ClinicalNLP 2024

**The 6th Workshop on Clinical Natural Language Processing (ClinicalNLP)**

**Proceedings of the Workshop**

June 21, 2024

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

Order copies of this and other ACL proceedings from:

# Preface

This volume contains papers from the 6th Workshop on Clinical Natural Language Processing (Clinical NLP), held at NAACL 2024.

Much of the information recorded in a clinical encounter is located exclusively in provider narrative notes, which makes them indispensable for supplementing structured clinical data in order to better understand patient state and care provided. The goal of this workshop is to bring together researchers interested in improving NLP technology to enable clinical applications, focusing on information extraction and modeling of narrative provider notes from electronic health records, patient encounter transcripts, and other clinical narratives. This year, we received a total of 48 submissions to the main workshop, of which 8 were accepted as oral presentations, and 21 were accepted as poster presentations.

ClinicalNLP 2024 also hosted four shared tasks, challenging researchers around the world to develop new approaches to solve clinical NLP problems: medical error detection and correction, multilingual and multimodal medical answer generation, text-to-SQL modeling, and chemotherapy timelines extraction. In addition to the four task description papers from the four shared task organizers, we received a total of 34 participant submissions to the shared tasks, of which 4 were accepted as oral presentations, and 30 were accepted as poster presentations.

# Organizing Committee

**General Chairs**

Tristan Naumann, Microsoft Research
Asma Ben Abacha, Microsoft
Steven Bethard, University of Arizona
Kirk Roberts, UTHealth Houston
Danielle Bitterman, Harvard Medical School

**MEDIQA-CORR 2024 Shared Task Organizers**

Asma Ben Abacha, Microsoft
Wen-wai Yim, Microsoft
Meliha Yetisgen, University of Washington
Fei Xia, University of Washington

**MEDIQA-M3G 2024 Shared Task Organizers**

Asma Ben Abacha, Microsoft
Wen-wai Yim, Microsoft
Meliha Yetisgen, University of Washington
Fei Xia, University of Washington
Martin Krallinger, Barcelona Supercomputing Center (BSC)

**EHRSQL 2024 Shared Task Organizers**

Edward Choi, KAIST
Gyubok Lee, KAIST
Sunjun Kweon, KAIST
Seongsu Bae, KAIST

**ChemoTimelines 2024 Shared Task Organizers**

Jiarui Yao, Boston Children's Hospital, Harvard Medical School
Guergana Savova, Boston Children's Hospital, Harvard Medical School
Harry Hochheiser, University of Pittsburgh
WonJin Yoon, Boston Children's Hospital, Harvard Medical School
Eli Goldner, Boston Children's Hospital, Harvard Medical School

# Program Committee

**Program Chairs**

Asma Ben Abacha, Steven Bethard, Danielle Bitterman, Tristan Naumann, Kirk Roberts

**ChemoTimelines Reviewers**

Nesrine Bannour, Asma Ben Abacha, Shohreh Haddadan, Yukun Tan, Liwei Wang, Wen-wai Yim, Xingmeng Zhao

**EHRSQL Reviewers**

Asma Ben Abacha, Satya Kesav Gundabathula, Sourav Bhowmik Joy, Sangryul Kim, Oleg Somov, Jerrin John Thomas, Wen-wai Yim

**MEDIQA-CORR Reviewers**

Jean-Philippe Corbeil, Aryo Pradipta Gema, Satya Kesav Gundabathula, Hyeon Hwang, Suramya Jadhav, Gyubok Lee, Juan Pajaro, Swati Rajwal, Nadia Saeed, Augustin Toma, Airat Valiev, Zhaolong Wu, Jiarui Yao

**MEDIQA-M3G Reviewers**

Marie Bauer, Ricardo Omar Chávez García, Gyubok Lee, Nadia Saeed, Jerrin John Thomas, Augustin Toma, Parth Vashisht, Jiarui Yao

**Main Workshop Reviewers**

Ashwag Alasmari, Emily Alsentzer, Ibtihel Amara, Hajer Ayadi, Leonardo Campillos-Llanos, Shan Chen, Hong-Jie Dai, Dmitriy Dligach, Richard Dufour, Jocelyn Dunstan, Naome A Etori, Matúš Falis, Yadan Fan, Aryo Pradipta Gema, Zelalem Gero, John Michael Giorgi, Natalia Grabar, Mei-Hua Hall, Ming Huang, Raphael Iyamu, Qiao Jin, Yoshinobu Kano, Yejin Kim, Yanis Labrak, Egoitz Laparra, Alberto Lavelli, Ulf Leser, Qiuhao Lu, Yuxing Lu, Diwakar Mahajan, Sérgio Matos, George Michalopoulos, Timothy A Miller, Aurélie Névéol, Ankur Padia, Satyajeet Raje, Pavithra Rajendran, Giridhar Kaushik Ramachandran, Frank Rudzicz, Ashwyn Sharma, Sonish Sivarajkumar, Arvind Krishna Sridhar, Dhananjay Srivastava, Behrad Taghibeyglou, Khushboo Thaker, Augustin Toma, Jinge Wu, Susmitha Wunnava, Dongfang Xu, Wen-wai Yim, WonJin Yoon, Paul Youssef, Weipeng Zhou

# Keynote Talk
# Dual Edges of Innovation: Risks and benefits of LLMs in LMICs

### David Restrepo
Universidad del Cauca

**Abstract:** Large language models (LLMs) have emerged as transformative forces within artificial intelligence, heralding new capabilities in numerous sectors, including healthcare. Yet, the dialogue about their risks and their potential to widen social disparities, particularly in low-resource settings, remains insufficiently explored. In this keynote, we will dissect the evolution and fundamental principles of natural language processing (NLP), with a focus on the advent of transformative transformer models and their implications for fairness and bias.

We will start by outlining basic NLP concepts, progressively delving into how transformer models have reshaped our understanding of human-language machine interactions. This discussion will serve as a foundation to address the significant, yet often subtle, challenges of fairness and bias that are inherent in these models. The pervasive integration of advanced NLP technologies in clinical applications carries risks of perpetuating, or even exacerbating, existing biases which could profoundly affect patient care and outcomes.

The discourse will then shift to explore the advantages and practical applications of LLMs, with a focus on use cases in the Latin American context. Through specific examples, we will illustrate how LLMs can be leveraged to bridge language barriers and improve healthcare delivery in low-resource settings. Additionally, we will examine case studies from clinical settings across Latin America, highlighting the critical need for vigilance and the implementation of corrective measures to ensure these powerful tools serve all communities equitably.

**Bio:** David Restrepo is an Electronics and Communications Engineer and Data Scientist from Colombia, currently serving as a researcher at MIT Critical Data. He has also conducted significant research at the Laboratory for Computational Physiology at MIT, USA, and the University of Cauca in Colombia.

David's research is primarily focused on the application of machine learning in healthcare. He is particularly dedicated to addressing health inequalities and biases by developing methods for bias detection and de-biasing in medical images, text, and electronic health records (EHR) data. Additionally, he is actively involved in open data initiatives and events that aim to build capacity in the field.

His technical expertise includes efficient multimodal deep learning techniques that integrate medical images, textual data, and tabular datasets. Beyond his research, David is committed to mentoring and plays a pivotal role in organizing global datathons. These events promote collaborative data science and foster a diverse and interdisciplinary ecosystem in healthcare settings.

# Table of Contents

# Program

12:00 - 12:30     *Oral Session IV: Low-Resource Settings*

*A Privacy-Preserving Corpus for Occupational Health in Spanish: Evaluation for NER and Classification Tasks*
Claudio Aracena, Luis Miranda, Thomas Vakili, Fabián Villena, Tamara Quiroga, Fredy Núñez-Torres, Victor Rocco and Jocelyn Dunstan

*Leveraging Prompt-Learning for Structured Information Extraction from Crohn's Disease Radiology Reports in a Low-Resource Language*
Liam Hazan, Naama Gavrielov, Roi Reichart, Talar Hagopian, Mary-Louise C. Greer, Ruth Cytter-Kuint, Gili Focht, Dan Turner and Moti Freiman

12:30 - 14:00     *Lunch*

14:00 - 14:22     *Oral Session V: MEDIQA-CORR*

*Overview of the MEDIQA-CORR 2024 Shared Task on Medical Error Detection and Correction*
Asma Ben Abacha, Wen-wai Yim, Yujuan Fu, Zhaoyi Sun, Fei Xia and Meliha Yetisgen

*PromptMind Team at MEDIQA-CORR 2024: Improving Clinical Text Correction with Error Categorization and LLM Ensembles*
Satya Kesav Gundabathula and Sriram R Kolar

14:22 - 14:45     *Oral Session VI: MEDIQA-M3G*

*Overview of the MEDIQA-M3G 2024 Shared Task on Multilingual Multimodal Medical Answer Generation*
Wen-wai Yim, Asma Ben Abacha, Yujuan Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen and Martin Krallinger

*WangLab at MEDIQA-M3G 2024: Multimodal Medical Answer Generation using Large Language Models*
Augustin Toma, Ronald Xie, Steven Palayew, Gary D. Bader and BO Wang

14:45 - 15:07     *Oral Session VII: EHRSQL*

*Overview of the EHRSQL 2024 Shared Task on Reliable Text-to-SQL Modeling on Electronic Health Records*
Gyubok Lee, Sunjun Kweon, Seongsu Bae and Edward Choi

*LG AI Research & KAIST at EHRSQL 2024: Self-Training Large Language Models with Pseudo-Labeled Unanswerable Questions for a Reliable Text-to-SQL System on EHRs*
Yongrae Jo, Seongyun Lee, Minju Seo, Sung Ju Hwang and Moontae Lee

15:07 - 15:30    *Oral Session VIII: ChemoTimelines*

*Overview of the 2024 Shared Task on Chemotherapy Treatment Timeline Extraction*
Jiarui Yao, Harry Hochheiser, WonJin Yoon, Eli T Goldner and Guergana K Savova

*LAILab at Chemotimelines 2024: Finetuning sequence-to-sequence language models for temporal relation extraction towards cancer patient undergoing chemotherapy treatment*
Shohreh Haddadan, Tuan-Dung Le, Thanh Duong and Thanh Q. Thieu

15:30 - 16:00    *Break*

16:00 - 17:30    *Poster Session*

# Exploring Robustness in Doctor-Patient Conversation Summarization: An Analysis of Out-of-Domain SOAP Notes

**Yu-Wen Chen, Julia Hirschberg**

Department of Computer Science, Columbia University, United States

yu-wen.chen@columbia.edu, julia@cs.columbia.edu

## Abstract

Summarizing medical conversations poses unique challenges due to the specialized domain and the difficulty of collecting in-domain training data. In this study, we investigate the performance of state-of-the-art doctor-patient conversation generative summarization models on the out-of-domain data. We divide the summarization model of doctor-patient conversation into two configurations: (1) a general model, without specifying subjective (S), objective (O), and assessment (A) and plan (P) notes; (2) a SOAP-oriented model that generates a summary with SOAP sections. We analyzed the limitations and strengths of the fine-tuning language model-based methods and GPTs on both configurations. We also conducted a Linguistic Inquiry and Word Count analysis to compare the SOAP notes from different datasets. The results exhibit a strong correlation for reference notes across different datasets, indicating that format mismatch (i.e., discrepancies in word distribution) is not the main cause of performance decline on out-of-domain data. Lastly, a detailed analysis of SOAP notes is included to provide insights into missing information and hallucinations introduced by the models.

## 1 Introduction

Automatically generated summary notes of doctor-patient conversations could improve the healthcare system. First, the generated notes serve as a valuable resource, allowing doctors to review and validate the information from the conversation with a patient, ensuring that vital information is noticed. In addition, the summary notes can be integrated into hospitalization risk prediction models (Song et al., 2022), empowering healthcare professionals with data-driven insights to make more precise clinical decisions.

However, summarizing doctor-patient conversations poses distinct challenges owing to its specialized domain. Specifically, medical conversations often involve highly specialized terminology that requires domain-specific knowledge to understand and summarize accurately. In addition, it is preferable to structure the generated note with **S**ubjective (information reported by the patient), **O**bjective (objective observations), **A**ssessment (doctor's evaluation), and **P**lan (future care plan) (SOAP). SOAP format is preferable because it is widely utilized by healthcare providers to document a patient's progress, providing an organized framework that reduces communication confusion among healthcare professionals. These challenges hinder the direct application of general-purpose summarization techniques to doctor-patient conversations, underscoring the need for a specialized model.

Doctor-patient conversation summarization has attracted significant attention recently (Joshi et al., 2020; Krishna et al., 2021; Zhang et al., 2021; Grambow et al., 2022; Abacha et al., 2023a). In 2023, the MEDIQA-Chat Challenge (Abacha et al., 2023a) attracted 120 registered teams from the academy and industry. Although various methods are proposed in MEDIQA-Chat, it remains a challenging field that needs further investigation. First, MEDIQA-Chat focuses on in-domain training and testing. However, cross-dataset analysis for doctor-patient conversation summarization is crucial because collecting in-domain training data is usually challenging given the constraints imposed by privacy and security concerns. Second, a detailed assessment of performance across SOAP note categories is essential. Such insights into the performance of each category can play a pivotal role in developing improved model structures and designing more effective evaluation metrics.

In this study, we investigate cross-dataset performance of state-of-the-art (SOTA) doctor-patient summarization models. Our focus is on generative summarization models because the real-world clinical notes are in an abstractive format. The experiments were conducted on English datasets as

the setting of most previous studies. The results of SOAP notes are evaluated separately to gain a deeper understanding of the strengths and limitations of the current models. We hope our result can offer new insights for future research in developing a robust doctor-patient summarization model for real-world scenarios.

## 2 Related Work

The MEDIQA-Chat challenge (Abacha et al., 2023a) separated doctor-patient conversation summarization into different tasks. Models designed for Task A predict the topic category of the conversation and then generate notes. The Task A models are closer to a general-purpose summarization model, producing notes without specifying distinct sections. In the top performance models, Wanglab (Giorgi et al., 2023) fine-tuned a FLAN-T5 model (Chung et al., 2022) for summarization and note classification. SummQA (Mathur et al., 2023) used BioBERT (Lee et al., 2020) to support the section classification, MiniLM (Wang et al., 2020) to select the prompt for GPT4, and GPT4 to predict the section class and generated the final note. The Cadence (Sharma et al., 2023) model fine-tuned BART-large on the SAMSum dataset, followed by fine-tuning on the augmented dataset. In addition, a N-pass summarization was employed to handle long conversations.

Models designed for Task B are SOAP-oriented, generating notes with SOAP sections. In the top performance models, WangLab used instructor (Su et al., 2023) to select the top-k conversation that is similar to the testing data, then used the selected conversations and notes as the in-context learning examples for GPT4. They also achieved top performance with the fine-tuned Longformer Encoder-Decoder (LED) (Beltagy et al., 2020). SummQA (Mathur et al., 2023) used the MiniLM (Wang et al., 2020) to select the prompt for the GPT4 in-context learning examples as their model for task A. GersteinLab (Tang et al., 2023) used GPT-4 with specifically designed instruction.

