

| data_id           | Dataset                                               | Domain    | Lang.   | Ontology  | Train | Licence            |
|-------------------|-------------------------------------------------------|-----------|---------|-----------|-------|--------------------|
| msnbc             | MSNBC (Cucerzan, 2007)                                | News      | English | Wikipedia |       | Unknown*           |
| aquaint           | AQUAINT (Milne and Witten, 2008)                      | News      | English | Wikipedia |       | Unknown*           |
| ace2004           | ACE2004 (Ratinov et al., 2011)                        | News      | English | Wikipedia |       | Unknown*           |
| kore50            | KORE50 (Hoffart et al., 2012)                         | News      | English | Wikipedia |       | CC-BY-SA 3.0       |
| n3-r128           | N3-Reuters-128 (Röder et al., 2014)                   | News      | English | Wikipedia |       | GNU AGPL-3.0       |
| n3-r500           | N3-RSS-500 (Röder et al., 2014)                       | RSS       | English | Wikipedia |       | GNU AGPL-3.0       |
| derczynski        | Derczynski (Derczynski et al., 2015)                  | Twitter   | English | Wikipedia |       | CC-BY 4.0          |
| oke-2015          | OKE-2015 (Nuzzolese et al., 2015)                     | News      | English | Wikipedia | Yes   | Unknown*           |
| oke-2016          | OKE-2016 (Nuzzolese et al., 2016)                     | News      | English | Wikipedia | Yes   | Unknown*           |
| wned-wiki         | WNED-WIKI (Guo and Barbosa, 2014)                     | Wikipedia | English | Wikipedia |       | Unknown            |
| wned-cweb         | WNED-CWEB (Guo and Barbosa, 2014)                     | Web       | English | Wikipedia |       | Apache License 2.0 |
| unseen            | WikilinksNED Unseen-Mentions (Onoe and Durrett, 2020) | Web       | English | Wikipedia | Yes   | CC-BY 3.0*         |
| tweeki            | Tweeki EL (Harandizadeh and Singh, 2020)              | Twitter   | English | Wikipedia | Yes   | Apache License 2.0 |
| reddit-comments   | Reddit EL (Botzer et al., 2021)                       | Reddit    | English | Wikipedia |       | CC-BY 4.0          |
| reddit-posts      | Reddit EL (Botzer et al., 2021)                       | Reddit    | English | Wikipedia |       | CC-BY 4.0          |
| shadowlink-shadow | ShadowLink (Provatorova et al., 2021)                 | Wikipedia | English | Wikipedia |       | Unknown*           |
| shadowlink-top    | ShadowLink (Provatorova et al., 2021)                 | Wikipedia | English | Wikipedia |       | Unknown*           |
| shadowlink-tail   | ShadowLink (Provatorova et al., 2021)                 | Wikipedia | English | Wikipedia |       | Unknown*           |
| zeshel            | Zeshel (Logeswaran et al., 2019)                      | Wikia     | English | Wikia     | Yes   | CC-BY-SA           |
| docred            | Linked-DocRED (Genest et al., 2023)                   | News      | English | Wikipedia | Yes   | CC-BY 4.0          |

Table 4: Public Entity Linking Datasets. \*Unknown licence information as provided in the original sources.

|                              | Parameters  |
|------------------------------|-------------|
| Seeds                        | 42          |
| Training Epochs              | 2           |
| Hard negatives               | 10          |
| Batch size (train)           | 32          |
| Batch size (eval)            | 256         |
| Max token length (context)   | 128         |
| Max token length (candidate) | 50          |
| Context window               | 500         |
| Learning rate                | 1e-5        |
| Gradient accumulation steps  | 4           |
| Scheduler                    | linear      |
| Optimizer                    | AdamW       |
| Warmup                       | 0.06        |
| Weight decay                 | 0.01        |
| Max grad norm                | 0.0         |
| Adam beta                    | [0.9, 0.98] |
| Adam epsilon                 | 1e-6        |

Table 5: Hyperparameters for Candidate Generator

## C Details of Command-Line Interface

Listing 4 illustrates the examples via CLI for training with ZELDA dataset and Dual-Encoder candidate generator. The training and evaluation processes for both the candidate retriever and reranker are instantiated through dedicated commands: train-retriever, eval-retriever, train-reranker, and eval-reranker. Additionally, eval-pipeline enables the end-to-end evaluation of the pipeline system, incorporating spaCy for mention detection. To facilitate rigorous model comparison, all training and evaluation workflows are integrated with Weights & Biases.

|                              | ZELDA       | AIDA-CoNLL  |
|------------------------------|-------------|-------------|
| Seeds                        | 42          | 42          |
| Training Epochs              | 1 (10)      | 5 (30)      |
| Candidates                   | 30          | 30          |
| Batch size (train)           | 8 (32)      | 8 (32)      |
| Batch size (eval)            | 32 (256)    | 32 (256)    |
| Max token length (context)   | 128         | 128         |
| Max token length (candidate) | 50          | 50          |
| Learning rate                | 2e-5 (5e-5) | 2e-5 (5e-5) |
| Gradient accumulation steps  | 4           | 4           |
| Scheduler                    | linear      | linear      |
| Optimizer                    | AdamW       | AdamW       |
| Warmup                       | 0.06        | 0.06        |
| Weight decay                 | 0.01        | 0.01        |
| Max grad norm                | 0.0         | 0.0         |
| Adam beta                    | [0.9, 0.98] | [0.9, 0.98] |
| Adam epsilon                 | 1e-6        | 1e-6        |

Table 6: Hyperparameters for Candidate Reranker. The values in parentheses represent the parameters for FEVRY.

```

1 # Model Training
2 entitylinkings-train-retrieval \
3 --retriever_id dualencoder \
4 --dataset_id zelda \
5 --dictionary_id_or_path zelda \
6 # Custom Corpus
7 # --train_file train.jsonl \
8 # --validation_file validation.jsonl \
9 # --dictionary_id_or_path
   dictionary.jsonl \
10 --output_dir save_model/ \
11 --num_hard_negatives 10 \
12 --num_train_epochs 2 \
13 --train_batch_size 8 \
14 --validation_batch_size 32 \
15 --config configs/dualencoder.yaml \
16 --gpu 0,1 \
17 --wandb

```

Listing 4: An CLI example of training, using Dual-Encoder as candidate generator and ZELDA as a dataset.

# ESG-KG: A Multi-modal Knowledge Graph System for Automated Compliance Assessment

Li-Yang Chang<sup>1</sup> Chih-Ming Chen<sup>1</sup> Hen-Hsen Huang<sup>2</sup>  
An-Zi Yen<sup>3</sup> Ming-Feng Tsai<sup>4,5</sup> Chuan-Ju Wang<sup>1</sup>

<sup>1</sup>Research Center for Information Technology Innovation, Academia Sinica, Taiwan

<sup>2</sup>Institute of Information Science, Academia Sinica, Taiwan

<sup>3</sup>Department of Computer Science, National Yang Ming  
Chiao Tung University, Taiwan

<sup>4</sup>Department of Computer Science, National Chengchi University, Taiwan

<sup>5</sup>Delta Electronics, Taiwan

## Abstract

We present ESG-KG, a system that automates ESG compliance assessment through multi-modal information extraction and knowledge graph construction. ESG-KG processes corporate sustainability reports containing diverse data formats—text, tables, figures, and infographics—and extracts ESG-related entities, relationships, and metrics into a structured knowledge graph. This KG-based architecture enables precise cross-modal information retrieval and provides verifiable evidence grounding for downstream analysis. Built upon this foundation, ESG-KG integrates retrieval-augmented generation (RAG) with LLM-based reasoning to automatically evaluate compliance against ESG frameworks and standards. Our demonstration showcases the system’s end-to-end pipeline, from multi-modal document processing to automated compliance scoring, highlighting its capability to handle real-world sustainability reports and generate interpretable assessment results with traceable evidence chains. To facilitate further research, we release our open-source Python toolkit for Automated Compliance Assessment at <https://github.com/cnclabs/website.kg.esg.demo.git>, and a live demonstration video is available at <https://youtu.be/Lj4Zp74J1nY>.

## 1 Introduction

The growing emphasis on Environmental, Social, and Governance (ESG) issues has created a complex landscape of sustainability reporting standards. Among these, the Global Reporting Initiative (GRI) Standards have emerged as the most widely adopted framework for ESG disclosure, providing structured reporting principles and standardized indicators for systematic assessment of organizational sustainability performance. However,

manually assessing compliance remains challenging due to the volume and heterogeneity of modern sustainability reports, which integrate textual, tabular, and visual elements.

Existing automated approaches for ESG compliance assessment typically rely on fact-based retrieval methodologies that decompose statements into individual claims and perform document-level evidence retrieval (Min et al., 2023). Recent advances in knowledge graph-based fact-checking (Chen et al., 2025) and multi-modal analysis (Wang et al., 2024b) have shown promise in handling complex document structures. However, these methods often struggle with the multi-modal nature of sustainability reports, where critical information is embedded not only in text but also in charts, graphs, and infographics. This limitation results in incomplete evidence gathering and compromises compliance verification accuracy.

To overcome these challenges, we present ESG-KG, a system that automates ESG compliance assessment through multi-modal information extraction and knowledge graph construction. Our approach extends beyond fact-level retrieval by incorporating layout analysis and visual data extraction to construct a comprehensive evidence base, capturing quantitative and contextual information from charts, tables, and infographics that text-only systems typically miss.

The extracted multi-modal data is structured into a knowledge graph that semantically aligns ESG disclosures with standard requirements, particularly the GRI Standards. This KG serves as a retrieval-augmented knowledge base, enabling accurate evidence extraction and providing a trusted foundation for LLM-based compliance assessment with traceable reasoning chains.

ESG-KG enables scalable, automated compli-

ance evaluation while substantially reducing manual verification effort. Our demonstration showcases how the system processes real-world sustainability reports end-to-end, bridging multi-modal content with structured ESG criteria to deliver interpretable compliance assessments. By integrating multi-modal extraction with knowledge graph reasoning, ESG-KG promotes greater transparency, reliability, and efficiency in corporate sustainability reporting and compliance verification.

## 2 Related Work

**LLM-based Compliance Assessment.** Automated compliance assessment has evolved from rule-based systems to approaches leveraging LLMs (Radford and Narasimhan, 2018), (Radford et al., 2019), (Brown et al., 2020), (Lewis et al., 2020), (Raffel et al., 2020). Early LLM applications leveraged semantic capabilities to interpret regulatory texts and summarize requirements directly from their parametric knowledge (Min et al., 2023). However, pure LLM approaches prove insufficient for high-stakes auditing due to their susceptibility to hallucination and inability to systematically cross-reference lengthy corporate disclosures against complex regulatory frameworks without external grounding.

**Retrieval-augmented Generation for Fact Verification.** To address these limitations, retrieval-augmented generation (RAG) approaches have emerged that ground LLM reasoning in retrieved evidence. Methods like FActScore (Min et al., 2023) decompose claims into atomic facts and retrieve supporting evidence, while knowledge graph-based fact-checking systems (Chen et al., 2025) provide structured reasoning paths. However, these approaches primarily focus on textual evidence and struggle with the multi-modal nature of corporate reports.

**Multi-modal Document Understanding.** Recent work has begun addressing multi-modal information extraction from complex documents. However, a critical gap remains: corporate sustainability reports often present crucial quantitative metrics—such as emissions and resource use—not in continuous text but embedded in charts, infographics, and complex tables (Gupta et al., 2025). Existing ESG benchmarks and text-centric systems explicitly acknowledge the “exclusion of visual elements” as a major limitation that directly impacts evidence coverage (He et al., 2025).