Task A in the MEDIQA-Chat challenge was evaluated on the MTS-Dialog dataset (Abacha et al., 2023b), which has a relatively shorter conversation and reference notes related to a specific category. Task B was focused on the ACI-BENCH (Yim et al., 2023) dataset, which has a relatively longer conversation and a long note with SOAP sections. Most top-performance teams in Task A used fine-tuning

language model (LM)-based methods, while most top-performance teams in Task B introduced GPT-based approaches. The results seem to indicate that the fine-tuning LLM-based method is more suitable for *short* dialogues with a specific category of information. In contrast, the GPT-based method is preferable for the *long* dialogue with detailed SOAP information (Abacha et al., 2023a). However, in real-world scenarios, conversations may vary in length and encompass one or multiple categories of information. Therefore, in this study, we aim to understand how these models perform in an cross-dataset settings and identify potential errors made by the models.

## 3 Data

We use two open-source doctor-patient conversation datasets, MTS-Dialog (Abacha et al., 2023b) and ACI-BENCH (Yim et al., 2023). Both datasets contain doctor-patient conversations, the corresponding note of the conversation, and the category of the note. Figure 1 illustrates the samples in the two datasets, and Table 1 summarizes the dataset statistics. The number of tokens is calculated using the *google/flan-t5-large* tokenizer[1].

Compared with the two datasets, the MTS-Dialog dataset contains relatively shorter conversation, and the reference note follows a concise format, comprising either a few words or a one-paragraph structure with a section header specifying the note category. In contrast, the conversations in the ACI-BENCH dataset are relatively longer, and the reference notes includes all SOAP sections.

|  | Train | Valid | Test |
| --- | --- | --- | --- |
| Number of samples | | | |
| MTS-Dialog | 1,201 | 100 | 200 |
| ACI-BENCH | 67 | 20 | 40 |
| Number of tokens of dialogue (mean/max) | | | |
| MTS-Dialog | 152.4 / 2343 | 129.27 / 820 | 144.2 / 793 |
| ACI-BENCH | 1931.49 / 4642 | 1814.95 / 2608 | 1824.4 / 3560 |
| Number of tokens of note (mean/max) | | | |
| MTS-Dialog | 59.63 / 1580 | 53.9 / 406 | 57.4 / 530 |
| ACI-BENCH | 663.22 / 1388 | 680.3 / 1176 | 647.7 / 1291 |

Table 1: Statistic of MTS-Dialog and ACI-BENCH dataset.

We categorized the note in the MTS-Dialog dataset and divided the note in ACI-BENCH dataset into S, O, or AP categories for analysis. Note that we merged A and P as AP because these

---

[1] https://huggingface.co/google/flan-t5-large

| MTS-Dialog | ACI-BENCH |
|---|---|
| Dialogue | |
| Doctor: Good afternoon, sir. My chart here says that you are a fifty one year old white male, is that correct?<br>Patient: Good afternoon, doctor. Yes, all of that is correct.<br>...<br>Doctor: Finally, your ECOG score is one according to the nurse, is that correct?<br>Patient: Yes, doctor. That's correct. | Doctor: hi, andrew. how are you?<br>Patient: hey, good to see you.<br>Doctor: i'm doing well, i'm doing well.<br>…<br>Doctor: let me know if your symptoms worsen and we can talk more about it, okay?<br>Patient: you got it.<br>Doctor: all right. hey, dragon. finalize the note. |
| Note | |
| **Section header**: GENHX<br>**Section text**: A 51-year-old white male diagnosed with PTLD in latter half of 2007. He presented with symptoms of increasing adenopathy, abdominal pain, weight loss, and anorexia. …. | **CHIEF COMPLAINT**<br>Upper respiratory infection.<br>**HISTORY OF PRESENT ILLNESS**<br>Andrew Campbell is a 59-year-old male with a past medical history significant for depression, … |

Figure 1: Dataset examples. Samples in the MTS-Dialog dataset have a section header that indicates the category of the annotation and the section text, which is the main content of the notes. The samples in the ACI-BENCH dataset have one full note, where each section is separated by bold title text.

are merged into AP in the ACI-BENCH dataset, making it difficult to separate them into A and P. Table 2 shows the mapping between original note categories and SOAP and the number of samples in each category.

| Dataset | Original section | # of samples |
|---|---|---|
| **Subjective** | | |
| MTS-Dialog | GENHX, FAM/SOCHX, PASTMEDICALHX, CC, PASTSURGICAL, ALLERGY, ROS, MEDICATIONS, IMMUNIZATIONS, GYNHX, PROCEDURES, OTHER_HISTORY, | 175 |
| ACI-BENCH | Subjective: CHIEF COMPLAINT, HISTORY OF PRESENT ILLNESS, and REVIEW OF SYSTEMS. | 40 |
| **Objective** | | |
| MTS-Dialog | EXAM, IMAGING, LABS | 7 |
| ACI-BENCH | Objective exam and objective result: RESULTS, PHYSICAL EXAMINATION, and VITALS REVIEWED. | 40 |
| **Assessment and plan** | | |
| MTS-Dialog | ASSESSMENT, DIAGNOSIS, DISPOSITION, PLAN, EDCOURSE | 18 |
| ACI-BENCH | Assessment and plan: ASSESSMENT AND PLAN | 40 |

Table 2: Mapping between original note categories and SOAP.

## 4 Methods

We divided the summarization model for doctor-patient conversation into general and SOAP-oriented configurations (illustrated in Figure 2). In this study, we investigate the current SOTA models of each configuration in a cross-dataset setting. Our research question is:

**RQ1: How do current SOTA doctor-patient conversation summarization models perform on**

out-of-domain datasets, and what causes the performance decline?



Figure 2: Illustration of the general and SOAP-oriented configurations.

### 4.1 Cross-dataset analysis of general model

We analyzed the limitations of directly applying a general configuration for doctor-patient conversation summarization. Because the model does not consider generating S, O, A, and P notes separate tasks, the model may emphasize some information more than others, thus leading to missing information issues in the generated note. Therefore, we examined the following research question:

**RQ2: What information is more likely to be missing in SOAP for model with a general configuration? (Figure 3)** Our hypothesis is that objective information can easily be excluded from summaries. Objective information usually includes numerical information that holds significant importance in medical contexts. The number could represent the quantity of medication administered to the patient or the values derived from their health examination report, serving as indispensable metrics for assessing the patient's overall health condition. However, numerical data is often considered as detailed information and thus omitted in summaries. In addition, objective information is closely associated with technical terms, making it more challenging for the LM.

### 4.1.1 Model and Data

We used the fine-tuned Flan-T5 model (Chung et al., 2022), which received the top rank in the MEDIQA-Chat challenge task A, as representative for model with general configuration. The Flan-T5 model was fine-tuned with the MTS-Dialog dataset, in which the reference notes focus only on one topic in the conversation. We also included the GPT results (gpt-3.5-turbo and gpt4) for comparison. Models with the general configuration are

Figure 3: Analysis of fine-tuning LM-based general model.

evaluated on the ACI-BENCH dataset. Because the conversations and reference notes in the this dataset contain all SOAP information, we can analyze what categories of information (i.e., S, O, A, or P) are missing from the generated note.

## 4.2 Cross-dataset analysis of SOAP-oriented model

The model with SOAP-oriented configuration aims to generate notes with S, O, A, and P sections. However, in real-world conditions, not all doctor-patient conversations include all of the S, O, A, and P information. For example, doctors might skip the objective information because they already have the record. They might also not mention assessments and plans because they only want to check the patient's condition. Therefore, we ask the following research question:

**RQ3: What SOAP-oriented model will generate if the input conversation does not include information related to a specific category? (Figure 4)** We hypothesize that the LM will have severe hallucination problems by generating information that does not exist in the conversation.



Figure 4: Analysis of fine-tuning LM-based SOAP-oriented model.

### 4.2.1 Model and Data

We used the fine-tuned LED model (Beltagy et al., 2020), which received top-rank performance in the MEDIQA-Chat challenge task B as representative of the SOAP-oriented model. The LED model was fine-tuned with the ACI-BENCH dataset that specifies notes into SOAP sections. We also included the GPT results for comparison. The GPT was prompted to generate a note with SOAP sections and was informed that it could skip the section if no relevant information was provided in the conversation. We evaluated the models on the MTS-Dialog dataset, in which conversations are short and usually do not contain information related to all SOAP categories.

## 5 Experiments

### 5.1 Model details

We used WangLab's FLAN-T5 and LED summarization models in the MEDIQA-Chat Challenge [23]. To evaluate the FLAN-T5 model on input longer than its training data, we modify the maximum token length from 1024 to 4096. Table 3 shows the prompts for all models in the experiments. The prompts of FLAN-T5 and LED follow WangLab's settings. For GPT models, we followed LED and FLAN-T5 prompts but removed the "including family history, diagnosis, past medical (and surgical) history, and known allergies" to prevent GPTs from specifically clarifying that certain information is not part of the conversation. Lastly, we designed a prompt to guide GPT in generating a summary with SOAP sections and a more parsable format.

### 5.2 Evaluation metrics

All models were evaluated using ROUGE-1 (Lin, 2004) and the average of ROUGE-1, BLEURT (Sellam et al., 2020), and BERTScore (Zhang et al., 2020) (referred to as an aggregate score). These automatic metrics have been shown to correlate highly with human judgments for the doctor-patient conversations in recent studies (Abacha et al., 2023c). The section headers in the reference and generated notes were excluded from the evaluation. We used the *en_core_sci_sm* model in scispacy[4] to identify the medical terms in the dialogue

---

[2] https://huggingface.co/wanglab/task-a-flan-t5-large-run-2

[3] https://huggingface.co/wanglab/task-b-led-large-16384-pubmed-run-3

[4] https://allenai.github.io/scispacy/

| LED |
|---|
| Summarize the following patient-doctor dialogue. Include all medically relevant information, including family history, diagnosis, past medical (and surgical) history, immunizations, lab results and known allergies. Dialogue: {dialogue} |

| FLAN-T5 |
|---|
| Summarize the following patient-doctor dialogue. Include all medically relevant information, including family history, diagnosis, past medical (and surgical) history, immunizations, lab results and known allergies. You should first predict the most relevant clinical note section header and then summarize the dialogue. Dialogue: {dialogue} |

| GPT-{3.5, 4}-general (MTS-Dialog) |
|---|
| Summarize the following patient-doctor dialogue. Include all medically relevant information. You should first predict the most relevant clinical note section header and then summarize the dialogue. Dialogue: {dialogue} |

| GPT-{3.5, 4}-general (ACI-BENCH) |
|---|
| Summarize the following patient-doctor dialogue. Include all medically relevant information. Dialogue: {dialogue} |

| GPT-{3.5, 4}-SOAP |
|---|
| Summarize the following patient-doctor dialogue and structure the summary into (1) Subjective, (2) Objective, (3) Assessment and Plan sections. Avoid including information that is not explicitly mentioned in the conversation. If no related information for the section is provided, skip the section. For example, if no specific subjective information is provided in the dialogue, write "N/A" in the subjective section. Dialogue: {dialogue} |

Table 3: Model prompts.

and notes. Lastly, Linguistic Inquiry and Word Count (LIWC) (Tausczik and Pennebaker, 2010) was used to analyze the word distribution in SOAP notes. LIWC is a text analysis tool that systematically examines and categorizes language based on psychologically meaningful dimensions. It aids in deciphering the linguistic characteristics of written or spoken text, providing insights into the emotional and cognitive dimensions of communication. Because emotional and cognitive words can reflect aspects of a person's health in certain situations, they play essential roles in the SOAP note.

## 6 Results

### 6.1 Cross-dataset Performance

We evaluated the cross-dataset performance of doctor-patient conversation summarization models. Performance on the ACI-BENCH dataset is presented in Table 4. The experimental results indicate a notable performance decrease in out-of-domain models compared to the in-domain baseline (i.e., LED). We also noticed that the general model performed particularly poorly on objective notes. When utilizing the general model for doctor-patient

summarization, adaptations are essential to preserve objective information. A potential approach involves treating the generation of objective notes as a distinct task. For example, the outcomes from gpt-SOAP models indicate that the performance of objective notes increases greatly by specifically instructing the model to generate notes with an objective section.

| Testing data / Model | S | O | AP |
|---|---|---|---|
| ROUGE-1 | | | |
| LED (In-domain) | 0.554 | 0.502 | 0.491 |
| gpt3.5-SOAP | 0.358 (-35%) | 0.420 (-16%) | 0.381 (-22%) |
| gpt4-SOAP | 0.373 (-33%) | 0.447 (-11%) | 0.379 (-23%) |
| FLAN-T5 | 0.339 (-39%) | 0.146 (-71%) | 0.265 (-46%) |
| gpt3.5-general | 0.349 (-37%) | 0.175 (-65%) | 0.352 (-28%) |
| gpt4-general | 0.370 (-33%) | 0.179 (-64%) | 0.363 (-26%) |
| Aggregate score | | | |
| LED (In-domain) | 0.569 | 0.538 | 0.546 |
| gpt3.5-SOAP | 0.494 (-13%) | 0.527 (-2%) | 0.520 (-5%) |
| gpt4-SOAP | 0.504 (-11%) | 0.552 (+2%) | 0.518 (-5%) |
| FLAN-T5 | 0.447 (-21%) | 0.350 (-35%) | 0.407 (-25%) |
| gpt3.5-general | 0.478 (-16%) | 0.384 (-29%) | 0.479 (-12%) |
| gpt4-general | 0.487 (-14%) | 0.395 (-27%) | 0.482 (-12%) |

Table 4: Model performance on the ACI-BENCH dataset. Testing data S, O, and AP means the evaluated reference note is the subjective, objective, and assessment and plan sections of the original reference note, respectively. The values in parentheses indicate the performance change compared with in-domain LED model (i.e., LED fine-tuned on ACI-BENCH). The FLAN-T5 model is fine-tuned on the MTS-Dialog dataset.

Table 5 shows performance on the MTS-Dialog dataset. Because the reference in the MTS-Dialog dataset only focuses on one category, we ignore unmatched sections of the generated note. For example, if the reference note has a subjective section header, we only compared the reference with the subjective section of the generated note (i.e., LED-S, gpt-3.5-SOAP-S, and gpt-4-SOAP-S). Results again reveal a notable performance decrease in out-of-domain models compared to the in-domain baseline (i.e., FLAN-T5). In addition, the performance of objective notes exhibits a relatively milder decline for the SOAP-oriented model.

**Finding 1 (RQ1)**: despite the high performance on the in-domain testing data, the fine-tuning LM-based summarization method suffers from overfitting issues, leading to a notable performance drop on out-of-domain data.

**Finding 2 (RQ2)**: When employing the general-purpose model for doctor-patient summarization, adaptation is essential to ensure the preservation of objective information, which is more prone to being excluded. Experimental results of gpt-SOAP

models indicate that the performance of objective notes can be greatly improved by specifically instructing GPT to generate notes with an objective section.

| Testing data / Model | S | O | AP |
|---|---|---|---|
| ROUGE-1 | | | |
| FLAN-T5 (In-domain) | 0.449 | 0.435 | 0.405 |
| gpt-3.5-general | 0.244 (-46%) | 0.266 (-39%) | 0.180 (-55%) |
| gpt4-general | 0.315 (-30%) | 0.298 (-31%) | 0.214 (-47%) |
| LED-S | 0.231 (-49%) | - | - |
| LED-O | - | 0.259 (-40%) | - |
| LED-AP | - | - | 0.112 (-72%) |
| gpt-3.5-SOAP-S | 0.225 (-50%) | - | - |
| gpt-3.5-SOAP-O | - | 0.357 (-18%) | - |
| gpt-3.5-SOAP-AP | - | - | 0.143 (-65%) |
| gpt-4-SOAP-S | 0.273 (-39%) | - | - |
| gpt-4-SOAP-O | - | 0.347 (-20%) | - |
| gpt-4-SOAP-AP | - | - | 0.184 (-55%) |
| Aggregate Score | | | |
| FLAN-T5 (In-domain) | 0.584 | 0.540 | 0.545 |
| gpt-3.5-general | 0.460 (-21%) | 0.465 (-14%) | 0.423 (-22%) |
| gpt4-general | 0.513 (-12%) | 0.480 (-11%) | 0.449 (-18%) |
| LED-S | 0.401 (-31%) | - | - |
| LED-O | - | 0.411 (-24%) | - |
| LED-AP | - | - | 0.334 (-39%) |
| gpt-3.5-SOAP-S | 0.408 (-30%) | - | - |
| gpt-3.5-SOAP-O | - | 0.482 (-11%) | - |
| gpt-3.5-SOAP-AP | - | - | 0.310 (-43%) |
| gpt-4-SOAP-S | 0.466 (-20%) | - | - |
| gpt-4-SOAP-O | - | 0.492 (-9%) | - |
| gpt-4-SOAP-AP | - | - | 0.406 (-26%) |

Table 5: Model performance on the MTS-Dialog dataset. Testing data S, O, and AP means that the evaluated reference note belongs to the subjective, objective, and assessment and plan categories, respectively. -S, -O, and -AP indicate the generated note in the subjective, objective, and assessment and plan sections, respectively. The values in parentheses indicate the performance change compared with the in-domain FLAN-T5 model (i.e., FLAN-T5 model fine-tuned on MTS-Dialog).

## 6.2 LIWC Analysis of SOAP Note

Experimental results presented in Section 6.1 reveal a notable decline in the performance of the fine-tuning language model-based method when applied to out-of-domain data. In this section, we investigate the characteristics of S, O, and AP samples in two datasets to better understand potential factors for performance degradation.

We computed LIWC features for S, O, and AP notes. Table 6 shows the example words in the selected LIWC categories, and Figure 5 visualizes the selected LIWC features for the ACI-BENCH and MTS-Dialog datasets. First, we find that LIWC shares similar patterns for S, O, and AP notes across the ACI-BENCH and MTS-Dialog datasets. Specifically, these datasets have corrections of 0.93, 0.95, and 0.77 for S, O, and AP notes, respectively. These results indicate that the SOAP notes in the

two datasets are structured in a similar way in terms of word category distribution.

We also observe a similarity in LIWC features between S and AP notes. This alignment is intuitive as S represents subjective information provided by the patient, whereas AP represents the *subjective* assessment and plan from the doctor. One difference between the S and AP notes is that, in S notes, negative emotion is higher than positive emotion, while in the AP notes, negative emotion is lower than positive emotion. This fits a typical scenario where a patient comes to the doctor because of concerns (negative emotion), and then the doctor makes an assessment and plans to address the patient's problem, introducing a more positive emotion.

**Finding 3**: LIWC features have characteristics that resonate with SOAP notes in real-world scenarios.

**Finding 4 (RQ1)**: Because LIWC features exhibit strong correlations for S, O, and AP notes across different datasets, format mismatch (i.e., discrepancies in word distribution) might not be the main cause of the model's performance decline on out-of-domain data.

| LIWC feature | Examples |
|---|---|
| pronoun | I, you, that, it |
| number | one, two, first, once |
| posemo (positive emotion) | good, love, happy, hope |
| negemo (negative emotion) | bad, hate, hurt, tired |
| anx (anxiety) | worry, fear, afraid, nervous |
| anger | hate, mad, angry, frustr* |
| sad | sad, disappoint*, cry |
| hear | heard, listen, sound |
| feel | touch, hold, felt |
| bio | eat, blood, pain |
| body | ache, heart, cough |
| health | medic*, patients, health |
| ingest (food) | food*, drink*, eat, dinner* |
| risk | secur*, protect*, pain, risk* |
| time | when, now, then, day |

Table 6: Selected LIWC features and example words.

## 6.3 Hallucination analysis

We examine the hallucination problem of SOAP-oriented models in scenarios where the input conversation might not include all SOAP information (Figure 6.) First, we compute the length of the generated note. Because Flan-T5 is fine-tuned with the in-domain data, the resulting note lengths are closer to the reference than other models. In contrast, the

Figure 5: LIWC analysis of SOAP notes. Note that this result is calculated using all samples (i.e., training, validation, and testing sets), rather than using only the testing set as experiments on model performance. In addition, for simplicity and visualization purpose, we only show that LIWC categories that have a higher association with healthcare. The correlations between the two data sets are 0.93, 0.95, and 0.77 in S, O, and AP, respectively.



Figure 6: Hallucination medical term ratio, the experiments were conducted on the MTS-Dialog dataset.

out-of-domain LED model generated notes much longer than the reference. In the case of the SOAP-oriented GPT models, each section (S, O, and AP) is shorter than the general model, but the combination of all sections (gpt-S + gpt-O + gpt-AP) is slightly longer than that of the general GPT model.

We then counted the number of unique medical terms that were not mentioned in the input dialogue but *were* generated in the note (i.e., hallucinated medical terms). Finally, we divided the number of hallucinated medical terms by the length

of the generated note to derive the hallucination ratio. We observed that LED has a notable hallucination problem for medical information especially on the objective section. For the SOAP-oriented GPT models, the AP sections (i.e., gpt3.5-SOAP-AP and gpt4-SOAP-AP) exhibit a relatively higher hallucination ratio, suggesting a higher tendency of generating hallucinated medical terms for assessment and plan section. Overall, the GPT-based models manifest a considerably lower hallucination ratio than LED in general.

Even the reference note may contain medical terms not present in the original dialogue. Reasons for this discrepancy are that the reference note is abstractive summarization and may use synonyms as substitutes for the original words. For example, the word "flu" in the conversation was replaced with "influenza" in the reference note. However, we believe it would be better to use exactly the same words as those in the conversation. Although LMs specialized in the medical domain might be aware of the similarity of medical terms, using the same medical terms as the conversation provide better consistency and can avoid confusion. In addition, as new medical terms emerge, the language model might not be updated with the latest information.

## 6.4 Case study of SOAP-oriented model

We conducted a case study to better understand the errors made by the summarization models (Table 7). First, we observed that the LED model tends to overfit its training data. In cases where the input conversation lacks sufficient information, the LED model generates unfaithful and irrelevant informa-

| Dialogue (sample in MTS-dialog dataset) |
|---|
| Doctor: Have you had your flu or Pneumonia vaccination this year? Patient: I got my flu shot this year, but I did not get the one for Pneumonia. |
| Reference (subjective: IMMUNIZATIONS) |
| Immunizations were up to date for influenza, negative for Pneumovax. |

| Model | Section | Generated note | |
|---|---|---|---|
| LED | S | (CHIEF COMPLAINT) **Pneumonia**. (MEDICAL HISTORY) **Patient reports history of pneumonia in the past**. (SOCIAL HISTORY) **patient reports that she has a family history of pneumonia**. (MEDICATIONS) Patient reports **she** received a flu shot this year. (REVIEW OF SYSTEMS) **Constitutional: Reports pneumonia. Respiratory: Denies pneumonia. Gastrointestinal: Reports abdominal pain. Genitourinary: Reportsuria. Musculoskeletal: Reports right lower extremity pain.** | **Hallucination**: Introducing information not present in the original conversation. **Contradictory information**: "report pneumonia" and "denies pneumonia" are contradictory. |
| | O | (PHYSICAL EXAM) Itegumentary. Examination: Normal gait. Examination of Abdomen: No enlarged lymph nodes. Palpation: Normal perfusion. Auscultation of Lungs: No evidence of lymphadenopathy. Cardiovascular Examination; normal perfusion, bilaterally. - Musculosity: Normal pulmonary function. gastrointestinal: Normal vaginal bleeding. | **Hallucination**: no objective information and future plans were mentioned in the conversation. |
| | AP | 1. Pneumonias. - Medical Reasoning: Based on the patient's examination and radiographic findings, we believe this is a case of Pneumona. - Patient Education and Counseling: We discussed the diagnosis and treatment options with the patient. We also discussed the importance of maintaining a high index of suspicion when treating patients with suspected pneumonia. - Additional Testing: We recommended that all patients be referred to a thoracic/GI/Lymphatic surgeon for further evaluation and possible treatment. - medical Treatment: Initiate antiviral treatment. Patient Agreements: The patient understands and agrees with the recommended medical treatment plan. (INSTRUCTIONS) The patient will follow-up as needed. | |

| Model | Section | Generated note | |
|---|---|---|---|
| gpt3.5 | S | N/A | **Incorrect category**: "patient received flu shot" is subjective information (information reported by the patient). |
| | O | The patient has received a flu vaccination this year, but they have not received a pneumonia vaccination | |
| | AP | The patient has not received a pneumonia vaccination. **Recommend the patient to get their pneumonia vaccination.** | |
| gpt4 | S | The patient reported having received their annual flu vaccination, but they have not been vaccinated for pneumonia this year. | **Hallucination**: no future plans were mentioned in the conversation. |
| | O | N/A | |
| | AP | The patient is up to date with their flu vaccination but lacks their pneumonia shot for the current year. **The plan is likely to facilitate the pneumonia vaccination for the patient.** | |

Table 7: Case study example

tion, trying to align the generated note more closely with its training data. In contrast, the gpt3.5-SOAP model has difficulty accurately classifying information. For example, "patient received flu shot" is information obtained from the patient and thus should be categorized as subjective information, but the gpt3.5 model incorrectly categorized this information as an objective fact. Moreover, although the prompts are instructed to "avoid including information that is not explicitly mentioned in the conversation (Table 3)", both gpt3.5-SOAP and gpt4-SOAP models produce hallucination results in the generated AP note. This aligns with our observation in Figure 6 that SOAP-oriented GPT models have a higher hallucination medical term ratio in the AP section. This result suggests that it is important to examine the assessment and plan section, as the model may have a higher tendency to generate hallucinated information in this category.

## 7 Limitations

One limitation of this study is that the SOAP data in the MTS-Dialog dataset is unbalanced, with most references focusing on subjective information. In addition, real-world doctor-patient conversations are complex in size and medical specialties and cannot be fully represented by two datasets. Another

issue lies in the generative model producing varied results in different runs, and the performance of the GPT model is affected by the prompt.

## 8 Conclusion

In this study, we evaluated the SOTA doctor-patient summarization models on out-of-domain data and investigated the challenges of using fine-tuning LM and GPT-based summarization models in real-world applications. For a model with a general configuration, the results indicate a high tendency of omitting objective information in the generated note. This concern can be alleviated by adopting the SOAP-oriented configuration, which orients the model to generate information relevant to all essential categories. Despite achieving the highest performance on in-domain data, the fine-tuned LM with SOAP-oriented configuration exhibits a significant hallucination issue. To generate a note closer to its training data, the model produces hallucinations when none or insufficiently related information is present in the conversation. In contrast, limitations of GPT-based models arise from a tendency to offer their own suggestions for the assessment and plan. We hope our results provide insights for future work toward creating more robust models for real-world settings.

# References

Asma Ben Abacha, Wen-wai Yim, Griffin Adams, Neal Snider, and Meliha Yetisgen-Yildiz. 2023a. Overview of the MEDIQA-Chat 2023 shared tasks on the summarization & generation of doctor-patient conversations. In *Proc. Clinical NLP Workshop 2023*, pages 503–513.

Asma Ben Abacha, Wen-wai Yim, Yadan Fan, and Thomas Lin. 2023b. An empirical study of clinical note generation from doctor-patient encounters. In *Proc. EACL 2023*, pages 2283–2294.

Asma Ben Abacha, Wen-wai Yim, George Michalopoulos, and Thomas Lin. 2023c. An investigation of evaluation metrics for automated medical note generation. In *Proc. ACL Findings 2023*, pages 2575–2588.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

John Giorgi, Augustin Toma, Ronald Xie, Sondra Chen, Kevin An, Grace Zheng, and Bo Wang. 2023. WangLab at MEDIQA-Chat 2023: Clinical note generation from doctor-patient conversations using large language models. In *Proc. Clinical NLP Workshop 2023*, pages 323–334.

Colin Grambow, Longxiang Zhang, and Thomas Schaaf. 2022. In-domain pre-training improves clinical note generation from doctor-patient conversations. In *Proc. NLG4Health 2022*, pages 9–22.

Anirudh Joshi, Namit Katariya, Xavier Amatriain, and Anitha Kannan. 2020. Dr. Summarize: Global summarization of medical dialogue by exploiting local structures. In *Proc. EMNLP 2020 Findings*, pages 3755–3763.

Kundan Krishna, Sopan Khosla, Jeffrey P Bigham, and Zachary C Lipton. 2021. Generating SOAP notes from doctor-patient conversations using modular summarization techniques. In *Proc. ACL 2021*, pages 4958–4972.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Yash Mathur, Sanketh Rangreji, Raghav Kapoor, Medha Palavalli, Amanda Bertsch, and Matthew R Gormley. 2023. SummQA at MEDIQA-chat 2023: In-context learning with GPT-4 for medical summarization. In *Proc. ClinicalNLP Workshop 2023*, pages 490–502.

Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proc. ACL 2020*, pages 7881–7892.

Ashwyn Sharma, David Feldman, and Aneesh Jain. 2023. Team Cadence at MEDIQA-Chat 2023: Generating, augmenting and summarizing clinical dialogue with large language models. In *Proc. Clinical NLP Workshop 2023*, pages 228–235.

Jiyoun Song, Mollie Hobensack, Kathryn H Bowles, Margaret V McDonald, Kenrick Cato, Sarah Collins Rossetti, Sena Chae, Erin Kennedy, Yolanda Barrón, Sridevi Sridharan, et al. 2022. Clinical notes: An untapped opportunity for improving risk prediction for hospitalization and emergency department visit during home health care. *Journal of biomedical informatics*, 128:104039.

Hongjin Su, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A Smith, Luke Zettlemoyer, Tao Yu, et al. 2023. One embedder, any task: Instruction-finetuned text embeddings. In *Proc. ACL 2023 Findings*, pages 1102–1121.

Xiangru Tang, Andrew Tran, Jeffrey Tan, and Mark Gerstein. 2023. GersteinLab at MEDIQA-Chat 2023: Clinical note summarization from doctor-patient conversations through fine-tuning and in-context learning. In *Proc. ClinicalNLP Workshop 2023*, pages 546–554.