Systems like SubstationAI (Wang et al., 2024b) demonstrate the value of processing visual elements alongside text, while recent document understanding approaches (Zhang et al., 2024) show progress in table structure recognition and chart-to-structured-data conversion. However, these multi-modal capabilities have not been systematically integrated with knowledge graph architectures for compliance assessment.

ESG-KG addresses this gap by combining specialized multi-modal extraction—including table structure recognition and chart data conversion—with knowledge graph construction, creating a unified framework where visual and textual evidence are jointly represented and retrievable for compliance verification.

## 3 The Proposed ESG-KG System

We present ESG-KG, a system for automated ESG compliance assessment that bridges unstructured regulatory documents with verifiable evidence retrieval from multi-modal corporate reports.

The system operates through three core components: (1) construction of a regulatory knowledge graph (KG) from GRI standards, (2) semantic refinement to ensure entity uniqueness and relational consistency, and (3) an online scoring pipeline that processes multimodal reports for evidence retrieval and compliance evaluation.

Figure 1 illustrates the system architecture.

### 3.1 Regulatory Knowledge Graph Construction

To transform the dense, unstructured text of GRI standards into a machine-interpretable format, we employ an LLM-driven extraction pipeline. We formalize the GRI standards as a Knowledge Graph (KG)  $\mathcal{G} = (\mathcal{E}, \mathcal{R})$ , where  $\mathcal{E}$  represents the set of entities and  $\mathcal{R}$  represents semantic relations between them.

The LLM parses raw standard documents into structured subgraphs following a strict schema. For each standard  $\mathcal{S}$ , extracted entities are decomposed into three categories:

$$\mathcal{E}_{\mathcal{S}} = \mathcal{T}_{\text{disc}} \cup \mathcal{R}_{\text{req}} \cup \mathcal{C}_{\text{score}}, \quad (1)$$

where:

- $\mathcal{T}_{\text{disc}}$  (**Disclosure Targets**): Specific metrics organizations must disclose (e.g., “Scope 1 GHG emissions”).

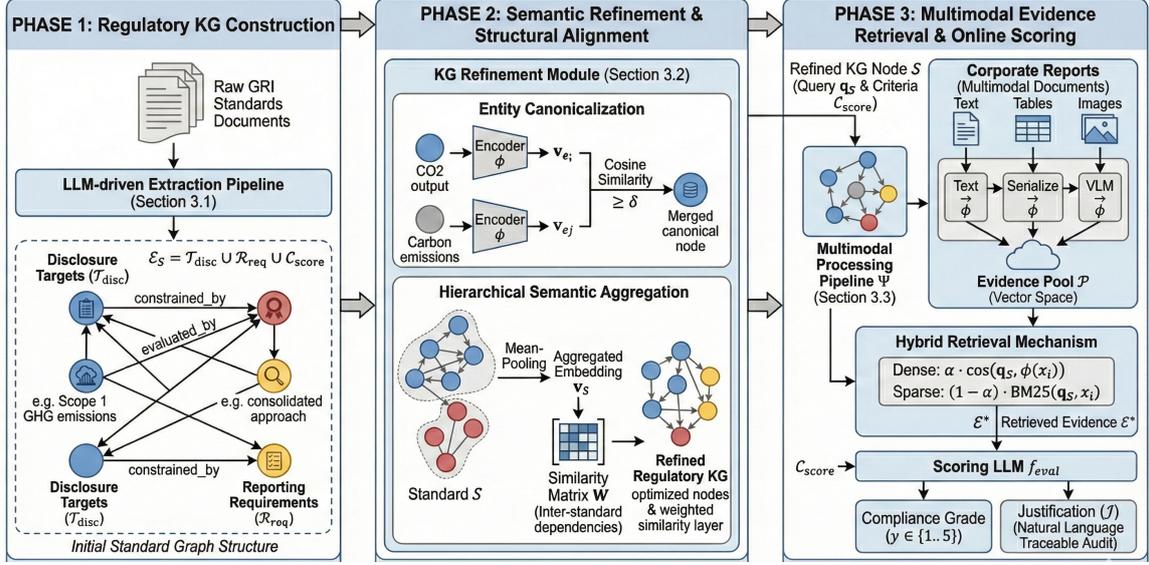


Figure 1: Multimodal Evidence Retrieval and Online Scoring Pipeline

- $\mathcal{R}_{req}$  (**Reporting Requirements**): Contextual constraints defining how disclosures should be prepared (e.g., “consolidated approach”).
- $\mathcal{C}_{score}$  (**Scoring Criteria**): Evaluative rubrics for assessing disclosure completeness and quality .

Relations  $r \in \mathcal{R}$  preserve the logical structure by linking these components. Requirement constrain targets via triples  $(e_t, constrained\_by, e_r)$  where  $e_t \in \mathcal{T}_{disc}$  and  $e_r \in \mathcal{R}_{req}$ . Evaluation logic is captured as  $(e_t, evaluated\_by, e_c)$  where  $e_c \in \mathcal{C}_{score}$ . This structured representation creates a “standard graph” that serves as the foundation for subsequent compliance verification.

### 3.2 Semantic Refinement and Structural Alignment

Raw triples extracted from natural language often contain redundancy and ambiguity. To construct a robust reasoning foundation, we implement a KG refinement module that addresses two key challenges: entity canonicalization and hierarchical semantic aggregation.

**Entity Canonicalization.** We employ a semantic encoder  $\phi$  (e.g., SBERT) to map the textual description of an entity to a dense vector representation  $\mathbf{v}_e = \phi(e)$ . To address synonymous concepts (e.g., “CO2 output” vs. “Carbon emissions”), we perform pairwise comparisons. Two entities  $e_i$  and  $e_j$  are merged into a single canonical node if their cosine

similarity exceeds the threshold  $\delta$ :

$$\text{merge}(e_i, e_j) \iff \frac{\mathbf{v}_{e_i} \cdot \mathbf{v}_{e_j}}{\|\mathbf{v}_{e_i}\| \|\mathbf{v}_{e_j}\|} \geq \delta.$$

**Hierarchical Semantic Aggregation.** Following node-level refinement, we compute holistic representations for each GRI Standard to capture inter-standard dependencies. For standard  $S$  with constituent nodes  $\mathcal{E}_S$ , we calculate the aggregated embedding  $\mathbf{v}_S$  via mean-pooling:

$$\mathbf{v}_S = \frac{1}{|\mathcal{E}_S|} \sum_{e \in \mathcal{E}_S} \mathbf{v}_e$$

This aggregation computes the semantic centroid of the standard’s requirements. We then construct a weighted similarity matrix  $\mathbf{W} \in \mathcal{R}^{K \times K}$  (where  $K$  is the total number of standards), with each entry  $W_{ij} = \cos(\mathbf{v}_{S_i}, \mathbf{v}_{S_j})$  represents semantic affinity between standards. This matrix serves as a retrieval prior, enabling the system to identify structurally related requirements beyond exact keyword matching, thereby enhancing evidence recall robustness..

### 3.3 Multimodal Evidence Retrieval and Online Scoring

The final component processes corporate reports and executes compliance audits, as illustrated in Figure 1(3). To address the “multimodal gap” where crucial evidence resides in non-textual formats, we implement a specialized preprocessing pipeline that projects visual and tabular data into a unified semantic space for retrieval.

**Multi-modal Document Processing.** An uploaded document  $\mathcal{D}$  is parsed into segments  $\mathcal{D} = \{u_1, u_2, \dots, u_N\}$ , where each segment  $u_i$  belongs to one of three modalities: text ( $\mathcal{U}_{\text{text}}$ ), tables ( $\mathcal{U}_{\text{tab}}$ ), or images ( $\mathcal{U}_{\text{img}}$ ). We apply modality-specific transformation  $\Psi$  to convert all segments into textual representations:

$$x_i = \Psi(u_i) = \begin{cases} u_i & \text{if } u_i \in \mathcal{U}_{\text{text}} \\ \text{Serialize}(u_i) & \text{if } u_i \in \mathcal{U}_{\text{tab}} \\ \text{VLM}(u_i) & \text{if } u_i \in \mathcal{U}_{\text{img}}, \end{cases}$$

where  $\text{VLM}(\cdot)$  denotes a vision-language model to generate descriptive captions for charts and infographics (preserving quantitative values), and  $\text{Serialize}(\cdot)$  linearizes table structures. Processed segments are then embedded using the semantic encoder  $\phi$  (defined in Sec. 3.2) to form an evidence pool  $\mathcal{P} = \{(\phi(x_i), x_i)\}_{i=1}^N$ .

**Hybrid Retrieval Mechanism.** For a given GRI standard node  $\mathcal{S}$ , the system formulates a query vector  $\mathbf{q}_{\mathcal{S}}$  and employs hybrid retrieval to identify relevant evidence  $\mathcal{E}^* \subset \mathcal{P}$ . The relevance score for candidate segment  $x_i$  combines semantic and lexical matching:

$$\text{Score}(\mathcal{S}, x_i) = \alpha \cdot \cos(\mathbf{q}_{\mathcal{S}}, \phi(x_i)) + (1 - \alpha) \cdot \text{BM25}(\mathbf{q}_{\mathcal{S}}, x_i)$$

where  $\alpha$  balances dense semantic similarity with sparse keyword matching via BM25.

**LLM-based Compliance Scoring.** Retrieved evidence  $\mathcal{E}^*$  and scoring criteria  $\mathcal{C}_{\text{score}}$  (from Eq. (1)) are provided to the scoring LLM, denoted as  $f_{\text{eval}}$ , to generate the compliance assessments:

$$(y, \mathcal{J}) = f_{\text{eval}}(\mathcal{E}^*, \mathcal{C}_{\text{score}}),$$

where  $y \in \{1, 2, 3, 4, 5\}$  represents the compliance grade and  $\mathcal{J}$  contains natural language justification. This design ensures a traceable audit trail, explicitly linking evidence to regulatory requirements through the knowledge graph.

## 4 Implementation and Demonstration

### 4.1 Python Toolkit

To facilitate reproducibility and foster community engagement, we have released the core components of our system as an open-source Python toolkit, available at our GitHub<sup>1</sup>. This toolkit is designed with a modular architecture, consisting of:

<sup>1</sup><https://github.com/cnclabs/website.kg.esg.demo.git>

**Data Ingestion Module.** Leveraging MinerU<sup>2</sup> ((Wang et al., 2024a), (Niu et al., 2025), (He et al., 2024)) for high-fidelity PDF parsing, this module performs advanced layout analysis to accurately extract and serialize multi-modal content—including complex tables, diagrams, and cross-page text flows—from unstructured regulatory documents and corporate reports.

**KG Construction Engine.** The engine implements the logic for building and refining the regulatory knowledge graph, serving as the foundational layer for structuring unstructured compliance standards. It employs semantic extraction algorithms to transform raw regulatory text (e.g., GRI Standards) into a structured graph format, identifying core entities—such as disclosure requirements and metric definitions—and establishing hierarchical relationships between them. Furthermore, the engine includes a refinement layer that resolves entity ambiguity and enforces schema consistency, ensuring a reliable knowledge base for downstream reasoning.

**Evaluation Pipeline.** This pipeline encapsulates the comprehensive logic for evidence retrieval and LLM-based compliance scoring. Crucially, it utilizes the explicit mapping between specific GRI standards and page numbers—voluntarily disclosed by companies in their sustainability reports (typically within the *GRI Content Index*)—as the primary reference for evidence localization. By anchoring the retrieval process to these self-disclosed page references, the system employs a Retrieval-Augmented Generation (RAG) workflow to precisely align corporate data with regulatory nodes. Beyond simple scoring, the pipeline manages the generation of natural language justifications and citation mapping, ensuring that every compliance assessment is transparent, auditable, and grounded in the specific document segments identified by the reporting entity.