Yla R Tausczik and James W Pennebaker. 2010. The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of language and social psychology*, 29(1):24–54.

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems*, 33:5776–5788.

Wen-wai Yim, Yujuan Fu, Asma Ben Abacha, Neal Snider, Thomas Lin, and Meliha Yetisgen. 2023. ACI-BENCH: a novel ambient clinical intelligence dataset for benchmarking automatic visit note generation. *Scientific Data*, 10(1):586.

Longxiang Zhang, Renato Negrinho, Arindam Ghosh, Vasudevan Jagannathan, Hamid Reza Hassanzadeh, Thomas Schaaf, and Matthew R Gormley. 2021. Leveraging pretrained models for automatic summarization of doctor-patient conversations. In *Proc. EMNLP 2021 Findings*, pages 3693–3712.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with BERT. In *Proc. ICLR 2020*.

# Efficient Medical Question Answering with Knowledge-Augmented Question Generation

**Julien Khlaut** [*1], **Corentin Dancette**[*1], **Elodie Ferreres**[*1],
**Alaedine Bennani**[*2], **Paul Hérent**[1], **Pierre Manceron**[1]

[1]Raidium

[2] Service de médecine vasculaire, Hôpital européen Georges Pompidou (HEGP),
AP-HP, Université Paris-Cité, Paris, France

[1]`first.last@raidium.fr`   [2]`alaedine.benani@aphp.fr`

## Abstract

In the expanding field of language model applications, medical knowledge representation remains a significant challenge due to the specialized nature of the domain. Large language models, such as GPT-4 (OpenAI, 2023), obtain reasonable scores on medical question answering tasks, but smaller models are far behind. In this work, we introduce a method to improve the proficiency of a small language model in the medical domain by employing a two-fold approach. We first fine-tune the model on a corpus of medical textbooks. Then, we use GPT-4 to generate questions similar to the downstream task, prompted with textbook knowledge, and use them to fine-tune the model. Additionally, we introduce ECN-QA, a novel medical question answering dataset containing "progressive questions" composed of related sequential questions. We show the benefits of our training strategy on this dataset. The study's findings highlight the potential of small language models in the medical domain when appropriately fine-tuned.

## 1 Introduction

Deep Learning led to a breakthrough in natural language processing, reaching human performances on many tasks like question answering or translation. However, their performances are still subpar in complex domains, such as medicine. This domain presents unique challenges, mainly due to its specialized vocabulary, complex concepts, and fast-changing medical literature. Language-based medical tasks, such as medical question answering, require vast knowledge and reasoning abilities to make correct diagnoses. Traditional language models (LMs), while effective in general language processing, struggle when faced with medical knowledge learning mainly because sufficient data for medical knowledge is not necessarily readily available for training. Moreover, in the context of language models, their number of parameters often plays a pivotal role in performances. Large models, although powerful, come with high computational costs and resource requirements, both for training and inference, making them less accessible and practical for widespread use. On the other hand, small models, which are more economical, face challenges in generalization and adapting to specialized domains like medicine. These models require careful fine-tuning to grasp the depth and breadth of medical knowledge effectively. The diversity of general, non-medical datasets on which LMs are trained poses another challenge. These datasets, encompassing a wide array of topics and styles, do not specifically cater to the medical domain. As a result, small models trained on such datasets might fail to develop the necessary understanding for answering more specialized medical questions.

Therefore, we tackle these issues for medical question answering tasks. First, we design a new dataset, ECN-QA. Existing medical question answering (QA) datasets such as MedQA (Jin et al., 2020) and others (Jin et al., 2019), (Pal et al., 2022) are usually single-question multiple answers, which do not encompass the complexity of making a medical diagnosis, which requires multiple turns of questions. Our dataset is based on the French medical residency examination and contains multiple related questions that require models to remember previous questions and reasoning over multiple steps. We then propose a method to train small to mid-size language models for medical question answering. We leverage a corpus of medical textbooks for pre-training. The pre-training set is enriched with specialized questions generated by large language models prompted with medical data from books. This helps to specialize the model on the target task with a small amount of original data.

---
[*]Equal Contribution

Our code will be made available online.

## 2   Datasets

### 2.1   ECN-QA Dataset

We design ECN-QA, a medical question answering dataset. The questions are collected from FreeCN[1], a website established by French medical students to facilitate ECN (*Examen Classant National*, the national ranking exam before medical residency), with their authorization. This website includes questions from past exams and additional questions ("custom" questions) to simulate exam conditions and aid in studying.

The ECN exams themselves consist of two parts. The first part, known as *Individual Questions* (IQ), features general medicine questions with 5 possible answers. Among these answers, one or multiple may be correct, and candidates must identify the true ones. We display an example in Table 1. The dataset contains 4481 IQ, 721 of which come from the historical data of previous exams. The rest, the "custom" subset, contains 3760 additional IQ-like questions created by the FreeCN team to help students prepare for the exam. The second part is known as *Progressive Questions* (PQ), which features clinical cases. Each PQ consists of an introduction followed by a series of successive questions. Similar to the IQ section, these questions also offer 5 possible answers, with 1 to 5 correct answers. A single PQ can contain numerous successive sub-questions, sometimes more than 20. We have 1050 sub-questions in all PQ. We show an example in Table 5 of Appendix E. We also show a whole progressive question in Appendix E.1. We use the accuracy as our evaluation metric. Each proposition in the question is answered separately and gets a score of 0 or 1. The accuracy is then averaged over the five propositions, i.e., for one question, the possible score can be 0, 0.2, 0.4, 0.6, 0.8, or 1.0. For example, in Table 1, if the model answers a, b, c, e as wrong and d as right, it would have one error since c is right. The accuracy would, therefore, be 0.8. If the model answers a and e as wrong and b, c, and d as right, it would also have an accuracy of 0.8.

All the original data is in French, but all models are pre-trained using mostly English data. Therefore, we translate all the questions and answers into

---

> **Question**: A woman of Martinican origin has just given birth. The child's father is also of Martinican origin. The child has a cleft lip and palate. With regard to regulatory newborn screening of this child, what is the exact proposal(s)? **Propositions**:
> (a) Phenylketonuria is the only disease of amino acid and organic acid metabolism currently being screened for newborn in France
> (b) General screening test can detect hypothyroidism of pituitary origin
> (c) **This couple can refuse the screening after information**
> (d) **Completion before 48 h of life decreases the sensitivity and/or specificity of the screening test**
> (e) Targeted screening for sickle cell disease is not indicated in this child

Table 1: Example of Individual Question (IQ) in the ECN, translated to English. Correct answers are in bold.

English using the Azure AI Translation API[2].

### 2.2   Medical Textbooks

Additionally, we use classical French medical textbooks designed for medical students, containing comprehensive medical knowledge and established protocols for managing various medical conditions. We detail in Section F how we extract sections from medical textbooks in PDF format.

In total, we worked with 17,509 PDF files. We grouped text in sections rather than pages, recognizing that a single topic might span multiple pages and should not be truncated. The sections are defined by the book titles and correspond to chapters or important parts. This approach resulted in a total of 234,495 sections. The full dataset is composed of 174,242,531 tokens (with the GPT-3 tokenizer). We detail how we extract sections from PDF files in Appendix F.

We use them for pre-training, and to generate additional questions, as explained in Section 3.

## 3   Method

We detail our training strategy in this section. The strategy is depicted in Figure 1. We detail related works in the Appendix A.

### 3.1   Baseline Model

For our baseline, we use the BioMedLM model (Bolton et al., 2022). This 2.7-billion-parameter model is built upon the GPT-2 architecture (Radford et al., 2019) and has been trained on a substantial corpus of medical and biological

---

[1] https://www.freecn.io

[2] https://learn.microsoft.com/en-us/azure/ai-services/translator/

data. BioMedLM's specialized biomedical tokenizer sets it apart, enhancing its comprehension of specialized terminology. BioMedLM's training data contains all PubMed abstracts and full documents from The Pile (Gao et al., 2020), ensuring a rich knowledge base. Notably, BioMedLM reported state-of-the-art scores on the MedQA (Jin et al., 2020) dataset.

However, this model does not possess the scale needed to achieve impressive zero-shot generalization on new tasks, and medical question answering datasets are limited in scale. Therefore, we aim to train it on specific high-quality data that resembles our benchmark. As our training dataset is small (4967 questions), we propose a method to augment it with question generation using a large language model prompted by some medical knowledge extracted from textbooks.



Figure 1: Our training strategy. Starting from an existing language model such as BioMedLM, we continue the pre-training on our corpus of medical textbooks. Then, we use GPT-4, prompted with knowledge from the textbooks, to generate clinical cases that are used to fine-tune the model.

## 3.2 Questions Generation

Our objective is to create cases that closely resemble genuine ECN cases, as this offers the most effective training for the model. The format we desire for these cases closely resembles that of the progressive questions: an introduction, a list of questions and their possible answers, with a label (true or false) for each answer.

To create our clinical cases, we concatenate several prompts using different approaches. The design of each prompt begins with adopting the prompt used by FreeCN, which primarily comprises an introduction to the task. We refer to this as the pre-prompt. Next, we compile a list of all the specific details we want the case to encompass. This list is informed by the insights of medical experts and the manner in which they typically structure questions

for the ECN. We refer to this list as the "constitution." When we initially applied this approach, we encountered somewhat disappointing results. The clinical cases exhibited two major shortcomings. First, they often had very similar subjects, causing the model to struggle to generate diverse cases. Additionally, the main issue was that the questions posed were consistently identical, revolving around topics such as "What is the diagnosis?" or "What tests would you conduct for confirmation?" and "How would you manage the patient?". To solve this issue and to introduce diversity in disease scenarios, we also supplied it with specific knowledge that could be utilized to construct these cases. This section was named the "knowledge part," and it drew upon information extracted from sections of medical books. An example can be seen in Annex 4.

We introduce an additional "justification" field. This component explains why a particular answer to a question is deemed suitable or not.

We build our pipeline using the OpenAI API, employing the GPT-4 model (Achiam et al., 2023) to generate clinical cases. We use the GPT-4 function calling JSON mode, which allows us to specify the output structure.

Following this approach, we generated a dataset with GPT-4, containing about 10,237,240 tokens. In some instances, the dataset underwent meticulous filtering procedures to rectify issues such as missing or alternative fields. This approach is inspired by phi (Gunasekar et al., 2023; Li et al., 2023) and Orca (Mitra et al., 2023; Mukherjee et al., 2023).

We gathered feedback and validation from the FreeCN team, composed of medical doctor students, for assistance and insights to ensure the quality of the questions. The prompt given to GPT-4 is displayed in Appendix B, and an example of a generated progressive question in Appendix G.

## 3.3 Pre-training

The initial phase involves pre-training the model on a dataset, partly composed of medical books and the additional generated questions. We start from BioMedLM's weights and use a next-token prediction loss to pre-train for three epochs. After training on the books, the model is further trained on the generated cases. The 160,889 generated questions are composed of 10,237,240 tokens. The training

is performed on one case at a time and the final loss is computed only on the model's answer and justification. Since the context length of BioMedLM is 2048, we truncate more prolonged cases. The training parameters are detailed in Appendix C.

## 3.4 Fine-tuning

Following the pre-training phase, the next step is fine-tuning the model on the ECN-QA dataset. For fine-tuning, the dataset is split into 90 % for training and 10% for testing set. There are multiple ways of getting the model to output an answer, for example, generating tokens with a specific format. However, since generating consistent word-by-word answers proved challenging for the model, often resulting in gibberish rather than accurate responses, we opted for a more traditional approach during fine-tuning. Similarly to previous work (Bolton et al., 2022), a classification head was added to the model. It operates at the proposition level: the model takes as input the question and a single proposition among the five. It then has to predict if the proposition is right or wrong, as a binary classification task. One possible approach to this binary classification involves predicting a single scalar value for each answer, training it with binary cross-entropy, and selecting a threshold value for inference. Another approach consists of adding the words "true" or "false" to the end of the sentence, feeding both sentences to the model, and selecting the answer with the highest score. Empirically, the second approach provided the best results. This modification allowed us to obtain more reliable responses from the model during evaluation.

## 4 Results

### 4.1 Evaluation of GPT models

We first evaluate the GPT models on our dataset to obtain baseline scores. For both GPT-3.5 and GPT-4 models, the 2023-12-01 version of the API is used (available on Azure).

We encountered occasional issues during evaluation, where specific prompts may have been blocked, possibly due to sensitive subjects like pediatric medicine. In such cases, we considered the model's response incorrect. The prompts were designed to be straightforward, typically asking the model to provide a true or false answer. Moreover, questions were asked in English using the translated dataset.

| Model | Accuracy |
|---|---|
| GPT-3.5 | 69.36 |
| GPT-4 | 79.04 |
| GPT-4-32k | 78.97 |
| GPT-4-32k 5 few shot | 81.42 |

Table 2: Results on the all evaluation dataset

The results are presented in Table 2. GPT-4's performances on our dataset are similar to those on MedQA and USMLE, reaching zero-shot performances of around 74% (Nori et al., 2023a). Overall, GPT-4 and its 32k-context variant is the strongest model. Additionally, we confirm (Nori et al., 2023b)'s findings that adding some questions in the prompt (*few shot*) increases the accuracy, in our case, by around 2.5 points.

### 4.2 Main Results

| Model | Accuracy |
|---|---|
| BioMedLM | 67.74 |
| BioMedLM + Books | 69.65 |
| BioMedLM + MQG | 68.62 |
| BioMedLM + Books + MQG | **70.56** |

Table 3: Final results for BioMedLM with various parameters. MQG stands for Medical Question Generation. The model is trained on books for three epochs and on MQG for two epochs. All models are then fine-tuned on ECN-QA.

The results of our experiments are shown in Table 3. We report the result of the original BioMedLM, as well as models pre-trained on the collection of books (*BioMedLM + Books*), pre-trained on the questions (*MQG* for Medical Question Generation) and our complete method (*BioMedLM + Books + MQG*). All models are fine-tuned on ECN-QA.

Including books as part of our training data improves the accuracy by approximately 2 points and the MQG method alone by 1 point. The best accuracy is achieved by combining the pre-training using books with the question-generation method. Overall, we significantly improve the baseline with our full method, getting +3 points in accuracy. We also surpass the GPT-3.5 model, as shown in Table 2.

In Figure 2, we display the number of questions for each score for our full method and GPT-4. We observe that our model still lags behind GPT-4.

Figure 2: Accuracy distribution by question (number of correct propositions divided by number of total propositions) on the FreeCN dataset of GPT-4 and BioMedLM + Books + MQG

Since the model answers all propositions independently and has no knowledge of its answer to other propositions, the model can contradict itself, which makes it harder to obtain a score of 1 (i.e. having the right answer to all propositions). More detailed statistics per subject are available in Figure 3. Our method appears less effective on subjects it has not been trained on, such as pediatrics.