Researchers and developers can utilize this toolkit to deploy their own local instances or extend the framework to support additional sustainability standards beyond GRI.

### 4.2 System Demonstration

Building upon the proposed toolkit, we developed a web-based system to demonstrate ESG-KG’s capabilities in a real-world scenario. We present a compliance assessment case study using the *Delta Electronics 2023 ESG Report*. Figure 2 illustrates

<sup>2</sup><https://github.com/opendatalab/MinerU>

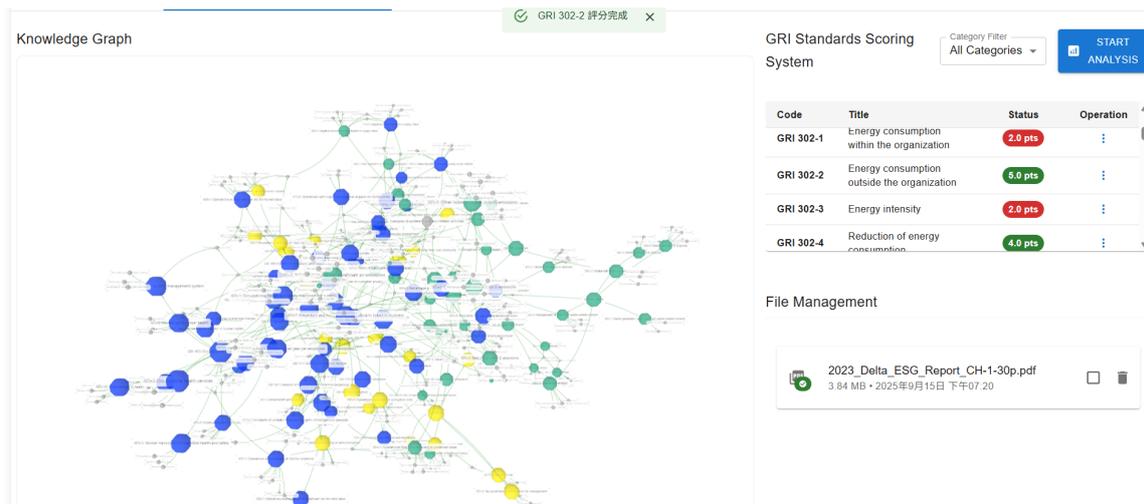


Figure 2: The ESG-KG system interface, featuring knowledge graph visualization for regulatory navigation and the automated compliance scoring panel for GRI standards

the system interface during the audit process, featuring a split-screen layout with the knowledge graph visualization on the left and the scoring control panel on the right.

**Document Upload and Processing.** As shown in the bottom-right “File Management” module, users upload reports such as the 2023\_Delta\_ESG\_Report\_CH-1-30p.pdf. The system automatically executes the multi-modal pipeline, decomposing the report into semantic units—text segments, tables, and infographics. These units are serialized, embedded, and indexed in a vector database to form the evidence pool.

**Knowledge Graph Visualization and Retrieval.** The left panel displays an interactive visualization of the regulatory Knowledge Graph, where nodes represent specific standards and compliance concepts. For compliance assessment against **GRI 302: Energy** standards, this graph serves as the navigational backbone. The system performs grounded retrieval to extract relevant evidence from the document while simultaneously querying the visualized GRI regulatory KG to retrieve corresponding logic, such as emission thresholds. This dual-graph alignment ensures that extracted evidence semantically maps to precise regulatory requirements.

**Automated Scoring with Traceable Evidence.** The top-right “GRI Standards Scoring System” presents the evaluation results. The model evaluates disclosure completeness and accuracy, assigning granular compliance scores as visible in the “Status” column: **2.0 pts** for GRI 302-1 (partial compliance), **5.0 pts** for GRI 302-2 (full com-

pliance), and **4.0 pts** for GRI 302-4. Real-time feedback is provided, as seen in the top overlay notification confirming the completion of the GRI 302-2 assessment. Users can filter results by category or initiate new evaluations using the “START ANALYSIS” button, enabling an efficient and interactive verification process.

**Interactive Exploration.** The interface allows users to explore the knowledge graph structure, trace evidence provenance, and review the reasoning chain connecting regulatory requirements to corporate disclosures. This transparency enables auditors to validate system decisions and identify areas requiring human expert review.

## 5 System Evaluation

### 5.1 Experimental Setup

To rigorously evaluate the retrieval performance of ESG-KG, we constructed a custom dataset named **ESG-50**. This dataset was collected by the authors and consists of publicly available sustainability reports from 50 representative companies in Taiwan. Published between 2024 and 2025, these reports cover a wide range of industrial sectors—including technology, finance, and manufacturing—and strictly adhere to the Global Reporting Initiative (GRI) standards, ensuring a diverse and standardized testbed for our experiments.

**Ground Truth Construction.** We leverage the *GRI Content Index* typically included in compliant reports, which explicitly maps each GRI disclosure item (e.g., GRI 302-1) to specific page numbers or sections. For a given standard  $\mathcal{S}$ , the ground

truth evidence set  $\mathcal{E}_{gt}$  consists of all text segments, tables, and charts located on the referenced pages.

**Baselines.** We compare ESG-KG against two established retrieval approaches:

- **BM25:** Keyword-based retrieval using exact term matching between GRI standard descriptions and document segments.
- **Dense retrieval:** Semantic search using a pre-trained embeddings (OpenAI text-embedding-3) without KG guidance or multi-modal processing.

## 5.2 Retrieval Performance

Table 1 presents the comparative results using Recall@K and NDCG@K metrics. BM25 achieves the lowest performance (Recall@5: 42.3%) due to the vocabulary mismatch between regulatory terminology and corporate reporting language. Dense retrieval improves substantially (Recall@5: 61.5%) by capturing semantic similarities, but struggles with quantitative data in tables and charts.

ESG-KG significantly outperforms both baselines, achieving Recall@5 of 84.1%. This improvement stems from two key capabilities: (1) **Multi-modal parsing** successfully retrieve evidence from tables, charts and infographics, which constitute a large portion of GRI data, and (2) **KG-guided hybrid retrieval** utilizes the structured “Standard Graph” to expand queries with semantically related requirements, ensuring retrieval captures logically relevant compliance evidence rather than merely semantically similar text.

Table 1: Retrieval Performance on ESG-50 Dataset.

| Metric    | BM25  | Dense | ESG-KG (Ours) |
|-----------|-------|-------|---------------|
| Recall@5  | 0.423 | 0.615 | <b>0.841</b>  |
| Recall@10 | 0.518 | 0.702 | <b>0.915</b>  |
| NDCG@10   | 0.387 | 0.594 | <b>0.812</b>  |

## 6 Conclusion

We presented ESG-KG, a system that automates ESG compliance assessment through multi-modal information extraction and knowledge graph construction. By integrating specialized multi-modal document processing with a structured GRI-based knowledge graph, ESG-KG enables precise, evidence-based compliance verification across textual, tabular, and visual content. Our demonstration

showcases how the system reduces manual assessment effort while enhancing transparency, accuracy, and auditability in ESG reporting. The system provides a practical foundation for scalable, standardized compliance validation aligned with global sustainability frameworks. Future work will extend coverage to additional ESG standards beyond GRI and incorporate user feedback mechanisms for iterative model refinement.

## 6.1 Limitations

While ESG-KG demonstrates effective automated compliance assessment, several limitations warrant consideration. First, the knowledge graph construction currently focuses on GRI Standards and does not cover all regional or industry-specific ESG frameworks (e.g., SASB, TCFD). Second, the multi-modal extraction pipeline may face challenges with highly unconventional document layouts or proprietary infographic formats. Third, although the LLM-based scoring provides generally consistent assessments, it may require human oversight for edge cases involving complex regulatory interpretations. Finally, the system has been primarily developed and evaluated on English-language reports; multilingual support remains an area for future development.

Future work will focus on expanding standard coverage, improving robustness to diverse document formats, and incorporating mechanisms for human-in-the-loop validation in ambiguous cases.

## 6.2 Ethics Statement

ESG-KG is designed as a decision-support tool to assist auditors and stakeholders in evaluating ESG disclosures, serving to augment rather than replace human judgment or definitive compliance determinations. We emphasize that automated compliance systems should support expert analysis, particularly in contexts involving complex regulatory interpretations or significant stakeholder impact.

We acknowledge the responsibility to mitigate potential biases inherited from training data or reporting standards. To this end, the system’s transparency features—including traceable evidence chains and explicit reasoning—are intentionally designed to enable human reviewers to validate automated assessments and identify potential errors. Ultimately, users must critically evaluate system outputs and retain full accountability for final compliance decisions.

## Acknowledgements

This research was supported in part by an industrial collaboration project with Delta Electronics, Inc.

## References

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 1877–1901.
- Yingjie Chen, Hao Liu, Yulu Liu, Jiawei Xie, Ruikang Yang, Hanming Yuan, Yanjun Fu, Peter Y. Zhou, Qi Chen, James Caverlee, and Irwin Li. 2025. GraphCheck: Breaking long-term text barriers with extracted knowledge graph-powered fact-checking. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, pages 14976–14995.
- Tanay Gupta, Tushar Goel, and Ishan Verma. 2025. [Exploring multimodal language models for sustainability disclosure extraction: A comparative study](#). In *The Sixth Workshop on Insights from Negative Results in NLP*, pages 141–149, Albuquerque, New Mexico. Association for Computational Linguistics.
- Chaoyue He, Xin Zhou, Yi Wu, Xinjia Yu, Yan Zhang, Lei Zhang, Di Wang, Shengfei Lyu, Hong Xu, Wang Xiaoqiao, Wei Liu, and Chunyan Miao. 2025. [ESGenius: Benchmarking LLMs on environmental, social, and governance \(ESG\) and sustainability knowledge](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 14623–14664, Suzhou, China. Association for Computational Linguistics.
- Conghui He, Wei Li, Zhenjiang Jin, Chao Xu, Bin Wang, and Dahua Lin. 2024. Opendatalab: Empowering general artificial intelligence with open datasets. *arXiv preprint arXiv:2407.13773*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves 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.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [FACTScore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100.
- Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, Zhenjiang Jin, Guang Liang, Rui Zhang, Wenzheng Zhang, Yuan Qu, Zhifei Ren, Yuefeng Sun, Yuanhong Zheng, Dongsheng Ma, Zirui Tang, Boyu Niu, Ziyang Miao, Hejun Dong, Siyi Qian, Junyuan Zhang, Jingzhou Chen, Fangdong Wang, Xiaomeng Zhao, Liqun Wei, Wei Li, Shasha Wang, Ruiliang Xu, Yuanyuan Cao, Lu Chen, Qianqian Wu, Huaiyu Gu, Lindong Lu, Keming Wang, Dechen Lin, Guanlin Shen, Xuanhe Zhou, Linfeng Zhang, Yuhang Zang, Xiaoyi Dong, Jiaqi Wang, Bo Zhang, Lei Bai, Pei Chu, Weijia Li, Jiang Wu, Lijun Wu, Zhenxiang Li, Guangyu Wang, Zhongying Tu, Chao Xu, Kai Chen, Yu Qiao, Bowen Zhou, Dahua Lin, Wentao Zhang, and Conghui He. 2025. [Mineru2.5: A decoupled vision-language model for efficient high-resolution document parsing](#).
- Alec Radford and Karthik Narasimhan. 2018. Improving language understanding by generative pre-training. In *OpenAI Technical Report*.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. pages 8–9.
- 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. In *Journal of Machine Learning Research*, pages 1–67.
- Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, Bo Zhang, Liqun Wei, Zhihao Sui, Wei Li, Botian Shi, Yu Qiao, Dahua Lin, and Conghui He. 2024a. [Mineru: An open-source solution for precise document content extraction](#).
- Jinhong Wang, Qixiu Song, Li Qian, Hang Li, Qinghua Peng, and Jianming Zhang. 2024b. [SubstationAI: Multimodal large model-based approaches for analyzing substation equipment faults](#). *CoRR*, abs/2412.17077.
- Qintong Zhang, Bin Wang, Victor Shea-Jay Huang, Junyuan Zhang, Zhengren Wang, Hao Liang, Conghui He, and Wentao Zhang. 2024. Document parsing unveiled: Techniques, challenges, and prospects for structured information extraction. *arXiv preprint arXiv:2410.21169*.