## 5 Conclusion

We introduced ECN-QA, a novel dataset for medical question answering that contains a novel type of exercise: progressive questions. We proposed a training strategy based on prompted question generation that improves results over our baseline model, enabling the model to surpass GPT-3.5 accuracy with a much lower parameter count.

Potential avenues for improving efficient medical question answering include increasing the size of the pre-training dataset and the number of generated questions and investigating retrieval-based answering (open-book exam). A model with significant capabilities in medical answering can aid in making informed decisions, especially in time-sensitive situations where rapid response is crucial. Such a model can offer up-to-date information, suggest potential diagnoses, and recommend treatment options based on the latest research and clinical guidelines.

## 6 Ethical Concerns

The model was trained on questions designed for students' examination, not for a real-world clinical setting. The generalization of this model to actual clinical settings is unknown. Indeed the model has potential biases and limitations in handling sensitive and complex medical cases and should not be used as so on real-world patients.

## 7 Acknowledgement

14

# References

OpenAI Josh Achiam, Steven Adler, Sandhini Agar-
wal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
man, Diogo Almeida, Janko Altenschmidt, Sam Alt-
man, Shyamal Anadkat, Red Avila, Igor Babuschkin,
Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-
ing Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake
Berdine, Gabriel Bernadett-Shapiro, Christopher Berner,
Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-
Luisa Brakman, Greg Brockman, Tim Brooks, Miles
Brundage, Kevin Button, Trevor Cai, Rosie Campbell,
Andrew Cann, Brittany Carey, Chelsea Carlson, Rory
Carmichael, Brooke Chan, Che Chang, Fotis Chantzis,
Derek Chen, Sully Chen, Ruby Chen, Jason Chen,
Mark Chen, Benjamin Chess, Chester Cho, Casey Chu,
Hyung Won Chung, Dave Cummings, Jeremiah Cur-
rier, Yunxing Dai, Cory Decareaux, Thomas Degry,
Noah Deutsch, Damien Deville, Arka Dhar, David
Dohan, Steve Dowling, Sheila Dunning, Adrien Ecof-
fet, Atty Eleti, Tyna Eloundou, David Farhi, Liam
Fedus, Niko Felix, Sim'on Posada Fishman, Juston
Forte, Isabella Fulford, Leo Gao, Elie Georges, Chris-
tian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh,
Raphael Gontijo-Lopes, Jonathan Gordon, Morgan Graf-
stein, Scott Gray, Ryan Greene, Joshua Gross, Shixi-
ang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han,
Jeff Harris, Yuchen He, Mike Heaton, Johannes Hei-
decke, Chris Hesse, Alan Hickey, Wade Hickey, Peter
Hoeschele, Brandon Houghton, Kenny Hsu, Shengli
Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,
Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin,
Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun,
Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kan-
itscheider, Nitish Shirish Keskar, Tabarak Khan, Logan
Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik
Kim, Hendrik Kirchner, Jamie Ryan Kiros, Matthew
Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew
Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen
Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy
Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li,
Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin,
Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Ade-
ola Makanju, Kim Malfacini, Sam Manning, Todor
Markov, Yaniv Markovski, Bianca Martin, Katie Mayer,
Andrew Mayne, Bob McGrew, Scott Mayer McKin-
ney, Christine McLeavey, Paul McMillan, Jake Mc-
Neil, David Medina, Aalok Mehta, Jacob Menick, Luke
Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
Monaco, Evan Morikawa, Daniel P. Mossing, Tong
Mu, Mira Murati, Oleg Murk, David M'ely, Ashvin
Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Nee-
lakantan, Richard Ngo, Hyeonwoo Noh, Ouyang Long,
Cullen O'Keefe, Jakub W. Pachocki, Alex Paino, Joe
Palermo, Ashley Pantuliano, Giambattista Parascan-
dolo, Joel Parish, Emy Parparita, Alexandre Passos,
Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe
de Avila Belbute Peres, Michael Petrov, Henrique Pondé
de Oliveira Pinto, Michael Pokorny, Michelle Pokrass,
Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris
Power, Elizabeth Proehl, Raul Puri, Alec Radford,
Jack Rae, Aditya Ramesh, Cameron Raymond, Francis
Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri
Roussez, Nick Ryder, Mario D. Saltarelli, Ted Sanders,
Shibani Santurkar, Girish Sastry, Heather Schmidt,
David Schnurr, John Schulman, Daniel Selsam, Kyla
Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker,
Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie
Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Ben-
jamin D. Sokolowsky, Yang Song, Natalie Staudacher,
Felipe Petroski Such, Natalie Summers, Ilya Sutskever,
Jie Tang, Nikolas A. Tezak, Madeleine Thompson, Phil
Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston
Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cer'on
Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea
Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang,
Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann,
Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian
Weng, Matt Wiethoff, Dave Willner, Clemens Win-
ter, Samuel Wolrich, Hannah Wong, Lauren Workman,
Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu,
Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,
Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia
Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk,
and Barret Zoph. 2023. Gpt-4 technical report.

Elliot Bolton, David Hall, Michihiro Yasunaga, Tony
Lee, Chris Manning, and Percy Liang. 2022. Biomedlm.

Dana Brin, Vera Sorin, Akhil Vaid, Ali Soroush, Ben-
jamin S. Glicksberg, Alexander W. Charney, Girish Nad-
karni, and Eyal Klang. 2023. Comparing chatgpt and
gpt-4 performance in usmle soft skill assessments. *Sci-
entific Reports*, 13.

Zeming Chen, Alejandro Hern'andez Cano, Angelika
Romanou, Antoine Bonnet, Kyle Matoba, Francesco
Salvi, Matteo Pagliardini, Simin Fan, Andreas Kopf,
Amirkeivan Mohtashami, Alexandre Sallinen, Alireza
Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz
Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne
Hartley, Martin Jaggi, and Antoine Bosselut. 2023.
Meditron-70b: Scaling medical pretraining for large
language models. *ArXiv*, abs/2311.16079.

Leo Gao, Stella Biderman, Sid Black, Laurence Gold-
ing, Travis Hoppe, Charles Foster, Jason Phang, Ho-
race He, Anish Thite, Noa Nabeshima, Shawn Presser,
and Connor Leahy. 2020. The Pile: An 800gb dataset
of diverse text for language modeling. *arXiv preprint
arXiv:2101.00027*.

Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio
César Teodoro Mendes, Allie Del Giorno, Sivakanth
Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo
de Rosa, Olli Saarikivi, et al. 2023. Textbooks are
all you need. *arXiv preprint arXiv:2306.11644*.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,
Hanyi Fang, and Peter Szolovits. 2020. What disease
does this patient have? a large-scale open domain ques-
tion answering dataset from medical exams. *ArXiv*,
abs/2009.13081.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W.
Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for
biomedical research question answering. In *Conference
on Empirical Methods in Natural Language Processing*.

Yuan-Fang Li, Sébastien Bubeck, Ronen Eldan, Alli-
son Del Giorno, Suriya Gunasekar, and Yin Tat Lee.

2023. Textbooks are all you need ii: phi-1.5 technical report. *ArXiv*, abs/2309.05463.

Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Codas, Clarisse Simoes, Sahaj Agarwal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, et al. 2023. Orca 2: Teaching small language models how to reason. *arXiv preprint arXiv:2311.11045*.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023a. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023b. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*.

OpenAI. 2023. GPT-4 Technical Report.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa : A large-scale multi-subject multi-choice dataset for medical domain question answering. In *ACM Conference on Health, Inference, and Learning*.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023a. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023b. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*.

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan,

Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288.

Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Towards building open-source language models for medicine. *arXiv preprint arXiv:2305.10415*, 6.

## A   Related Work

### A.1   Medical QA Datasets

Various datasets have been developed in medical question answering (QA). Among these, the MedQA dataset (Jin et al., 2020) stands out for its comprehensive coverage of multiple-choice questions derived from professional medical board exams. This dataset is particularly significant because it encompasses many questions, totaling 12,723 items. It aims to evaluate the depth of medical knowledge encoded in AI models.

Another dataset is PubMedQA Dataset for Biomedical Questions, (Jin et al., 2019). This dataset uniquely focuses on questions generated from article titles and abstracts within the biomedical literature, excluding conclusions, and provides answers in a format conducive to yes/no/maybe evaluations.

Further expanding the landscape, the MedMCQA dataset (Pal et al., 2022) is a large-scale, multi-subject repository of medical multiple-choice questions. This dataset has a large scope and relevance, covering many medical subjects.

### A.2   Medical QA Models

Several strategies aim to construct a good model with high accuracy and reliability of responses on those medical tests. One method involves leveraging large language models (LLM) such as GPT-4. Through prompt engineering, (Brin et al., 2023) or (Nori et al., 2023a) have demonstrated excellent results on MedQA.

Further exploration into the efficacy of large-scale models has been conducted, with (Singhal et al., 2023a) and (Singhal et al., 2023b). These studies have assessed the performance of such models on MedQA and across a diverse array of medical datasets.

Moreover, the landscape of medical QA has been enriched by initiatives to fine-tune pre-existing LLMs. For instance, adaptations of Llama 2 (Touvron et al., 2023) have been proposed (Wu et al., 2023; Chen et al., 2023). These efforts signify a targeted move towards refining the capabilities of LLMs to meet the demands of the medical domain, illustrating a focus on customizing general models for specialized tasks.

In the context of smaller-scale models, (Bolton et al., 2022) has been recognized for its superior performance. This model stands out as a testament to the effectiveness of more compact models in handling medical QA tasks, offering an alternative to the larger, more complex systems.

## B   Question Generation

Table 4 shows the prompt we used to generate questions with GPT-4. The prompt is appended with a section coming from a medical textbook.

---

You are a French professor of medicine. You seek to test the level of medicine of your students. Your task is to generate 1 to 2 different clinical cases requiring the highest medical understanding. Each clinical case consists of an Introduction and 4-10 multiple-choice questions. They must be formatted as follows: Introduction, Propositions. Propositions contain several proposals with a justification and a field to know if they are correct. The clinical case needs to be very very hard and accurate. The level of difficulty is 10 out of 10. It should be very hard even for the best students. And you should have a very detailed justification. The case should be long with detailed questions and detailed justification.
The criteria to be met are:
1. The introduction is common to all questions.
2. There must be 4-10 different questions.
3. A question can have 5-10 possible choices.
4. One or more proposals may be fair.
5. Justification must be specific, justified and sourced. It is very important to have a very good and long justification. It should be at least 3 lines long.
6. Uses the highest medical level possible.
7. Questions must be diversified to a minimum of 4. They must deal with the patient's disease but also with the examinations to be carried out, the follow-up and the possible developments of the case. They will make the case both nuanced and complex.
8. The case must be precise or even quantitative. It is a question of providing as much information as possible, and the solution to the questions may be found in detail.
9. Cases must be pedagogical and the questions must be linked to build a complete reasoning.
10. Responses should be directed to prioritize severe and frequent cases.
11. The student's expected behaviour is above all to avoid medical misconduct.
12. The student's method must be a probabilistic approach.
13. A language model must be able to answer questions. For example, do not ask the wizard to create images or audio.
14. The case must be written in English.
15. All fields must be completed.
16. The MA for the drug and the recommendations of the HAS and ANSM must be respected. In the absence of recommendations from HAS and ANSM, the current practices recommended by French speciality colleges and learned societies will be applied.
###
To do that you can use the following information: [Extract of a medical book]

---

Table 4: Prompt used to generate progressive questions with GPT-4

Figure 3: Accuracy per subject of BioMedLM and GPT-4

## C   Implementation and training details

For training, we use a node with 4 NVIDIA V100 gpus. The model is pre-trained on books for three epochs on generated cases for two epochs, and we fine-tune the final model for 20 epochs. We use a learning rate of 1.e-4 for pre-training and 2.e-6 for fine-tuning.

## D   Additional Results

We show in Figure 3 the accuracies for each topic of the ECN exam for GPT-4 and our model. Our model is close to GPT-4 for most subjects and performs worse on subjects that GPT-4 often refuses to generate questions, like pediatrics. These questions are from the test set but only come from the additional questions provided by the FreeCN team.

## E   Example of Cases

In Table 5, we display the first question of a progressive question from the ECN-QA dataset, with the propositions of answer.

**Introduction**: A 67-year-old man consults for right calf pain occurring after a walk that the patient estimates to be 350 meters away. He is a retired and sedentary taxi driver. This patient has been smoking a pack of cigarettes a day since the age of 30. You follow it for high blood pressure discovered by a systematic and balanced examination by perindopril. Blood sugar is normal as well as lipid balance.
**Question**: What is your main diagnostic hypothesis ?
**Propositions:**
- **Obliterating arterial disease of the lower limbs**
- Narrow lumbar canal
- Lombosciatica
- Hypokalemia
- Deep vein thrombosis

*For the example the following questions are:*
- You suspect arterial disease obliterating the lower limbs. Which of the following semiological elements will guide the diagnosis towards this hypothesis?
- The interrogation confirms the appearance of a pain when walking with a cramp localized in the right calf. The pain manifests itself early when the patient climbs a slope, thus supporting your diagnostic hypothesis of arterial obliterating disease of the lower limbs. [...] On the data of this clinical examination, which is(are) the arterial atheromatous lesion(s) that you should suspect?

Table 5: Example of PQ in the ECN-QA dataset. This particular PQ consists of 16 questions, and both this PQ and the previous IQ section are derived from the 2020 ECN. Correct answers are in bold.

## E.1 Full Progressive Question

Here, we display a full progressive question with all possible answers. Correct answers are in bold font.

---

**Full example of a progressive question**

**Introduction:** A 54-year-old man, a long-term smoker who has been hypertensive for 12 years (calcium channel blocker treatment), consults his attending physician for an isolated episode of total gross hematuria, without a clot. His other history has been an appendectomy in childhood. The blood count is as follows: Hb 10.4 g/dL (MCV 78 μm3), GB 8 G/L, blisters 247 G/L. Creatinine is 110 μmol/L (estimated glomerular filtration rate of 65 ml/min/1.73 m2). A renal ultrasound showed a hyperechoic mass of 7 cm on the right kidney.

**Questions:** What are the elements (present or to be sought at the interrogation and clinical examination) that can evoke a malignant tumor of the kidney? (one or more correct answers)

- **Smoking**
- **Chronic high blood pressure**
- Long-term calcium channel blocker treatment
- A family history of multiple endocrine neoplasia
- **Low back pain**

Which exam(s) are you asking for as a first line?

- **Urinary cytology with pathological examination**
- **Cytobacteriological examination of urine**
- Serum erythropoietin assay
- **Abdominopelvic CT scan with and without contrast injection**
- Ultrasound-guided puncture of the mass

On the cut shown below, what are the True propositions? (one or more correct answers)

- **This is an abdominal CT scan with injection**
- This is a coronal cup
- Structure number 1 is the inferior vena cava
- **The cut passes through the third duodenum**
- The number 2 corresponds to the inferior mesenteric artery

What are the real propositions? (one or more correct answers)

- **The patient must receive red blood cells**
- The patient must receive platelet pellets
- In case of transfusion of red blood cells, you would prescribe O-negative pellets
- A search result for irregular agglutinins less than 48 h old must be available
- Since 2003, there has been no risk of transmission of infectious pathogens through red blood cell transfusion

What is the real proposal(s)?

- **This is acute renal failure**
- The glomerular filtration rate must be recalculated
- An obstacle on the contralateral kidney is likely
- **It may be functional renal failure**
- **An ionogram should be prescribed on a urine sample**

What are the exact proposals? (one or more correct answers)

- He has moderate chronic renal failure
- **His antihypertensive treatment must include an inhibitor of the renin-angiotensin system**

---

- The LDL cholesterol target to be achieved is 1.3 g/L
- He must follow a diet containing no more than 1.5 g/kg of protein weight
- It is necessary to advocate a diet low in fast sugars

What risk(s) does he run?

- **Gradual decrease in diuresis**
- **Increased cardiovascular risk**
- **Hyperphosphoremia**
- **Erectile dysfunction**
- **Contralateral kidney cancer**

What is the True answer(s)?

- The ALD file is completed by the patient and validated by the medical specialist
- **The attending physician must specify in the request the protocol of care envisaged including treatments, examinations and consultations**
- **The medical officer of the Health Insurance must validate the care protocol**
- In case of coverage in ALD, remains the responsibility of the patient only the co-payment
- The third-party payer is the part of the care paid by the insured whether or not he is registered in ALD

What is your interpretation of the electrocardiogram below?

- **Sinus rhythm**
- Sino-auricular block
- T-waves suggestive of hyperkalemia
- Expanded QRS Complexes
- **Left ventricular hypertrophy**

To reduce edematous syndrome, what do you recommend at this stage? (one or more correct answers)

- A low-salt diet (less than 6 g/d)
- Water restriction
- **A loop diuretic (furosemide)**
- A thiazide diuretic (hydrochlorothiazide)
- Blood ultrafiltration (start of hemodialysis)

What are the possible cause(s) in the context of the new biological abnormality observed?

- **Excessive calcium intake**
- Taking furosemide
- **Chronic renal failure**
- Secondary hyperparathyroidism
- **Bone metastases from kidney cancer**

What additional examination(s) do you recommend to explore this biological anomaly?

- **Ionized serum calcium**
- Test de PAK
- **PTH assay**
- PTHrp assay
- **Bone scintigraphy**

Which proposals are correct? (one or more correct answers)

- Metastatic cancer is a contraindication to dialysis
- Haemodialysis confers survival advantage over peritoneal dialysis
- The preparation of an arteriovenous fistula (AVF) is contraindicated given the prognosis
- **A tunneled central venous catheter may be placed to initiate hemodialysis**
- A transplant from a cadaveric donor must be discussed

In general, regarding living donors, what are the real proposals? (one or more correct answers)

- Only people with a genetic link to the recipient can be donors
- **Transplantation can be done in incompatible ABO condition**
- Rhesus compatibility must be respected

- HLA incompatibility between donor and recipient is a formal contraindication
- The donor is remunerated on a basis proportional to the recipient's waiting time

## F   Extraction of text from PDF files

We use the Azure AI Document Intelligence API[3] to extract text sections from PDF files. The API returns paragraphs of texts and titles, sorted in reading order, along with tables and figures. We remove all tables and figures and implement text filtering algorithms to remove useless paragraphs, like headers, footers, or paragraphs that are just a few characters long.

We then regroup paragraphs in sections based on the titles and divide them into subsections, or regroup them, to have sections of similar lengths, between 500 and 1000 words.

## G   Generated Progressive Question

Below, we present an example of a progressive question generated by GPT-4. Correct answers are in bold font.

---

**Introduction:** A 45-year-old female is being admitted to the emergency department. She is complaining of a severe and sudden headache unlike any she has ever experienced before. The headache was followed by episodes of vomiting and photophobia. Her Glasgow Coma Scale (GCS) score on admission is 14 and her physical examination is unremarkable. Computed Tomography (CT) of the brain reveals subarachnoid hemorrhage (SAH).

**Questions** What is the most likely diagnosis in this patient?

- **Ruptured cerebral aneurysm** (The presence of SAH in the setting of a severe, sudden-onset headache (often described as a 'thunderclap headache'), especially when associated with vomiting and photophobia, is highly suggestive of a ruptured cerebral aneurysm. Other common symptoms of SAH due to a ruptured aneurysm include neck stiffness, altered mental status and cranial nerve palsies.)
- Migraine (While migraines can cause severe headaches, they do not typically lead to SAH. Moreover, the acute onset described in the vignette is uncharacteristic of migraines which usually have a more gradual onset.)
- Tension headache (A tension headache often presents as a dull, constant and bilateral headache. SAH or sudden severe headache is not typical of a tension headache.)
- Cluster headache (Cluster headaches are usually unilateral, severe, and around the eye, and are associated with autonomic symptoms such as tearing, nasal congestion and ptosis. They do not typically

cause vomiting or lead to SAH.)
- Brain tumor (While brain tumors can cause headaches, they do not typically present acutely. Moreover, they generally result in neurological symptoms and signs that localize to the area of the brain where the tumor is located rather than leading to SAH.)

What is the next step in the management of this patient?

- **Cerebral angiography** (Cerebral angiography is the gold standard for diagnosing cerebral aneurysms. It can demonstrate the size, location and morphology of the aneurysm and help plan intervention. Cerebral angiography is also able to diagnose other causes of SAH such as arteriovenous malformation.)
- Magnetic resonance imaging (MRI) (While MRI may be used to evaluate for other causes of SAH such as arteriovenous malformations (AVM) if the angiogram is inconclusive cerebral angiography is the preferred primary imaging modality in this clinical scenario.)
- Lumbar puncture (Lumbar puncture may be performed when the clinical suspicion for SAH is high but initial CT scan is negative. Given the positive CT findings in this case lumbar puncture would not be the best next steps.)
- Blood cultures (Blood cultures are indicated in patient with suspected infective endocarditis that can lead to mycotic aneurysms. In this case there is no clinical suspicion of infective endocarditis.)
- Repeat CT scan (A repeat CT scan may be used in the surveillance of SAH treatment and complications but it is not the best choice for the acute management of suspected ruptured cerebral aneurysm which requires immediate intervention.)

If this patient's cerebral angiography reveals a vasospasm of the cerebral arteries, what is the most appropriate step that should be taken?

- Platelet antiaggregation (Platelet antiaggregation is not the standard management for cerebral vasospasm. Studies showed that antiplatelet medication would lead to a higher rate of rebleeding in cerebral aneurysm.)
- Anticoagulation (Anticoagulation is not the standard management for cerebral vasospasm as it could increase the risk of bleeding.)
- **Calcium channel blockers** (Nimodipine, a calcium channel blocker, is the only medication proven to improve neurologic outcomes in the case of SAH. It doesn't prevent or treat the vasospasm itself but reduces the severity of neurologic deficits and improves outcome.)
- Statins (Although statins have been shown to have some vasoprotective effects they are not the first-line treatment for cerebral vasospasm.)
- Analgesics (While managing pain is an important component of care in this patient population analgesics themselves do not treat vasospasm).

---

# Gemini Goes to Med School:
# Exploring the Capabilities of Multimodal Large Language Models
# on Medical Challenge Problems & Hallucinations

**Ankit Pal,    Malaikannan Sankarasubbu**

Saama AI Research, Chennai, India

{ankit.pal, malaikannan.sankarasubbu}@saama.com

## Abstract

Large language models have the potential to be valuable in the healthcare industry, but it's crucial to verify their safety and effectiveness through rigorous evaluation. In our study, we evaluated LLMs, including Google's Gemini, across various medical tasks. Despite Gemini's capabilities, it underperformed compared to leading models like MedPaLM 2 and GPT-4, particularly in medical visual question answering (VQA), with a notable accuracy gap (Gemini at 61.45% vs. GPT-4V at 88%). Our analysis revealed that Gemini is highly susceptible to hallucinations, overconfidence, and knowledge gaps, which indicate risks if deployed uncritically. We also performed a detailed analysis by medical subject and test type, providing actionable feedback for developers and clinicians. To mitigate risks, we implemented effective prompting strategies, improving performance, and contributed to the field by releasing a Python module for medical LLM evaluation and establishing a leaderboard on Hugging Face for ongoing research and development. Python module can be found at github.com/promptslab/RosettaEval

## A.1   Introduction

Large language models (LLMs) that can understand and generate text that is similar to human language have shown remarkable progress across domains such as language (Brown, 2020) and code (Baptiste Rozière, 2024). Models like GPT-3 (Brown, 2020) and PaLM (Aakanksha Chowdhery, 2022) have been pre-trained on massive text datasets and demonstrate an ability to recognize linguistic patterns. The rapid innovations in artificial intelligence, driven by the continual development of more powerful LLMs, promise to accelerate discovery and enhance research in specialized domains. Capabilities have improved systematically alongside increases in model size, data, and computation. Many of these advanced models leverage



Figure A.1: The MultiMedQA score of the Med-PaLM 2, GPT-4 and Gemini Pro, where the detailed performance of MultiMedQA in Section A.4.1

the transformer architecture (Vaswani et al., 2017), which is well-suited for linguistic applications and are further enhanced through self-supervised learning techniques for textual data.

The application of LLMs in medicine is not only innovative but essential. These models can parse vast amounts of medical literature, synthesize information, and offer insights, which could be a breakthrough in an industry where knowledge evolves rapidly. Researchers have begun assessing how LLMs may assist medicine by augmenting human capabilities (Karan Singhal, 2023; Singhal et al., 2022). The deployment of LLMs within the medical domain presents both promising opportunities and significant challenges. Critical open questions persist - can LLMs demonstrate expert-level medical comprehension? Do they make potentially unsafe errors beyond their competence limits? Assessing these capabilities and limitations will be critical as we explore responsible ways to harness

the power of language models to advance medicine.

Recent research into benchmarks has revealed how LLMs absorb clinical knowledge (Liévin et al., 2023), indicating potential ways for improving medical practices. Google's Gemini model (Gemini Team, 2023) is at the forefront of multimodal language modelling, designed to comprehend and generate content from text, images, audio, and video inputs. With its architecture promising deep comprehension and contextual awareness, Gemini seems well-suited to navigating the complexities of medical data. This study seeks to analyze Gemini's capabilities by comparing it with other models in order to identify its strengths and limitations within the medical domain through investigation of several key questions:

- *How accurately can Gemini solve complex medical reasoning problems in different modalities, including textual and visual information processing?*

- *Does Gemini hallucinate and produce false medical information without appropriate safeguards? When faced with difficult questions, does Gemini guess or admit the limits of its knowledge?*

Our research focuses on evaluating Google's Gemini within the medical domain. Using three benchmarks: MultiMedQA, Med-HALT (Pal et al., 2023), and Medical Visual Question Answering (Jin et al., 2024). We rigorously assess Gemini's proficiency in medical reasoning, susceptibility to hallucination, and comparative performance against open-source and commercial models. The addition of the Medical VQA task aims to evaluate Gemini's capacity to interpret medical imagery and comprehend complex visual questions, representing a critical aspect of clinical diagnostics and patient care.

Our findings reveal that while Gemini demonstrates a robust understanding across various medical subjects, it also exhibits certain limitations, particularly in areas requiring intricate reasoning or specialized knowledge. Through extensive testing across diverse medical datasets, we highlight Gemini's strengths in synthesizing medical literature and pinpoint areas where it falls short. For example, in handling complex diagnostic questions and avoiding misinformation.

In brief, the contributions of this study are as follows

- **First Rigorous Multi-Modal Evaluation of Gemini's Medical Competencies:** We provide a detailed assessment of Google Gemini's performance across the VQA & MultiMedQA benchmark. We employ various advanced prompting techniques such as direct few-shot, chain-of-thought, self-consistency, and ensemble refinement to evaluate Gemini's understanding and reasoning in the medical domain.

- **Probing Safety & Hallucination Risks through Med-HALT:** Our research presents an in-depth evaluation of Gemini on the Med-HALT benchmark to systematically assess hallucination tendencies in medical LLMs. By exploring both reasoning-based and memory-based hallucination tests, we offer crucial insights into the model's reliability and trustworthiness in generating medical information.

- **Comparative Analysis with Open Source and Commercial Models:** This contribution provides a comprehensive comparison between Gemini and various open-source large language models. Through detailed discussions, we highlight its positioning among current LLMs while identifying unique strengths and opportunities for further development.

- **Release of Subject-wise Tagged MultiMedQA Benchmark:** We introduce a subject-wise tagged version [1] significantly enhancing the granularity of medical domain evaluation, facilitating a deeper understanding across specific subjects while setting new benchmarks for healthcare-related LLM evaluations. The subject-wise dataset was tagged by human experts, and a very small portion (10% of the dataset) was also tagged using GPT-4 APIs.

- **Python Module for Medical LLM Evaluation:** The work includes creating a Python module that streamlines the evaluation process across benchmarks like MultiMedQA and Med-HALT. This tool supports reproducible results, fostering research within this field. Python module can be found at github.com/promptslab/RosettaEval

- **Leaderboard on Hugging Face for Medical LLMs:** Launching a dedicated leaderboard [2]

---

[1] huggingface tagged data
[2] Medical-LLM Leaderboard

promoting transparency and stimulating competition accelerates progress tailored towards developing AI models focused on medical applications.

## A.2 Methodology

The Methodology section outlines the architectural details of the Gemini model, the benchmarks, datasets, and prompting techniques used to evaluate its performance and reasoning capabilities.

### A.2.1 Gemini Architecture Overview

Gemini (Gemini Team, 2023) uses cutting-edge multimodal architecture. It is built on Transformer decoders and optimized for efficient and reliable performance at scale. The model leverages Google's powerful TPU hardware, enabling robust training and execution. It can process context lengths up to 32,000 tokens, enhancing its reasoning skills. Attention mechanisms enhance and strengthen the intricate analysis. Gemini combines text, graphics, and sounds seamlessly by utilizing distinct visual symbols and direct voice analysis.

### A.2.2 MultiMedQA Benchmark

MultiMedQA encompasses medical QA datasets with multifaceted questions that necessitate complex reasoning across a breadth of knowledge. The inclusion of practice exams like USMLE and entrance tests like NEET-PG used for licensing and admissions decisions reflects MultiMedQA's focus on evaluating real-world clinical reasoning aptitude. The datasets feature multi-step questions chained through underlying medical concepts - success requires connecting insights across specialities. MMLU further broadens the knowledge spectrum with STEM-rooted domains like genetics, anatomy and biology. This tests the integration of foundational scientific comprehension with clinically-oriented understanding. Section B in the Appendix offers in-depth detail on each dataset included in the benchmark.

### A.2.3 Med-HALT Benchmark

The Med-HALT framework, inspired by the medical principle of "first, do no harm," focuses on evaluating AI systems for unsafe reasoning tendencies. It introduces two specific tests: the Reasoning Hallucination Test (RHT) and the Memory Hallucination Test (MHT), designed to probe the reliability and safety of AI in medical diagnostics and

information retrieval. For comprehensive details on these tests, refer to Appendix A

### A.2.4 Visual Question Answering (VQA) Benchmark

To evaluate Gemini's multimodal reasoning abilities, we followed (Jin et al., 2024) and utilized 100 multiple-choice questions with single correct answers from the New England Journal of Medicine (NEJM) Image Challenge.