# BanSuite: A Unified Toolkit and Software Platform for Low-Resource NLP in Bangla

Md. Abu Sayed<sup>1\*</sup>, Faisal Ahamed Khan<sup>1\*</sup>, Jannatul Ferdous Tuli<sup>1\*</sup>,  
Nabeel Mohammed<sup>3</sup>, Mohammad Ruhul Amin<sup>4</sup>, Mohammad Mamun Or Rashid<sup>5</sup>

<sup>1</sup>Giga Tech Limited, Dhaka, Bangladesh, <sup>3</sup>North South University, Dhaka, Bangladesh  
<sup>4</sup>Fordham University, New York, USA, <sup>5</sup>Bangladesh Computer Council, Dhaka, Bangladesh

Correspondence: [faisal.cse06@gigatechltd.com](mailto:faisal.cse06@gigatechltd.com) \*These authors share first authorship

## Abstract

Bangla is one of the world’s most widely spoken languages, yet it remains significantly under-resourced in natural language processing (NLP). Existing efforts have focused on isolated tasks such as Part-of-Speech (POS) tagging and Named Entity Recognition (NER), but comprehensive, integrated systems for core NLP tasks including Shallow Parsing and Dependency Parsing are largely absent. To address this gap, we present BanSuite, a unified Bangla NLP ecosystem developed under the EBLICT project. BanSuite combines a large-scale, manually annotated Bangla Treebank with high-quality pretrained models for POS tagging, NER, shallow parsing, and dependency parsing, achieving strong in-domain baseline performance (POS: 90.16 F1, NER: 90.11 F1, SP: 86.92 F1, DP: 90.27 UAS). The system is accessible through a Python toolkit (Bkit) and a Web Application, providing both researchers and non-technical users with robust NLP functionalities, including tokenization, normalization, lemmatization, and syntactic parsing. In benchmarking against existing Bangla NLP tools and multilingual Large Language Models (LLMs), BanSuite demonstrates superior task performance while maintaining high efficiency in resource usage. By offering the first comprehensive, open, and integrated NLP platform for Bangla, BanSuite lays a scalable foundation for research, application development, and further advancement of low-resource language technologies. A demonstration video is provided to illustrate the system’s functionality in <https://youtu.be/3pcfiUQfCoA>

## 1 Introduction

Bengali, or Bangla, is a classical Indo-Aryan language of the Indo-European family, native to the Bengal region of South Asia (Cardona and Jain, 2007). It has about 242 million native speakers, and 43 million second-language speakers (Eberhard et al., 2024), making it one of the world’s

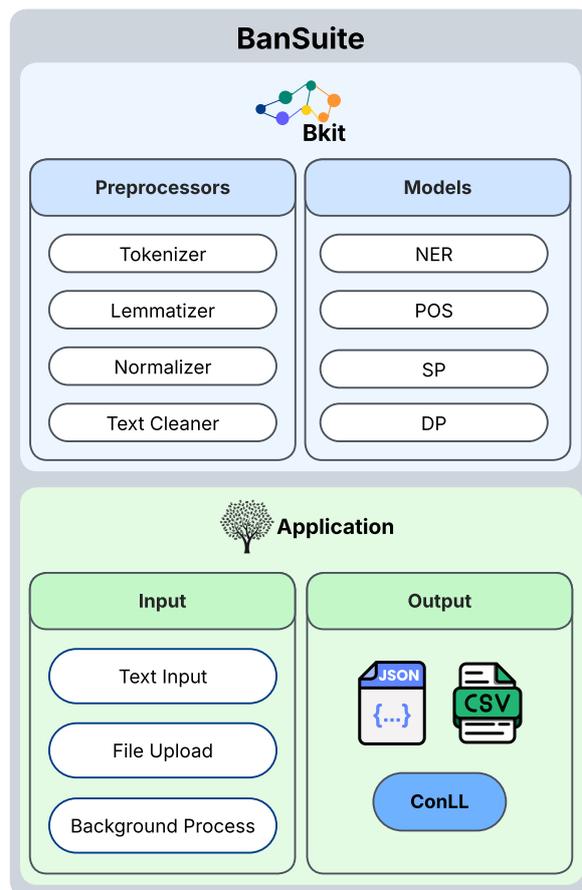


Figure 1: **BanSuite** comprises two complementary components. The **Bkit** package delivers core NLP functionalities. The **Application** provides an interactive NLP prediction platform, enabling users to upload files and obtain results in multiple downloadable formats, thereby facilitating both research and practical usage.

most spoken languages (Eberhard et al., 2024). It is the official language of Bangladesh, where around 98% of the population speaks it natively (Eberhard et al., 2024). In India, Bengali is the second most spoken language and serves as an official language in West Bengal, Tripura, and the Barak Valley of Assam. Since 2011, it has also been recognised as the second official language of Jharkhand.

Despite this vast number of speakers, the Bengali language is unfortunately still a low-resource language in Natural Language Processing (NLP). Although some important work on various core components of the Bengali language has been done (Arora, 2020) (Sarker, 2021) sporadically at the research level over the past decade, a fundamental gap has always been evident. Unlike English, which has powerful and integrated ecosystems like Stanford CoreNLP (Manning et al., 2014) or spaCy (Honnibal et al., 2020) that serve as reliable platforms for researchers, academics, and technology enthusiasts interested in the language, a single, open, and integrated platform for Bengali has been absent until now. The Comparison of popular NLP libraries by their Bangla support illustrated in **Table 1**.

We are introducing this integrated Bengali comprehensive Core NLP ecosystem **Figure 1**, built under the Enhancement of Bangla Language in ICT through Research and Development (EBLICT) project, to the academic and research community for the first time. In this article, we have discussed in detail how our system was prepared, what its architecture is, and what core components and resources have been made available for use under this platform.

The process of preparing our system was completed in several steps. First, an annotation platform was created where gold-standard data was annotated by experienced and certified linguists. This annotated data was used to train our core Bengali NLP models, such as Named Entity Recognition (NER), Part-of-Speech (POS) tagging, Shallow Parsing (SP), and Dependency Parsing (DP). The trained models were then deployed into the software system <sup>1</sup> and integrated into the Bangla Text Processing Kit (Bkit) <sup>2</sup>. Bkit is a comprehensive Python library that provides a wide range of functionalities for Bangla NLP, including word tokenization, sentence tokenization, normalization, text cleaning, and lemmatization.

The landscape analysis highlights a significant functional gap in accessible, application-ready advanced systems for the Bangla language. Existing major Bangla toolkits (Sarker, 2021; Arora, 2020) either focus primarily on foundational pre-processing or provide limited support for advanced

tasks, without offering high-performance, unified models that are readily usable for complex NLP applications. The BanSuite system is designed to explicitly address this gap by providing robust, dedicated models for resource-intensive, high-value tasks through a consolidated interface.

- We present the first open and unified Bengali core NLP platform, integrating all core NLP components under a single framework.
- Our Pretrained Language Models (PLMs) were trained on a linguist-certified Bengali Treebank. The inter-annotator agreement (IAA) scores, measured using Fleiss’  $\kappa$  (Fleiss, 1971), are: POS:  $\kappa = 0.77$ , NER:  $\kappa = 0.92$ , SP:  $\kappa = 0.68$ , and DP:  $\kappa = 0.68$ .
- The developed Bangla NLP models achieve n-domain baseline performance scores of POS: 90.16% F1, NER: 90.11% F1, SP: 86.92% F1, and DP: 90.27% UAS on test sets, demonstrating strong baseline results for Bengali NLP tasks.
- We provide access for both technical and non-technical users through two interfaces: the Python API and the Web Application.
- Bridging the low-resource gap creates a foundational ecosystem akin to NLTK, spaCy, or Stanford CoreNLP for Bengali, offering broader core NLP functionalities.

## 2 Related Work

Early efforts in NLP tool development, exemplified by the Natural Language Toolkit (NLTK) (Bird et al., 2009), provided essential computational interfaces for core linguistic operations such as tokenization, stemming, and access to lexical resources. As the field increasingly shifted toward deep learning, libraries like spaCy (Honnibal et al., 2020) and research-oriented frameworks such as AllenNLP (Gardner et al., 2018) standardized the use of high-performance neural architectures for routine NLP tasks, thereby enabling rapid NLP experimentation.

The Stanford CoreNLP (Manning et al., 2014), widely regarded as a contemporary gold standard for system demonstrations. CoreNLP introduced a unified, language-agnostic, fully neural pipeline for comprehensive text analysis across multiple languages, covering tasks such as parser, tagger. A Python version of CoreNLP, Stanza (Qi et al.,

<sup>1</sup><https://corpus.bangla.gov.bd/corpus/ml-model/tree-bank>

<sup>2</sup><https://github.com/eblict-gigatech/bangla-text-processing-kit>

| Name                 | Stanza                    | NLTK                      | AllenNLP                  | Flair                     | SpaCy                     | Stanford CoreNLP          | iNLTK                           | BNLP             | BanSuite              |
|----------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------------|------------------|-----------------------|
| <b>Language</b>      | Multilingual (Bangla N/A) | Multilingual (Bangla Supported) | Dedicated Bangla | Dedicated Bangla      |
| <b>Accessibility</b> | Python Package            | App.                      | Python Package                  | Python Package   | Python Package + App. |

Table 1: Comparison of popular NLP libraries by their Bangla support: dedicated, multilingual, or none (Sarker, 2021; Arora, 2020; Qi et al., 2020; Gardner et al., 2018). Although (Honnibal et al., 2020) lists Bangla support, we found no Bangla-specific functionality. The table also distinguishes toolkit-based libraries from standalone applications.

| Feature       | BNLP | iNLTK | Bkit (Ours) | App. (Ours) |
|---------------|------|-------|-------------|-------------|
| Tokenizer     | ✓    | ✓     | ✓           |             |
| Embedding     | ✓    | ✓     | ✓           |             |
| Text Cleaning | ✓    | ✓     | ✓           |             |
| Normalizer    | ✓    |       | ✓           |             |
| Lemmatizer    |      |       | ✓           |             |
| POS           | ✓    |       | ✓           | ✓           |
| NER           | ✓    |       | ✓           | ✓           |
| SP            |      |       | ✓           | ✓           |
| DP            |      |       | ✓           | ✓           |
| Visualization |      |       | ✓           | ✓           |

Table 2: Comparison of key features in existing Bangla NLP tools (Sarker, 2021; Arora, 2020) and our toolkit, showing coverage of tokenization, embeddings, text cleaning, normalization, lemmatization, and models.

2020), was later developed as a comprehensive library that not only provides core NLP models but also supports fundamental tasks, including tokenization, multi-word token expansion, and lemmatization for multiple languages.

An Indic language-based work is the iNLTK library (Arora, 2020), which aims to lower the barrier for applied research across 13 Indic languages, including Bangla. (Arora, 2020) provides pre-trained models and out-of-the-box support for essential functionalities such as tokenization, word embeddings, textual similarity, and data augmentation. However, (Arora, 2020) does not provide core NLP models and is limited to tasks such as embeddings and sentence similarity.