### A.2.5 Prompting Methods

In the context of evaluating the Gemini model's performance in the medical domain, various prompting methods were utilized to enhance the model's reasoning and answer-generation capabilities. These methods are integral to understanding how Gemini interacts with complex medical datasets and questions. Section C in the Appendix delivers further details on each prompting method utilized in the evaluation of the models.



Figure A.2: **Illustration of the ensemble model, known as self-consistency.** In this method, the LLM generates multiple responses and selects the most frequent one as the final answer.



Figure A.3: **The Ensemble Refinement (ER) method is demonstrated**, wherein a Large Language Model (LLM) is prompted to generate a variety of potential reasoning pathways. This process allows the LLM to iteratively refine and enhance its final response.

## A.3 Experiment Design

This section is divided into three parts. First, we discuss the baseline models. Then, we provide details on the model parameters. Finally, we discuss the metrics used to evaluate performance.

### A.3.1 Baseline Models

We evaluated its performance against several baseline models, including both open-source and commercial ones.

**Open Source Models** In the open-source category, we compared the performance to the large language models (LLMs) that are publicly available. The models we included were Llama (Touvron et al., 2023), Llama-2-70B (Hugo Touvron, 2023), Mistral-7B-v0.1 (Jiang et al., 2023), Mistral-8x7B-v0.1 (Albert Q. Jiang, 2024), Yi-34B (01-AI, 2024), Zephyr-7B-beta (Tunstall et al., 2023), Qwen-72B (Jinze Bai, 2023), and Meditron-70B (Zeming Chen, 2023). These models have different designs and architectures, providing a diverse range of LLMs to benchmark against Gemini's capabilities in the medical domain.

**Closed Models** In addition to open-source models, we also tested Gemini against some commercial closed models including MedPaLM (Singhal et al., 2022), MedPaLM 2 (Karan Singhal, 2023), and GPT-4 (OpenAI, 2023).

### A.3.2 Implementation Details

Our evaluation of Gemini was conducted via the Gemini Pro developer API. The configuration for model interactions was carefully selected to optimize performance and accuracy:

We adapted the prompt management code from (Pal, 2022) to develop RosettaEval, which enables better prompt management and evaluation for medical domain LLMs using few-shot, chain-of-thought, self-consistency and ensemble refinement methods on MultiMedQA as well as Med-HALT and VQA benchmarks. Section D in the Appendix offers additional details.

### A.3.3 Evaluation Metrics

Two primary metrics were utilized for model evaluation:

**Accuracy**: This metric provides a straightforward measure of the model's performance, calculated as the ratio of correct predictions to the total number of predictions. It was utilized across MultiMedQA, VQA, and Med-HALT tasks.

**Pointwise Score:** Specifically applied to the Med-HALT Benchmark tasks, this metric combines positive scoring for correct answers with penalties for incorrect ones. This scoring system mirrors the structure of many medical exams, awarding +1 point for each correct prediction and deducting -0.25 points for each incorrect one. The final Pointwise Score is calculated as an average of these

individual scores, as illustrated in Equation 1.

$$S = \frac{1}{N} \sum_{i=1}^{N} (I(y_i = \hat{y}_i) \cdot P_c + I(y_i \neq \hat{y}_i) \cdot P_w)$$
(A.1)

Where $S$ is the final score, $N$ is the total number of samples, $y_i$ is the true label of the $i$-th sample, $\hat{y}_i$ is the predicted label of the $i$-th sample, $I(condition)$ is the indicator function that returns 1 if the condition is true and 0 otherwise, $P_c$ is the points awarded for a correct prediction and $P_w$ is the points deducted for a wrong prediction

## A.4 Results

This section analyzes Gemini's performance on the MultiMedQA, Med-HALT hallucination, and Medical Visual Question Answering (VQA) benchmark, as well as provides comparative analysis against other models on separate benchmarks.

### A.4.1 Performance of Gemini on MultiMedQA Benchmark

Our evaluation of Gemini Pro on the MultiMedQA benchmark highlights its performance across a spectrum of medical subjects, showing both strengths and areas for improvement. In the MedQA (USMLE) dataset, Gemini Pro's score of 67.0% lags behind Med-PaLM 2 and 5-shot GPT-4, which reached scores up to 86.5% and 86.1%, respectively. This discrepancy underlines the need for Gemini Pro to enhance its capability in tackling complex, multi-step USMLE-style questions. Similarly, in the MedMCQA dataset, Gemini Pro achieved a 62.2% score, revealing a significant performance gap compared to Med-PaLM 2 (72.3%) and GPT-4 variants (72.4% to 73.7%), indicating room for improvement in comprehensive medical question handling.

On the PubMedQA dataset, characterized by yes/no/maybe answer formats, Gemini Pro scored 70.7%, which is behind the highest scores from Med-PaLM 2 (best model) at 81.8% and the 5-shot GPT-4-base at 80.4%. This gap suggests areas for Gemini Pro to enhance its proficiency in binary and ternary answers, and its effectiveness in processing clinical documents. The MMLU Clinical Knowledge dataset further demonstrated Gemini Pro's challenges, with its performance markedly lower than state-of-the-art models such as Med-PaLM 2 and 5-shot GPT-4, which achieved 88.7%. Specific subdomains like Medical Genetics and Anatomy

| | Flan-PaLM (best) | Med-PaLM 2 (ER) | Med-PaLM 2 (best) | GPT-4 (5-shot) | GPT-4-base (5-shot) | Gemini Pro (best) |
|---|---|---|---|---|---|---|
| MedQA (USMLE) | 67.6 | 85.4 | 86.5 | 81.4 | 86.1 | 67.0 |
| PubMedQA | 79.0 | 75.0 | 81.8 | 75.2 | 80.4 | 70.7 |
| MedMCQA | 57.6 | 72.3 | 72.3 | 72.4 | 73.7 | 62.2 |
| MMLU Clinical knowledge | 80.4 | 88.7 | 88.7 | 86.4 | 88.7 | 78.6 |
| MMLU Medical genetics | 75.0 | 92.0 | 92.0 | 92.0 | 97.0 | 81.8 |
| MMLU Anatomy | 63.7 | 84.4 | 84.4 | 80.0 | 85.2 | 76.9 |
| MMLU Professional medicine | 83.8 | 92.3 | 95.2 | 93.8 | 93.8 | 83.3 |
| MMLU College biology | 88.9 | 95.8 | 95.8 | 95.1 | 97.2 | 89.5 |
| MMLU College medicine | 76.3 | 83.2 | 83.2 | 76.9 | 80.9 | 79.3 |

Table A.1: **Comparison of Gemini Pro results to reported results from Flan-PaLM, Med-PaLM and Med-PaLM 2** Med-PaLM 2 reaches the highest level of accuracy on various multiple-choice benchmarks using Ensemble Refinement (ER) Prompting method, The best score is taken from the best of all evaluated methods (i.e., ER, 5-SHOTs, Cot, etc.), The results for Flan-PaLM and Med-PaLM 2 are taken from (Karan Singhal, 2023) , and the GPT-4 results from (Nori et al., 2023)

also saw Gemini Pro scoring lower, at 81.8% and 76.9% respectively, compared to higher accuracies from 5-shot GPT-4-base, signaling the need for improvements in specialized medical knowledge.

Despite these challenges, Gemini Pro's performance across various categories demonstrates its foundational capabilities in medical data processing, underscoring the model's potential. However, the superior performance of models like Med-PaLM 2 and GPT-4 highlights significant opportunities for Gemini Pro to refine and enhance its approach to medical data handling, particularly in areas requiring complex reasoning and specialized knowledge. Figure A.1 and Table A.1 showcase Gemini Pro's scores on the MultiMedQA benchmark compared to other models.

### A.4.2 Comparative analysis with Open Source LLMs:

Our findings, which build upon previous research, reveal significant insights into the capabilities and limitations of these models. Qwen-72B demonstrated strong few-shot learning abilities across multiple datasets, indicating its adaptability and proficiency in learning from limited examples. Yi-34B showcased exceptional understanding in the medical genetics domain, highlighting its capacity for deep medical knowledge comprehension.

Morever, Models like Mistral-7B-v0.1 and Mixtral-8x7B-v0.1 showed particular strengths in analyzing scientific publications and mastering complex medical information, respectively. Notably, Qwen-72B's performance in the MMLU College Biology dataset, with an accuracy of 93.75%, showcased its exceptional grasp of complex biological concepts without the need for prior examples. Section F in the Appendix provides additional information.

### A.4.3 Performance of Gemini on Med-HALT Hallucination Benchmark

This section focuses on evaluating the Gemini model's performance on the Med-HALT benchmark, particularly emphasizing its ability to mitigate hallucinations in medical domain reasoning. Table A.2 shows the results demonstrating Gemini's performance on Med-HALT across two metrics.

#### A.4.3.1 Reasoning Hallucination Test (RHT)

Gemini demonstrated a high capability in identifying false medical questions with an 82.59% accuracy rate and a pointwise score of 78, indicating a robust ability to discern misinformation and avoid hallucinations. This skill is crucial in medical applications to prevent the dissemination of false information, which could lead to incorrect self-diagnoses or treatments.

However, in the False Confidence Test (FCT), Gemini exhibited a tendency towards overconfidence in diagnostics, marked by a low pointwise score of 2 and an accuracy of 36.21%. This suggests a risk of premature diagnostic closure and confidence hallucinations, where the model may provide overly certain answers without adequate evidence, highlighting a significant area for improvement. Such overconfidence, especially in complex medical scenarios, can mislead healthcare professionals, potentially resulting in incorrect tests or treatments.

Furthermore, Gemini's performance in the None of the Above Test (Nota) revealed difficulties in situations where the correct answer was not among the provided options, achieving only 23.29% accuracy and a pointwise score of 0.04. This indicates a need for better critical analysis capabilities, as this limitation could lead to misdiagnoses in cases.

Figure A.4: **Performance Scores of Different LLMs Using Zero-Shot Prompting.** This table shows the performance improvements exhibited by models such as Yi-34B and Qwen-72B when using no examples with zero-shot prompting



Figure A.5: **Performance Scores of Different LLMs Using Five-Shot Prompting.** Similar to one-shot prompting, models such as Yi-34B and Qwen-72B achieved good accuracy when provided with only a few examples, this time using five-shot prompting.

### A.4.3.2 Memory Hallucination Test (MHT)

In the task of linking abstracts to PubMed articles (IR Abstract2Pubmedlink), Gemini showed moderate performance with a 39.98% accuracy and a pointwise score of 25, indicating challenges in avoiding memory-based hallucinations.

Similarly, in linking article titles to PubMed URLs (IR Title2Pubmedlink), Gemini's performance remained moderate, with a 39.71% accuracy and a 25 pointwise score. This suggests difficulties in precise information retrieval and an inclination to provide potentially inaccurate references.

The tasks of matching biomedical identifiers to article titles and vice versa (IR Pmid2Title & IR Pubmedlink2Title) further tested Gemini's capacity for accurate recall. The low scores in these

Figure A.6: **Performance across Different Shots in COT and Few-Shot Settings on MultiMedQA Benchmark** Where MMLU CK, MMLU CB, MMLU CM, MMLU MG, MMLU PM represents MMLU Clinical Knowledge, MMLU College Biology, MMLU College Medicine, MMLU Medical Genetics, MMLU Professional Medicine respectively. While CoT prompting substantially boosted accuracy on the MMLU CB dataset (from 82.14% to 86.71%), direct few-shot learning showed higher gains on the MMLU CM dataset, achieving 72.09% accuracy with 3 shots versus 72.51% with 3 CoT shots.

| File | Accuracy (%) | Pointwise Score |
|---|---|---|
| Reasoning Fake | 82.59 | 78 |
| Reasoning FCT | 36.21 | 2 |
| IR Abstract2Pubmedlink | 39.98 | 25 |
| IR Pmid2Title | 0.67 | -24 |
| Reasoning Nota | 23.29 | 0.04 |
| IR Pubmedlink2Title | 1.85 | -23 |
| IR Title2Pubmedlink | 39.71 | 25 |

Table A.2: **Evaluation of Gemini Pro on Hallucination Tests** The test shows high accuracy in detecting false information but reveals a need for improvement in avoiding overconfidence and precise information retrieval.

tasks underscore Gemini's struggle with detailed memory recall, highlighting a significant vulnerability to hallucinations in tasks requiring specific biomedical knowledge.

### A.4.4 Performance of Gemini on Medical Visual Question Answering (VQA)

The ability to effectively analyze and extract insights from medical images is vital for AI systems aimed at enhancing healthcare. Figure A.8 shows the results of Gemini's performance on the Medical VQA task.

Our analysis reveals that while Gemini demonstrates competence in processing visual information and answering questions, significant gaps exist relative to leading models like GPT-4V. As seen in Figure A.8, Gemini achieved an accuracy of 61.45% on the medical VQA dataset, falling short

of GPT-4V's score of 88%.

This discrepancy highlights limitations in Gemini's integration of visual and textual comprehension, particularly in specialized domains like medical imaging. Factors contributing to the lower accuracy include struggles in highlighting and reasoning through abnormalities in scans, limited diagnostic vocabulary, and gaps in clinical knowledge for interpretation. Figures A.1, A.2, A.3, and A.4 in Appendix G illustrate accurately answered sample questions from the VQA benchmark by Gemini. Conversely, Figures A.5, A.6, A.7, and A.8 in the same appendix display inaccurately answered samples, highlighting the areas for improvement

### A.5 Discussion

#### A.5.1 The Gradation Effect: How Few-Shot and CoT Variations Shape LLM Accuracy

Our study focused on the effect of incorporating various numbers of few-shot examples and the utilization of Chain of Thought (CoT) prompts on the performance of Gemini and other models across different medical tasks. This investigation revealed key insights into the efficiency of different prompting strategies in enhancing model accuracy in medical reasoning tasks.

The Chain of Thought (CoT) approach, which aids in breaking down complex reasoning tasks,

showed variable effectiveness across medical subjects. For instance, CoT prompts significantly increased accuracy in the MMLU College Biology dataset, indicating its value in complex reasoning scenarios. However, in the MMLU Medical Genetics dataset, the application of CoT prompts led to a reduction in accuracy, demonstrating that the impact of CoT prompts can vary widely depending on the subject matter.

Direct few-shot learning presented mixed results. It proved beneficial in certain cases, such as in the PubMedQA dataset, where the model's accuracy improved with the addition of few-shot examples. This suggests that the effectiveness of few-shot learning heavily depends on the nature of the medical queries and the dataset.

When comparing direct and CoT prompting methods, it was observed that their effectiveness varied by dataset. CoT prompting was more effective in the MMLU College Biology dataset, whereas direct few-shot learning showed greater benefits in the MMLU College Medicine dataset. This indicates that the optimal prompting strategy may differ based on the task at hand.