BNLP (Sarker, 2021) is a dedicated Bangla toolkit that consolidates core preprocessing and tagging capabilities, offering pre-trained models for tokenization, word and document embeddings, POS tagging, and NER. Its primary strength lies in providing essential linguistic analysis for higher-level applications. However, BNLP does not include syntactic parsers, limiting its functionality to foundational tagging tasks. **Table 2** demonstrates the comparison of key features in existing Bangla NLP tools.

### 3 Development of BanSuite

#### 3.1 Treebank

We use a large-scale, in-house annotated Bangla treebank to train all models presented in this work. The treebank provides multiple layers of linguistic annotation, NER, POS, SP, and DP, and is designed to ensure broad coverage and consistent annotation quality across diverse Bangla text genres. As shown in **Table 3**, the treebank contains roughly 170K sentences and 2M words for each task, along with high inter-annotator agreement (IAA) scores across annotation layers.

| Metric         | NER    | POS    | SP    | DP    |
|----------------|--------|--------|-------|-------|
| Sentence Count | 169K   | 172K   | 174K  | 174K  |
| Word Count     | 2.0M   | 2.0M   | 2.0M  | 2.0M  |
| IAA Score      | 92.60% | 77.30% | 68.15 | 68.54 |

Table 3: Overall statistics of our Bangla Treebank.

To ensure comprehensive linguistic and domain coverage, the dataset was sourced from a diverse collection of text types, including news, social media, literature, and conversational dialogues. The distribution of these domains across the NER, POS, SP, and DP datasets is illustrated in Figure 2. This diversity enables BanSuite models to generalize effectively across both formal and informal Bangla.

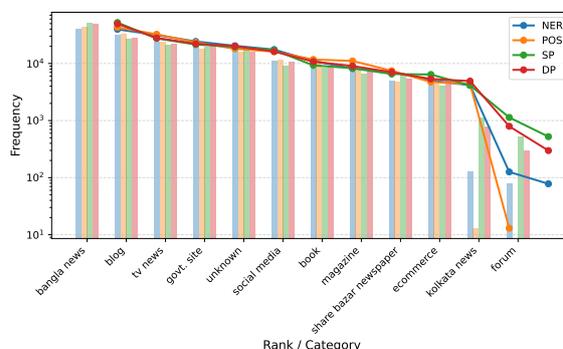


Figure 2: Distribution of source types across NER, POS, SP, and DP datasets. Bars represent absolute frequencies, while lines show log-scale rank–frequency trends.

### 3.2 Annotation and Validation

Across all four tasks, NER, POS, SP, and DP, we adopted a unified and systematic annotation framework (Islam et al., 2023), as also used by (Mahtab et al., 2025). Each dataset was annotated by trained annotators and subsequently validated by domain experts to ensure accuracy and consistency. Annotations were completed in small groups, where every sentence was independently labeled by three annotators following task-specific instructions. A majority-vote mechanism was used to determine the initial label for each instance. The assigned validator then examined the annotations, resolved inconsistencies, and finalized the labels visualized in **Figure 3**. A more detailed process is available in (Mahtab et al., 2025).

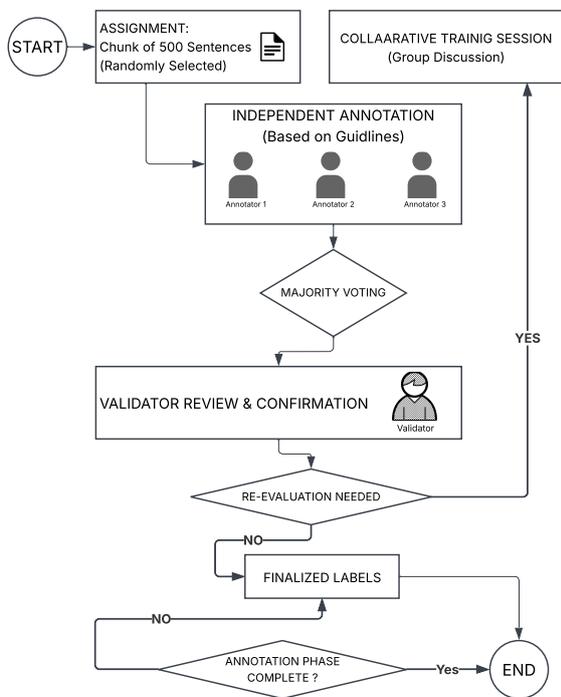


Figure 3: The annotation process involves assigning a chunk of 500 sentences for independent annotation (Mahtab et al., 2025). If re-evaluation is needed, collaborative training sessions are held for group discussions. The process includes majority voting and validator review to confirm labels, which leads to finalized labels or further re-evaluation. This cycle continues until the annotation phase is complete.

The number of annotation groups, assigned instances, validated instances, and acceptance rates are visualized in **Table 4**. These metrics reflect the strong quality-control measures used throughout the dataset construction process.

| Metric              | NER     | POS     | SP      | DP      |
|---------------------|---------|---------|---------|---------|
| Number of Groups    | 48      | 147     | 108     | 45      |
| Assigned Data       | 183,201 | 202,335 | 179,009 | 194,604 |
| Annotated Data      | 180,859 | 186,791 | 177,959 | 192,649 |
| Validated Data      | 170,261 | 173,123 | 174,784 | 174,415 |
| Rejected Data       | 10,595  | 11,015  | 3,174   | 13,328  |
| Annotators          | 86      | 234     | 184     | 95      |
| Validators          | 7       | 15      | 13      | 13      |
| Acceptance Rate (%) | 94.2    | 92.7    | 98.2    | 90.5    |

Table 4: Annotation workflow statistics for NER, POS, SP, and DP tasks. Acceptance Rate is computed as (Validated / Annotated).

### 3.3 Models

For part-of-speech tagging and named entity recognition, we adopt the noisy-label learning framework of (Liu et al., 2021), which explicitly models annotation uncertainty through confidence estimation. By applying confidence-aware regularization, the model mitigates the effect of incorrect labels and improves robustness and generalization under imperfect supervision.

For shallow parsing (chunking), we employ the self-attentive constituency parsing architecture proposed in (Kitaev and Klein, 2018) together with its multilingual extension (Kitaev et al., 2019). The use of global contextual representations via self-attention enables more accurate phrase boundary detection while supporting cross-lingual transfer.

For dependency parsing, we utilize a sequence-based generative formulation inspired by (Lin et al., 2022), which linearizes dependency trees into token sequences and predicts them using autoregressive decoding. This formulation simplifies structured prediction while maintaining strong syntactic accuracy and effectively modeling long-distance dependencies.

### 3.4 Training Hyperparameters

For each task, we selected a suitable pre-trained model based on prior literature and preliminary experiments: (Raffel et al., 2023) for NER, POS, and DP tasks, and (Devlin et al., 2019) for SP. We optimized hyperparameters using a combination of grid search and early stopping. Key hyperparameters tuned included the learning rate, number of epochs, gradient clipping norm, and maximum input sequence length. Gradient clipping with a maximum norm of 1.0 was applied to stabilize training, and gradient accumulation steps were used to simulate larger batch sizes within memory limits. The maximum input sequence length was set to 512

tokens for all tasks. Early stopping based on development set performance was applied to terminate training when no improvement was observed for a pre-defined number of evaluation steps. The final hyperparameter configurations, which yielded the best performance on the development sets, are shown in **Table 5**.

| Hyperparameter        | NER    | POS    | SP   | DP         |
|-----------------------|--------|--------|------|------------|
| Pre-trained Model     | T5 Enc | T5 Enc | BERT | T5 Seq2Seq |
| Learning Rate         | 1e-5   | 5e-5   | 3e-5 | 5e-5       |
| Epochs                | 40     | 40     | 10   | 25         |
| Max Grad Norm         | 1.0    | 1.0    | 1.0  | 1.0        |
| Gradient Accum. Steps | 5      | 5      | 1    | 1          |
| Sequence Length       | 512    | 512    | 512  | 512        |

Table 5: Hyperparameter configurations used for training the models across all tasks. Values shown are those that achieved the best development-set performance for NER, POS, SP, and DP

### 3.5 Bkit: Python Package

Bkit is a unified Python package for Bangla NLP that integrates preprocessing and model inference via a modular, consistent API. It supports text cleaning, tokenization, normalization, and lemmatization (Afrin et al., 2023), producing consistent lemmas while handling orthographic variability. Models can be initialized with a single class, automatically leveraging CPU/GPU resources. Built-in visualization tools facilitate interpretation of syntactic and semantic structures, especially for shallow and dependency parsing (Figure 4). Usage examples and API details are provided in **Appendix B**.

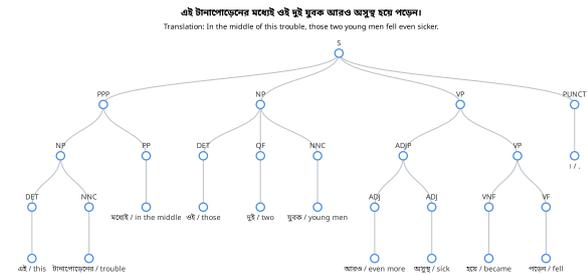


Figure 4: Shallow parsing visualization of a Bengali sentence. Nodes represent syntactic categories, leaf nodes correspond to words, and edges indicate parent-child relationships.

### 3.6 Web Application

BanSuite provides a web-based inference interface **Figure 5**. Users can perform real-time single-sentence inference or asynchronous batch processing via CSV uploads. Requests are managed through an API Gateway, routing real-time tasks

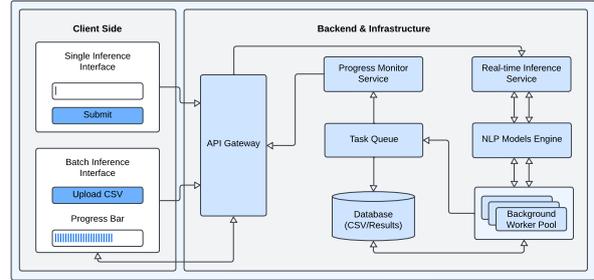


Figure 5: Overview of the Inference Service through application, Task Queue, and Background Worker, with core processing in the Models Engine and results stored in a Database for easy access.

to a synchronous service and batch jobs to worker pools, with results stored for access. Unlike traditional tools like Stanford CoreNLP (Manning et al., 2014), BanSuite requires no local setup and is accessible online up to a usage limit. Further application details are in **Appendix A**.

## 4 Result Analysis

### 4.1 In-domain Baseline Performance

Table 6 reports the in-domain performance of our NER, POS, SP, and DP models. The system achieves strong and consistent results, with development scores ranging from 86–91% and test set performance closely aligned. These results provide a reliable upper bound under matched training and evaluation conditions and serve as a reference for assessing generalization in zero-shot scenarios.

| Dataset      | NER   | POS   | SP    | DP    |
|--------------|-------|-------|-------|-------|
| Dev Set (%)  | 91.44 | 89.91 | 86.19 | 90.68 |
| Test Set (%) | 91.11 | 91.16 | 86.92 | 90.27 |

Table 6: In-domain performance of our models on the gold dataset.

### 4.2 Zero-shot Evaluation

We evaluate BanSuite in a zero-shot setting across four core Bengali NLP tasks: NER, POS, SP, and DP. Performance is compared against open-source Bangla toolkits (Sarker, 2021) and state-of-the-art multilingual LLMs, including Gemini 2.5F (Google DeepMind, 2025), Gemma 3 (27B) (Team et al., 2025), GPT-OSS-20B (OpenAI et al., 2025), and Llama 4 Scout (Meta AI, 2025). To explore cost-efficient alternatives, we also evaluate locally runnable variants of Gemma3<sup>3</sup>,

<sup>3</sup><https://ollama.com/library/gemma3:27b>

Llama4<sup>4</sup>, and GPT-OSS<sup>5</sup> via Ollama, which supports efficient execution and KV-caching (Pope et al., 2022).

Table 7 reports zero-shot results (F1 for NER, POS, SP; UAS for DP). BanSuite consistently achieves the highest scores: 88.78 F1 on NER, substantially outperforming Gemini 2.5F (56.29), Gemma 3 (35.59), and Llama 4 Scout (27.05). On POS, it averages 79.93 F1, ahead of BNLN (50.94–55.10), GPT-OSS-20B (64.01), and other open models. For syntactic tasks, BanSuite leads with 52.56 F1 on SP and 72.19 UAS on DP, while Gemini 2.5F and GPT-OSS-20B achieve moderate scores, and Llama 4 Scout trails. BNLN provides a lightweight competitive baseline, but BanSuite establishes a clear upper bound on zero-shot performance across Bangla NLP tasks.

#### 4.2.1 Few-Shot Analysis

LLMs are known to produce different outputs depending on the prompt (Zhuo et al., 2024). To study this effect, we evaluate Gemma3 27B under both zero-shot and few-shot settings. Figure 6 shows that providing a few examples in the prompt improves performance by  $\approx 11.37\%$  across all tasks. These results highlight that carefully designed prompts can influence the LLM’s behavior across different tasks.



Figure 6: Comparison of Gemma 3 27B performance in zero-shot and few-shot settings across four tasks. Adding 3 examples improves NER and SP performance, has a limited impact on POS, and big improvement in DP performance.

### 4.3 Inference Cost Analysis

The trade-offs between model performance and power consumption across the evaluated systems are shown in Figure 7. Large models such as

<sup>4</sup><https://ollama.com/library/llama4:16x17b>

<sup>5</sup><https://ollama.com/library/gpt-oss:20b>

Llama 4 Scout ( $\approx 41\%$ ), Gemma 3 27B ( $\approx 43\%$ ), and GPT-OSS 20B ( $\approx 51\%$ ) achieve moderate task performance but require over 200 W on average, reflecting substantial energy demands. In contrast, BanSuite (T5) delivers 72% accuracy while consuming under 100 W, demonstrating a favorable balance between performance and efficiency, see Appendix C

At the low-power extreme, BNLN (Lafferty et al., 2001) runs entirely on CPU using only  $\approx 30$  W, though its task performance is lower at  $\approx 27\%$ . Overall, this comparison highlights how model design influences both energy efficiency and effectiveness, allowing informed choices depending on deployment constraints.

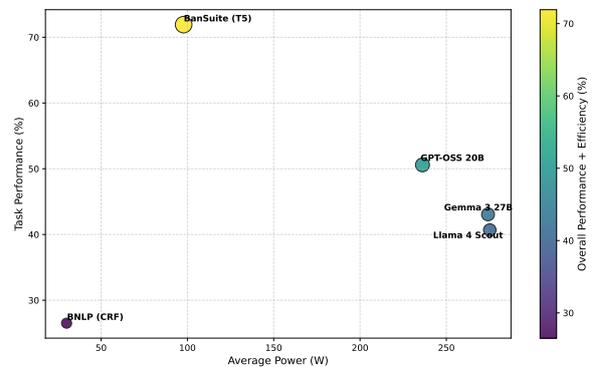


Figure 7: Trade-off between power usage and task performance across models. Point size and color indicate efficiency, with larger and darker points representing higher efficiency.

### Limitations

While BanSuite provides a unified toolkit and platform for Bangla NLP, several limitations remain. First, the system relies primarily on in-house annotated data, which, despite their size and diversity, may not capture all linguistic variations, dialects, or domain-specific usage in real-world Bangla text. Second, the current models are trained on standard and semi-formal text, so performance on highly informal or code-mixed data may be lower than reported benchmark results. Third, some complex linguistic phenomena, such as long-range dependencies in dependency parsing or nested entities in NER, remain challenging due to sequence length and model architecture limitations. Fourth, to evaluate existing models like (Sarker, 2021) for POS and NER, we mapped their data tags to our annotation scheme, which can introduce slight inconsistencies or misalignments. In contrast, no tag mapping was applied when evaluating LLMs, which

| Tasks | Datasets              | BNLP  | Llama4 Scout | Gemma3 27B | GPT-OSS-20B | Gemini 2.5F | BanSuite     |
|-------|-----------------------|-------|--------------|------------|-------------|-------------|--------------|
| NER   | (Mahtab et al., 2025) | 52.65 | 27.05        | 35.59      | 40.64       | 56.29       | 88.78        |
|       | (Haque et al., 2023)  | 66.14 | 50.76        | 50.53      | 53.86       | 54.76       | 79.80        |
|       | (Mhaske et al., 2023) | 46.50 | 65.86        | 66.25      | 73.60       | 77.65       | 82.50        |
|       | <b>Average</b>        | 55.10 | 47.22        | 50.12      | 55.97       | 60.20       | <b>83.03</b> |
| POS   | (Bali et al., 2010)   | 52.51 | 55.01        | 64.93      | 63.03       | 72.50       | 79.95        |
|       | (Kharagpur, 2005)     | 49.36 | 53.04        | 65.10      | 65.00       | 72.93       | 79.90        |
|       | <b>Average</b>        | 50.94 | 54.03        | 65.02      | 64.01       | 72.72       | <b>79.93</b> |
| SP    | (Mishra et al., 2024) | –     | 17.03        | 22.74      | 28.53       | 29.81       | <b>52.56</b> |
| DP    | (Jannat et al., 2021) | –     | 44.44        | 34.38      | 53.85       | 57.69       | <b>72.19</b> |

Table 7: Zero-shot evaluation (%) of available Bangla unified toolkit and LLMs across different tasks. F1 scores are reported for NER, POS, and SP tasks, while UAS is reported for DP. Scores are reported on multiple benchmark datasets, with averages computed where applicable. A dash (-) indicates that the feature is not available for the corresponding model. BanSuite consistently achieves the highest performance across all tasks.

may limit the comparability of their performance against BanSuite and (Sarker, 2021). Moreover, LLM performance can often be improved with carefully designed prompts (see Section 4.2.1), suggesting that their reported results may not fully reflect their potential with optimized prompting strategies. Finally, BanSuite currently focuses exclusively on Bangla and does not provide support like (Manning et al., 2014; Qi et al., 2020; Honnibal et al., 2020) for other low-resource languages, which limits its applicability to multilingual NLP research or cross-lingual tasks.

Future work could focus on expanding annotated corpora, improving support for code-mixed and informal text, refining tag mapping procedures, and adding built-in deployment features to Bkit for effortless deployment.

## Acknowledgments

We sincerely thank the *Bangla Syntactic Treebank Corpus with Processing Pipeline and Distribution Platform* project, part of the Enhancement of Bangla Language in ICT through Research and Development (EBLICT)<sup>6</sup>, supported by the Bangladesh Computer Council<sup>7</sup> under the ICT Division of the Government of Bangladesh<sup>8</sup>.

Our gratitude extends to the development consultant team at Giga Tech Limited<sup>9</sup> and the Department of Communication Disorders, University of

Dhaka<sup>10</sup> for their essential annotation support.

We also thank the Institute of Information Technology (IIT), University of Dhaka, for their valuable contribution as the testing team.

We gratefully acknowledge the financial support from the People’s Republic of Bangladesh, which enabled this research, and we thank the support teams at Facebook<sup>11</sup> and YouTube<sup>12</sup> for their assistance with accessibility and content integration.

## References

- Sadia Afrin, Md. Shahad Mahmud Chowdhury, Md. Islam, Faisal Khan, Labib Chowdhury, Md. Mahtab, Nazifa Chowdhury, Massud Forkan, Neelima Kundu, Hakim Arif, Mohammad Mamun Or Rashid, Mohammad Amin, and Nabeel Mohammed. 2023. *Ban-Lemma: A word formation dependent rule and dictionary based Bangla lemmatizer*. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3695–3710, Singapore. Association for Computational Linguistics.
- Gaurav Arora. 2020. *iNLTK: Natural language toolkit for indic languages*. In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, pages 66–71, Online. Association for Computational Linguistics.
- Kalika Bali, Monojit Choudhury, and Priyanka Biswas. 2010. Indian language part-of-speech tagset: Bengali (l1dc2010t16). Linguistic Data Consortium, Philadelphia. <https://catalog.ldc.upenn.edu/LDC2010T16>. Web download.
- <sup>6</sup><http://eblict.gov.bd/>
- <sup>7</sup><https://bcc.gov.bd/>
- <sup>8</sup><https://ictd.gov.bd/>
- <sup>9</sup><https://gigatechltd.com/>
- <sup>10</sup><https://www.du.ac.bd/body/LIN>
- <sup>11</sup><https://www.facebook.com/help/1020633957973118>
- <sup>12</sup><https://support.google.com/youtube/answer/9783148>

- Steven Bird, Edward Loper, and Ewan Klein. 2009. *Natural Language Processing with Python*. O'Reilly Media Inc.
- George Cardona and Dhanesh Jain. 2007. *The Indo-Aryan Languages*. Routledge.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **Bert: Pre-training of deep bidirectional transformers for language understanding**. *Preprint*, arXiv:1810.04805.
- David M Eberhard, Gary F Simons, and Charles D Fenig. 2024. *Ethnologue: Languages of the World*, 27th edition. SIL International, Dallas, Texas.
- Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2018. Allennlp: A deep semantic natural language processing platform. *arXiv preprint arXiv:1803.07640*.
- Google DeepMind. 2025. Gemini 2.5 flash. <https://ai.google.dev/gemini-api/>. Accessed: 2025-11-21.
- Md. Zahidul Haque, Sakib Zaman, Jillur Rahman Saurav, Summit Haque, Md. Saiful Islam, and Mohammad Ruhul Amin. 2023. **B-ner: A novel bangla named entity recognition dataset with largest entities and its baseline evaluation**. *IEEE Access*, 11:45194–45205.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spacy: Industrial-strength natural language processing in python. <https://spacy.io>.
- Md. Ekramul Islam, Labib Chowdhury, Faisal Ahamed Khan, Shazzad Hossain, Md Sourave Hossain, Mohammad Mamun Or Rashid, Nabeel Mohammed, and Mohammad Ruhul Amin. 2023. **Sentigold: A large bangla gold standard multi-domain sentiment analysis dataset and its evaluation**. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '23, page 4207–4218, New York, NY, USA. Association for Computing Machinery.
- Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. 1991. **Adaptive mixtures of local experts**. *Neural Computation*, 3(1):79–87.
- Siratun Jannat, Mizanur Rahoman, Shafi Sourov, Janatul Ferdaousi, Syeda Shahzadi, and Daniel Zeman. 2021. Ud bengali bru (universal dependencies v2.9). [https://github.com/UniversalDependencies/UD\\_Bengali](https://github.com/UniversalDependencies/UD_Bengali). Accessed: 2025-11-19; License: CCBY-SA4.0.
- IIT Kharagpur. 2005. Indian language pos-tagged corpus. <https://www.nltk.org/api/nltk.corpus.reader.indian.html>. Accessed: 2025-11-11.
- Nikita Kitaev, Steven Cao, and Dan Klein. 2019. **Multilingual constituency parsing with self-attention and pre-training**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3499–3505, Florence, Italy. Association for Computational Linguistics.
- Nikita Kitaev and Dan Klein. 2018. **Constituency parsing with a self-attentive encoder**. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2676–2686, Melbourne, Australia. Association for Computational Linguistics.
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*, ICML '01, page 282–289, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Boda Lin, Zijun Yao, Jiabin Shi, Shulin Cao, Binghao Tang, Si Li, Yong Luo, Juanzi Li, and Lei Hou. 2022. **Dependency parsing via sequence generation**. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7339–7353, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Kun Liu, Yao Fu, Chuanqi Tan, Moshua Chen, Ningyu Zhang, Songfang Huang, and Sheng Gao. 2021. **Noisy-labeled NER with confidence estimation**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3437–3445, Online. Association for Computational Linguistics.
- Md. Motahar Mahtab, Faisal Ahamed Khan, Md. Ekramul Islam, Md. Shahad Mahmud Chowdhury, Labib Imam Chowdhury, Sadia Afrin, Hazrat Ali, Mohammad Mamun Or Rashid, Nabeel Mohammed, and Mohammad Ruhul Amin. 2025. **BanNERD: A benchmark dataset and context-driven approach for Bangla named entity recognition**. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 6807–6828, Albuquerque, New Mexico. Association for Computational Linguistics.
- Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. **The Stanford CoreNLP natural language processing toolkit**. In *Association for Computational Linguistics (ACL) System Demonstrations*, pages 55–60.
- Meta AI. 2025. **Llama 4 scout: A natively multimodal mixture-of-experts model**. Version 17B-16E, released April 5, 2025.

Arnav Mhaske, Harshit Kedia, Sumanth Doddapaneni, Mitesh M. Khapra, Pratyush Kumar, Rudra Murthy, and Anoop Kunchukuttan. 2023. [Naamapadam: A large-scale named entity annotated data for Indic languages](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10441–10456, Toronto, Canada. Association for Computational Linguistics.

Pruthwik Mishra, Vandan Mujadia, and Dipti Misra Sharma. 2024. [Multi task learning based shallow parsing for indian languages](#). *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 23(9).

OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sebastien Bubeck, and 108 others. 2025. [gpt-oss-120b gpt-oss-20b model card](#). *Preprint*, arXiv:2508.10925.

Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Anselm Levskaya, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. 2022. [Efficiently scaling transformer inference](#). *Preprint*, arXiv:2211.05102.

Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A Python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Preprint*, arXiv:1910.10683.

Siddharth Samsi, Dan Zhao, Joseph McDonald, Baolin Li, Adam Michaleas, Michael Jones, William Bergeron, Jeremy Kepner, Devesh Tiwari, and Vijay Gadeppally. 2023. [From words to watts: Benchmarking the energy costs of large language model inference](#). *Preprint*, arXiv:2310.03003.

Sagor Sarker. 2021. [Bnlp: Natural language processing toolkit for bengali language](#). *Preprint*, arXiv:2102.00405.

Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. [Gemma 3 technical report](#). *Preprint*, arXiv:2503.19786.

Jingming Zhuo, Songyang Zhang, Xinyu Fang, Haodong Duan, Dahua Lin, and Kai Chen. 2024. [ProSA: Assessing and understanding the prompt sensitivity of LLMs](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages

1950–1976, Miami, Florida, USA. Association for Computational Linguistics.

## Appendix

### A Web Application

#### A.1 Single Inference

BanSuite provides an interactive web interface that allows users to perform single-sentence inference directly from a text input field. The interface supports real-time model execution and immediate visualization of results, making it suitable for quick experimentation and model inspection. See **Figure 8**.

Figure 8: Web-based interface for single-sentence inference.

**Named Entity Recognition** After processing the input sentence, the output is presented in a highly intuitive format where the original sentence is displayed, and each recognized entity is visually tagged and highlighted directly within the text. For example, a geographical location might be tagged as GPE (Geopolitical Entity), an event as EVENT, and a person’s name as PER. See **Figure 9**.

Figure 9: NER analysis on a sample Bangla sentence within the Tree Bank module: the output highlights each identified entity with its assigned tag

**Part-of-Speech** tagging, which assigns grammatical categories to every word in the input sentence, such as noun, verb, adjective, or determiner. The

resulting structured output is fundamental for facilitating deeper linguistic analysis, including syntactic analysis and word sense disambiguation. See **Figure 10**.

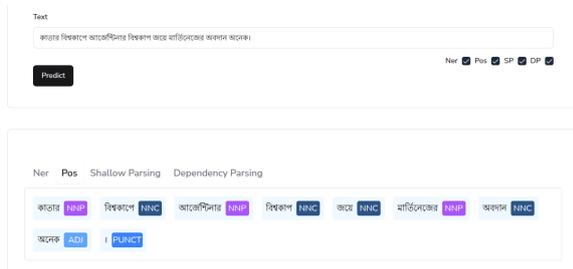


Figure 10: PoS tagging output showing each word in the sentence classified by its grammatical function

**Shallow Parsing** can also be done on the input text. This intermediate level of syntactic analysis aims to identify and segment the sentence into non-overlapping phrases or chunks. This process is instrumental in identifying relationships between major elements in a sentence and aids in both sentence comprehension and more complex downstream tasks. See **Figure 11**.

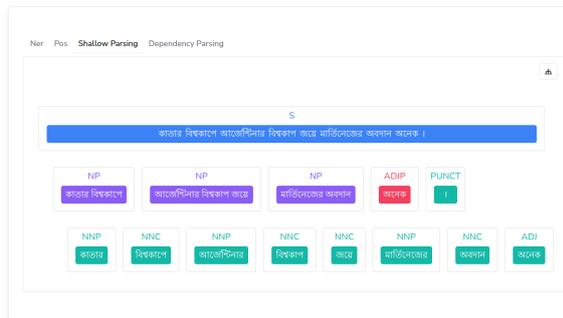


Figure 11: Visualization of Shallow Parsing where the recognized chunks, such as noun phrases (NP) and Adjective Phrases (ADJP), are marked right above the sequence of PoS-tagged words.

**Dependency Parsing** is the last layer of linguistic analysis provided by the treebank application. This method examines the grammatical links between words in a sentence. It frames the phrase as a directed graph, where each word is a node, and the relationship is represented by a named, directed arc from a head word to its dependent word. See **Figure 12**.

## A.2 Batch Inference

For large-scale processing, the BanSuite application supports batch inference through file uploads.

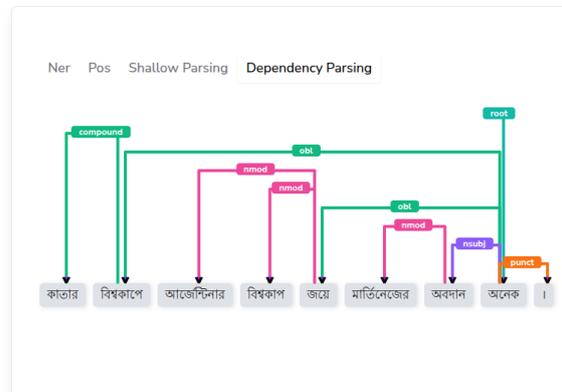


Figure 12: Visualization of a Dependency Parsing tree where the arcs are connecting the words, along with the specific relationship labels

Users can submit CSV files containing multiple input samples, which are processed asynchronously in the background. This feature enables the efficient handling of extensive datasets without requiring manual repetition. See **Figure 13**.

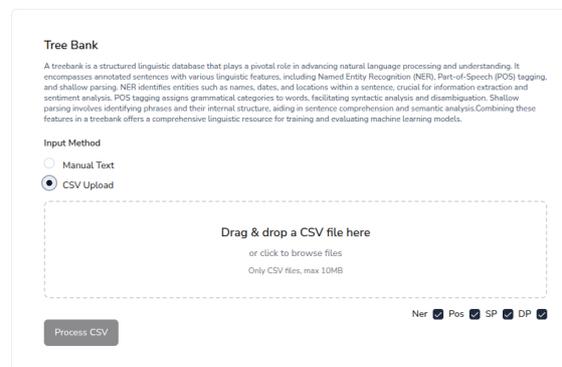


Figure 13: Option to upload a CSV file for batch inference.

During processing, the interface provides live progress feedback, allowing users to monitor ongoing inference tasks and retrieve results once computation is complete. See **Figure 14**.

## B Bkit

Bkit exposes a unified Python package that integrates preprocessing and model inference under a consistent API. It enables researchers to compose reproducible Bangla NLP pipelines with minimal

| File Name               | Task Type  | Progress                                                         | Status    | Processed | Created At            | Last Updated          | Actions                      |
|-------------------------|------------|------------------------------------------------------------------|-----------|-----------|-----------------------|-----------------------|------------------------------|
| bangla_sentence_100.csv | Dependency | <div style="width: 100%;"><div style="width: 100%;"></div></div> | Completed | 100 / 100 | Oct 5, 2025, 11:32 AM | Oct 5, 2025, 11:32 AM | <a href="#">View Details</a> |
| bangla_sentence_100.csv | Pos        | <div style="width: 100%;"><div style="width: 100%;"></div></div> | Completed | 100 / 100 | Oct 5, 2025, 11:32 AM | Oct 5, 2025, 11:32 AM | <a href="#">View Details</a> |
| bangla_sentence_100.csv | Ner        | <div style="width: 100%;"><div style="width: 100%;"></div></div> | Completed | 100 / 100 | Oct 5, 2025, 11:32 AM | Oct 5, 2025, 11:32 AM | <a href="#">View Details</a> |
| bangla_sentence_100.csv | Shallow    | <div style="width: 100%;"><div style="width: 100%;"></div></div> | Completed | 100 / 100 | Oct 5, 2025, 11:32 AM | Oct 5, 2025, 11:32 AM | <a href="#">View Details</a> |

Figure 14: Dashboard view of uploaded file being processed in the background.

code while ensuring explicit configuration and deterministic behavior.

## B.1 Text Cleaning

The cleaning module standardizes text by removing unwanted symbols and normalizing Unicode. It supports configurable filters, allowing transparent preprocessing for reproducibility.

```
import bkit
bkit.transform.clean_text(
    text,
    remove_punctuations=True,
    remove_digits=True,
    remove_emojis=True,
    remove_non_bangla=True,
)
```

## B.2 Tokenizer

The tokenizer applies Bangla-specific rules and supports rule-based segmentation. It returns token spans and offsets for precise mapping to model outputs.

```
import bkit
text = "Some Bangla text here"
tokens = bkit.tokenizer.tokenize(text)
```

## B.3 Normalizer

Normalization mitigates orthographic variability in Bangla by harmonizing visually similar characters, diacritics, and spacing conventions.

```
import bkit
normalizer = bkit.transform.Normalizer(
    normalize_characters=True,
    normalize_zw_characters=True,
    normalize_halant=True,
    normalize_vowel_kar=True,
    normalize_punctuation_spaces=True
)
normalizer(text)
```

## B.4 Lemmatizer

The lemmatizer combines a rule- and dictionary-based system (Afrin et al., 2023) with lexical

lookup to generate consistent lemmas while preserving meaning-critical morphemes.

```
import bkit
lemmatized = bkit.lemmatizer.lemmatize(text)
```

## B.5 Inference Pipeline

The API is designed to be both simple and user-friendly. By importing a single class, users can initialize the model without additional configuration.

```
from bkit.ner import Infer
model = Infer('ner-noisy-label')
predictions = model(text)
```

```
from bkit.pos import Infer
model = Infer('pos-noisy-label')
predictions = model(text)
```

```
from bkit.shallow import Infer
model = Infer('pos-noisy-label')
predictions = model(text)
```

```
from bkit.dependency import Infer
model = Infer('dependency-parsing')
predictions = model(text)
```

## C Inference Cost Analysis

LLMs impose substantial computational demands during inference (Samsi et al., 2023), making efficiency a critical factor for deployment in low-resource or real-time scenarios. To quantify these costs, we measured power consumption and GPU utilization for four model families: Gemma 3 27B, Llama 4 Scout, GPT-OSS, and our proposed BanSuite (T5) model running identical inference workloads on 2× Nvidia A40 GPUs.

As shown in **Figure 15**, Gemma 3 and Llama 4 Scout exhibit the highest energy requirements, averaging 274.14 W and 275.18 W, respectively. GPT-OSS also maintains a substantial draw at 236.18 W. In contrast, BanSuite operates at just 97.81 W, reducing power consumption by a factor of  $\approx 2.4\times$  relative to GPT-OSS and  $\approx 2.8\times$  relative to Gemma 3 and Llama 4 Scout.

**Figure 16** illustrates the corresponding utilization. Gemma 3 drives the GPUs to consistently high load, averaging 84.37% utilization with brief synchronization-related dips. Llama 4 Scout (Meta AI, 2025) shows considerably lower but more variable usage, averaging 28.48%, reflecting the intermittent expert activation characteristic

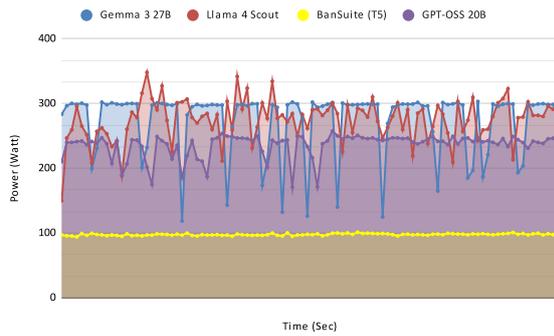


Figure 15: Power consumption of different models during NER inference. Gemma 3 and Llama 4 Scout have the highest usage, GPT-OSS is moderate, and BanSuite is the most efficient.

of Mixture of Experts (MoE) architectures (Jacobs et al., 1991). GPT-OSS 20B (OpenAI et al., 2025), which also follows this active expert (Jacobs et al., 1991), displays similar utilization at 76.24%.

BanSuite, despite its competitive performance, maintains a stable and low utilization profile, averaging 23.50%. This smooth utilization trace mirrors its compact 0.24B parameter footprint and efficient T5 (Raffel et al., 2023).

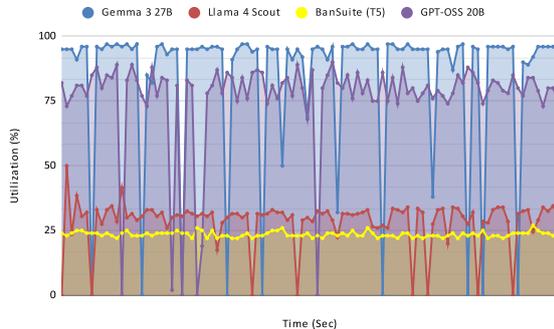


Figure 16: GPU utilization over time for Gemma 3 27B, Llama 4 Scout, and BanSuite (T5). Gemma 3 maintains consistently high load with brief dips, Llama 4 Scout shows moderate but variable usage due to its MoE architecture, and BanSuite (T5) exhibits stable low utilization, reflecting its compact and efficient design.

# Author Index

- Abramov, German, 297  
Abulkhanov, Dmitry, 297  
Agarwal, Shiven, 417  
Ahmad, Irfan, 225  
Ahnert, Georg, 537  
Akti, Seymanur, 175  
Al-Badrashiny, Mohamed, 341  
Al-shaibani, Maged S., 225  
Alheraki, Mais, 225  
Alomari, Nawaf, 225  
Alturki, Mustafa, 225  
Alumäe, Tanel, 575  
Alyafeai, Zaid, 225  
Amin, Mohammad Ruhul, 609  
Anees, Abdul Basit, 341  
Anvekar, Tejas, 417  
Ashraf, Ahmed, 225  
Aßenmacher, Matthias, 282
- Bach, Stephen, 349  
Barnes, Jeremy, 550  
Baroudi, Séverin, 320  
Batalov, Artem, 297  
Belanec, Robert, 188  
Benassi, Riccardo, 86  
Berkane, Thomas, 583  
Bielikova, Maria, 188  
Bojarskaja, Emilia, 271  
Bontcheva, Kalina, 154  
Bordoloi, Priyanuj, 417  
Budenny, Semen, 271  
Burdisso, Sergio, 320  
Bölücü, Necva, 428
- Cabezudo, Marco Antonio Sobrevilla, 407  
Cabrio, Elena, 1  
Chang, Li-Yang, 602  
Chen, Chih-Ming, 602  
Chen, Meng, 33  
Chen, Yiyang, 320  
Chenene, Mohamed, 1  
Chintha, Sai Sreenivas, 445  
Clematide, Simon, 203  
Contalbo, Michele Luca, 86  
Cyrta, Pawel, 320
- Daniélou, Jean, 1  
Dautov, Almaz, 297
- David, Jones, 129  
de Ioris, Roberto, 550  
Dejl, Adam, 480  
Dobrovoljc, Kaja, 75  
Du, Mingzhe, 46  
Duenser, Andreas, 428
- Eidt, Aaron Louis, 111  
Elneima, Ashraf Hatim, 341  
Emelyanov, Anton, 139  
Esfandiarpour, Reza, 349  
Estevanell Valladares, Ernesto Luis, 527  
Estève, Louis, 75
- Feldhus, Nils, 111, 163  
Fenogenova, Alena, 139  
Ferro, Markel, 550  
Feurer, Matthias, 282  
Fischer, Peter, 101  
Flores-Herr, Nicolas, 101  
Foster, Michael, 154  
Frei, Johann, 163  
Frontull, Samuel, 215  
Fujita, Tsuyoshi, 591
- Gahegan, Mark, 397  
Gashkov, Mikhail, 297  
Gein, Pavel, 297  
Ghosh, Himel, 261  
Ghosh, Shreya, 129  
Gollapalli, Sujatha Das, 46  
Grasso, Veronica, 492  
Grünert, David, 320  
Guerra, Francesco, 86  
Gupta, Vivek, 417
- Hachcham, Aymane, 492  
Hakam, Mouad, 46  
Haller-Seeber, Simon, 215  
Hamzeh, Mohammed, 46  
Heiß, Timo, 282  
Hewavitharana, Sanjika, 341  
Hilgert, Lukas, 175  
Huang, Hen-Hsen, 602  
Huber, Christian, 175
- Ijurco, Oier, 550  
Ionov, Timur, 21

Irons, Jessica, 428  
 Jagdale, Avani, 359  
 Jin, Brian, 428  
 Khan, Faisal Ahamed, 609  
 Khatib, Khalid Al, 563  
 Kochmar, Ekaterina, 457  
 Kokh, Vladimir, 271  
 Koneru, Sai, 175  
 Kong, Jiaming, 575  
 Kononov, Vasily, 21  
 Korzanova, Anna, 21  
 Kramer, Frank, 163  
 Kreutner, Maximilian, 537  
 Kruglikov, Vladislav, 297  
 Kudriashov, Sergei, 139  
 Kumar, Sachin, 359  
 Labrak, Yanis, 320  
 Lacalle, Oier Lopez De, 550  
 Lamsiyah, Salima, 512, 527  
 Latypov, Ramil, 297  
 Lee, Changhyun, 428  
 Leippold, Markus, 492  
 LI, Minghao, 397  
 Liu, Ruoqi, 359  
 Luis, Javier Alonso Villegas, 407  
 Luqman, Hamzah, 225  
 Madikeri, Srikanth, 320  
 Majumder, Maimuna S., 583  
 Maram, Durga Prasad, 445  
 Marxer, Ricard, 320  
 Maurer, Maximilian, 61  
 Maurya, Kaushal Kumar, 457  
 McCallum, Andrew, 445  
 Medvedev, Aleksandr, 297  
 Meyer, Alexander, 163  
 Michail, Andrianos, 203  
 Milner, Rosanna, 154  
 Milosevic, Nikola, 101  
 Mitkov, Ruslan, 527  
 Moeller, Lucas, 203  
 Mohammed, Nabeel, 609  
 Motliceck, Petr, 320  
 Msemo, Nakiete, 492  
 Munoz, Rafael, 527  
 Murtazin, Ruslan, 271  
 Murugaraj, Keerthana, 512  
 Muñoz, Iñigo Vilá, 550  
 Nabhani, Sara, 563  
 Naeem, Numaan, 457  
 Ng, See-Kiong, 46  
 Niehues, Jan, 175  
 Nikita, Surkov, 297  
 Nikolaev, Evgenii, 21  
 Novopoltsev, Maxim, 271  
 Olev, Aivo, 575  
 Opitz, Juri, 203  
 Padó, Sebastian, 203  
 Paganelli, Matteo, 86  
 Park, Jiwoo, 359  
 Patel, Dhruvesh, 445  
 Pearson, Jonathan, 480  
 Pederzoli, Sara, 86  
 Petukhova, Kseniia, 457  
 Picazo-Izquierdo, Alicia, 527  
 Podziubanchuk, Bohdan, 575  
 Porcellini, Valentin, 154  
 Potapov, Anatolii, 297  
 Raithel, Lisa, 163  
 Ranasinghe, Tharindu, 527  
 Rashid, Mohammad Mamun Or, 609  
 Razuvayevskaya, Olesya, 154  
 Refai, Dania, 225  
 Retkowski, Fabian, 175  
 Roberts, Ian, 154  
 Roller, Roland, 163  
 Rouhier, Jeanne, 1  
 Rozonoyer, Benjamin, 445  
 Rudat, Max, 101  
 Ruland, Timm Heine, 101  
 Rupprecht, Jens, 537  
 Rybinski, Maciej, 428  
 Sakai, Yusuke, 591  
 Sakhovskiy, Andrey, 271  
 Salaberria, Ander, 550  
 Salem, Ahmed, 537  
 Santacroce, Marta, 86  
 Sarkar, Mihir, 1  
 Savkin, Maksim, 21  
 Sawada, Yuya, 591  
 Sawaf, Hassan, 341  
 Sayed, Md. Abu, 609  
 Schlager, Moritz, 282  
 Schymura, Christopher, 101

Senni, Chiara Colesanti, 492  
Shah, Yash, 417  
Shekhar, Ashish Raj, 417  
Shen, Chang, 33  
Srba, Ivan, 188  
Srisuwananukorn, Andrew, 359  
Steinigen, Daniel, 101  
Stoianov, Dmitrii, 297  
Strohmaier, Markus, 537

Tan, Neseet, 397  
Taranets, Danil, 297  
Teucher, Roman, 101  
Teyssou, Denis, 154  
Theobald, Martin, 512  
Tsai, Ming-Feng, 602  
Tsymboi, Olga, 297  
Tuli, Jannatul Ferdous, 609  
Turkstra, Frieso, 563

Ugan, Enes Yavuz, 175  
Ulitin, Ivan, 271

Vaghefi, Saeid, 492  
Valeria, Venturelli, 86  
Vallabhaneni, Sunith, 583

Villatoro-tello, Esaú, 320

Waibel, Alexander, 175  
Wan, Stephen, 428  
Wang, Chuan-Ju, 602  
Wang, Jiayuan, 550  
Watanabe, Taro, 591  
Werner, Nick Elias, 261  
Wu, Kuizong, 33

Xu, Long, 33

Yang, Huichen, 428  
Yen, An-Zi, 602  
Yuan, Shaozu, 33  
Yuksel, Kamer Ali, 341

Zehle, Tom, 282  
Zelenkovskiy, Viktor, 297  
Zhang, Ping, 359  
Zhao, Jing, 359  
Ziletti, Angelo, 101  
Zuo, Max, 349