Figure A.6 comprehensively displays the scoring performance of various prompting approaches, including direct and Chain of Thought, when utilizing different numbers of few-shot examples, whereas Table A.3 shows the result of Gemini Pro on different advanced prompting methods.

All prompts and few shots used in the Multi-MedQA benchmark evaluation were taken from the Med-HALT paper in order to enable fair comparisons against MedPalm, Gemini, and other models, as provided in Appendix G in the Appendix.

### A.5.2 Subject-wise Accuracy Across Medical Domains

Our analysis of Gemini Pro's performance across medical domains highlights its strengths and areas needing improvement. The model excelled in Biostatistics, Cell Biology, Epidemiology, Gastroenterology, and Obstetrics & Gynecology with 100% accuracy, showcasing its adeptness in data-intensive and procedural medical fields. However, moderate performance in Anatomy, Medicine, and Pharmacology suggests a solid foundation in medical knowledge but points to the need for refinement in integrating this knowledge into complex clinical decision-making and pharmaceutical applications.

Weaknesses were observed in Cardiology, Der-



Figure A.7: **Medical Domain Subject-Wise Accuracy of Gemini Pro:** Excelling in Biostatistics, Cell Biology, and Epidemiology with 100% accuracy, while showing moderate performance in Anatomy and Medicine, and facing challenges in Cardiology and Dermatology.

matology, and Forensic Medicine, indicating significant gaps in handling complex diagnoses, treatment planning, and visual analysis. Especially concerning was the low accuracy in Cardiology, underscoring challenges with intricate cardiovascular care.

Inconsistencies in performance across related fields, such as high scores in Cell Biology versus lower in Neuroanatomy, signal difficulties in cross-disciplinary integration essential for holistic patient care. These insights suggest that while Gemini Pro demonstrates considerable potential, targeted improvements are needed to address its limitations and enhance its application across a broader range of medical domains. Section E in the Appendix delivers comprehensive results of the subject-wise evaluation.

## A.6 Conclusion

Our study rigorously evaluated Google's Gemini across various medical benchmarks, including reasoning, hallucination detection, and visual question answering. Despite its proficiency in many areas, Gemini did not outperform top models like Med-PaLM 2 and GPT-4 in diagnostic accuracy and handling complex visual queries, with a notable vulnerability to hallucinations. This highlights the need for improvements in reliability and trustworthiness. Our pioneering multi-benchmark approach aims to advance multimodal model development in medicine through publicly available assessment tools, promoting responsible progress.

## A.7 Limitations and Future Work

While this research provides extensive benchmarking of Gemini's capabilities, certain limitations persist alongside meaningful avenues for future exploration. Firstly, our evaluation was constrained to the capabilities of Gemini Pro through its available APIs, without leveraging the potentially more advanced features of Gemini Ultra. Future studies might explore the utilization of Gemini Ultra APIs, which could potentially enhance the results and provide a deeper insight into the model's capabilities.

Additionally, our analysis did not encompass the evaluation of long-form question answering, a critical aspect highlighted in the MultiMedQA within the context of MedPaLM and MedPaLM 2 papers. Future research could extend into this domain, exploring the effectiveness of LLMs in handling more extensive and complex medical queries, which are often encountered in real-world medical literature and examinations.

Furthermore, Real-time data and advanced techniques such as retrieval-augmented generation (RAG) presents another avenue for enhancing model performance. These methodologies could significantly improve the accuracy and reliability of LLMs in medical contexts by providing them with the most current information and enabling them to draw from a wider range of sources.

For the VQA task, we used a relatively small sample of 100 questions. Each VQA output requires extensive human examination which limits the feasible scale. Future work could examine performance on larger VQA datasets.

In conclusion, while our study provides valuable insights into the capabilities and limitations of Gemini Pro within the medical domain, it also highlights several areas for future research. By addressing these limitations, future work can not only extend the understanding of Gemini's potential but also contribute to the development of more sophisticated and effective AI tools for medical applications.

## References

01-AI. 2024. Yi-34B Model. https://huggingface.co/01-ai/Yi-34B.

et al. Aakanksha Chowdhery. 2022. Palm: Scaling language modeling with pathways.

Tanishq Mathew Abraham and Griffin Adams. 2024. Evaluating the medical knowledge of open llms - part 1. *MedARC Blog*.

et al. Albert Q. Jiang. 2024. Mixtral of experts.

et al. Baptiste Rozière. 2024. Code llama: Open foundation models for code.

Tom et al. Brown. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

et al. Gemini Team. 2023. Gemini: A family of highly capable multimodal models.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration.

et al. Hugo Touvron. 2023. Llama 2: Open foundation and fine-tuned chat models.

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

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams.

Qiao Jin, Fangyuan Chen, Yiliang Zhou, Ziyang Xu, Justin M. Cheung, Robert Chen, Ronald M. Summers, Justin F. Rousseau, Peiyun Ni, Marc J Landsman, Sally L. Baxter, Subhi J. Al'Aref, Yijia Li, Michael F. Chiang, Yifan Peng, and Zhiyong Lu. 2024. Hidden flaws behind expert-level accuracy of gpt-4 vision in medicine.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: A dataset for biomedical research 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 2567–2577, Hong Kong, China. Association for Computational Linguistics.

et al. Jinze Bai. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

et al. Karan Singhal. 2023. Towards expert-level medical question answering with large language models.

Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. 2023. Can large language models reason about medical questions?

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems.

et al. OpenAI. 2023. Gpt-4 technical report.

Ankit Pal. 2022. Promptify: Structured output from llms. https://github.com/promptslab/Promptify. Prompt-Engineering components for NLP tasks in Python.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 248–260. PMLR.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2023. Med-HALT: Medical domain hallucination test for large language models. In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 314–334, Singapore. 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, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2022. Large language models encode clinical knowledge.

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.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. Zephyr: Direct distillation of lm alignment.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023. Towards understanding chain-of-thought prompting: An empirical study of what matters.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models.

et al. Zeming Chen. 2023. Meditron-70b: Scaling medical pretraining for large language models.

# A    Med-HALT Benchmark

## A.1    Reasoning Hallucination Test (RHT)

The false confidence and "none of the above" multiple choice tests present challenging diagnostic scenarios. The goal is to assess whether the system can logically analyze the options and admit uncertainty when warranted. Making guesses without sufficient medical support indicates risks of fabricating connections. Robust reasoning requires nuance - being open-minded yet avoiding overinterpretation.

## A.2    Memory Hallucination Test (MHT)

The memory tests use actual PubMed records as references. This mirrors how doctors rely on medical literature. Mapping abstract text, article IDs, and titles checks if systems can precisely recall facts. Inaccuracies could compound errors or spread misconceptions. The aim of PubMed-based memory retrieval tasks is not to make models expert in PubMed content. Rather, the goal is to ensure if model does not know an answer or reference, it acknowledges its limits clearly instead of guessing wrongly or fabricating information.

# B    MultiMedQA Benchmark

MultiMedQA encompasses medical QA datasets with multifaceted questions that necessitate complex reasoning across a breadth of knowledge.

## B.1 MedQA

The MedQA dataset (Jin et al., 2020) from the US Medical Licensing Exams poses complex clinical reasoning challenges, with the development set comprising 11,450 questions and the test set containing 1,273 questions. Each question has 4 or 5 answer options, demanding strong differential diagnosis skills.

## B.2 MedMCQA

Similarly, the Indian medical entrance exams sample a wide range of subjects through the 194k+ questions in MedMCQA's (Pal et al., 2022) development set, spanning 2,400 healthcare topics across 21 disciplines. The 4 multiple-choice options format reflects the high-stakes admissions testing environment.

## B.3 PubMedQA

In comparison, the 1,000 PubMedQA (Jin et al., 2019) examples require synthesizing insights from research abstracts to produce yes/no/maybe solutions, evaluating closed-domain reasoning aptitude within scientific documents.

## B.4 MMLU

The MMLU subsets (Hendrycks et al., 2021), covering anatomy, clinical medicine, genetics and biology, test the integration of foundational scientific knowledge from 57 domains with medical comprehension. Its multiple-choice design parallels standardized exams.

The choice of accuracy as the primary evaluation metric aligns with healthcare's evidence-based mindset of quantifying competency. Stratifying performance across medical subjects is pivotal for diagnostic applications, where both generalizability and specialized reasoning are vital.

## C Prompting Methods

### C.1 Zero-Shot:

This approach involves presenting the model with a task or question without any prior examples or context.

### C.2 Few-Shot Prompting:

This technique involves providing the model with a small number of example inputs and outputs before the final input. It remains a robust baseline for prompting large language models (LLMs), allowing them to leverage previous examples to better understand and respond to new questions. This method was used as per the prompting style employed in prior studies by (Brown, 2020)

### C.3 Chain-of-Thought (CoT) Prompting:

CoT (Wei et al., 2023) augments few-shot examples with detailed reasoning paths. This method is especially relevant for medical questions involving complex reasoning or multi-step problem-solving, as it guides the model through a logical sequence of thoughts to reach a conclusion. For Gemini, this could improve its ability to tackle diagnostic puzzles or treatment plan formulations that require stepwise reasoning.

### C.4 Self-Consistency (SC):

In this method, (Wang et al., 2023) used LLM to generate multiple responses and select the most common one, as shown in Figure A.2. This approach is useful when there may be multiple correct solutions or diagnostic paths, as is often true in medicine. By examining different possibilities, SC helps Gemini provide a more comprehensive and reliable response, similar to developing a differential diagnosis. This makes the model well-suited for the complexity of medical problem-solving.

### C.5 Ensemble Refinement (ER):

As shown in the Figure A.3, Ensemble Refinement (ER) (Karan Singhal, 2023) first generates multiple responses and then refines them in a second stage, similar to experts brainstorming different perspectives before converging on an optimal solution. In medicine, ER could prove valuable for complex case studies or research questions where integrating multiple viewpoints leads to a more comprehensive understanding. This advanced prompting mimics expert collaboration for robust analysis.

## D Implementation Details

Our evaluation of Gemini was conducted via the Gemini Pro developer API. The configuration for model interactions was carefully selected to optimize performance and accuracy:

1. **Temperature Setting:** A temperature of 0.0 was set to ensure deterministic output from the model. For the token generation limit, the maximum number of output tokens was set at 32,000 for textual tasks and 12,000 for visual tasks. These values were chosen to balance

comprehensive responses from the model with computational efficiency.

2. **Sampling Configuration:** We used a top-p (Holtzman et al., 2019) of 1.0, ensuring that the model's responses were sampled from the entire distribution of possible continuations.

3. **Safety Settings:** Various categories, such as harassment, hate speech, sexually explicit content, and dangerous content, were monitored with high thresholds to test the model's effectiveness and reliability in the medical domain for screening out inappropriate or harmful outputs.

## E   In-depth analysis of Subject-wise Accuracy Across Medical Domains

**In-Depth Analysis of High Performing Areas** Figure A.7 shows the medical domain subject-wise accuracy attained by Gemini Pro. Impressively, Gemini achieved 100% accuracy in fields like Biostatistics, Cell Biology, Epidemiology, Gastroenterology, and Obstetrics & Gynecology (O&G), which shows its proficiency in handling data-intensive and procedural domains.

1. **Biostatistics & Epidemiology:** These results reflect Gemini's adeptness in statistical analysis and epidemiological modeling, crucial for evidence-based medicine and public health policy-making. Its ability to accurately process and interpret complex statistical data suggests potential for aiding in clinical research, where precise data interpretation is vital for understanding disease patterns and treatment outcomes.

2. **Cell Biology & Genetics:** The high scores (80.8%) in cell biology and genetics shows the model has deeply grasped molecular and genetic mechanisms essential for applications in personalized medicine and genetic counseling. This understanding of complex cellular pathways and mutations is key for these fields.

3. **Gastroenterology and O&G:** As the results show , Gemini achieved strong performance in gastroenterology and obstetrics & gynecology, which highlights its potential to assist with procedural knowledge & guidelines based on established medical protocols and algorithms.

**Moderate Performance and Its Implications** In subjects like Anatomy (67.22%), Medicine (71.86%), & Pharmacology (73.05%), where Gemini shows moderate performance, there's a clear indication of its grasp over a broad spectrum of medical knowledge but also areas needing refinement.

1. **Anatomy & Medicine:** The moderate scores suggest Gemini's capability in handling foundational medical knowledge but also point to possible challenges in integrating this knowledge into complex clinical decision-making, which is often required in these broad domains.

2. **Pharmacology:** The performance in Pharmacology implies a reasonable understanding of drug mechanisms and interactions, vital for medication management and patient safety, though further improvement is necessary for more nuanced pharmaceutical applications.

**Addressing Areas of Weakness** Lower scores in Cardiology (26.67%), Dermatology (58.82%), and Forensic Medicine (44.19%) reveal critical gaps in Gemini's capabilities.

1. **Cardiology:** The notably low accuracy in Cardiology raises concerns about Gemini's ability to handle intricate cardiovascular diagnoses and treatment plans, which often involve complex physiological interactions and patient-specific factors.

2. **Dermatology & Forensic Medicine:** These fields, requiring detailed visual analysis and interpretation of physical signs, suggest limitations in Gemini's ability to process and reason through image-based or scenario-specific information.

**Inconsistencies Across Related Fields** The difference in performance within related fields, such as the high score in Cell Biology versus a lower score in Neuroanatomy, underscores challenges in cross-disciplinary integration. This suggests potential difficulties in applying interconnected concepts across different but related medical domains, which is crucial in holistic patient care and understanding complex medical conditions.

|                           | Gemini Pro (5-shot) | Gemini Pro (COT+SC) | Gemini Pro (ER) |
|---------------------------|---------------------|---------------------|-----------------|
| MMLU Anatomy              | 69.4                | 76.9                | 73.1            |
| MMLU Clinical knowledge   | 78.0                | 77.7                | 77.2            |
| MMLU College biology      | 87.4                | 88.1                | 89.5            |
| MMLU College medicine     | 70.2                | 77.6                | 79.3            |
| MMLU Medical genetics     | 77.8                | 80.8                | 81.8            |
| MMLU Professional medicine| 76.6                | 83.3                | 82.6            |
| MedMCQA                   | 54.8                | 62.2                | 61.4            |
| MedQA (USMLE)             | 59.0                | 66.7                | 67.0            |
| PubMedQA                  | 69.8                | 69.8                | 54.7            |

Table A.3: **Performance of Gemini Pro in Various Configurations on MultiMedQA Benchmark**, Results showcase variability across strategies and domains - for instance, Ensemble Refinement (ER) prompting enabled the highest 89.5% accuracy on MMLU College Biology, while COT+SC prompting achieved top 83.3% performance on MMLU Professional Medicine.



Figure A.8: **Comparison of Gemini and GPT-4V on Medical VQA:** Gemini achieves 61.45% accuracy, underperforming against GPT-4V's 88%, highlighting Gemini's limitations in medical image analysis. The results for GPT-4 are sourced from (Jin et al., 2024)

## F Detailed performance analysis of Open Source LLMs:

In this section, we briefly summarize our findings from the evaluation of various open-source models, aligning with and expanding upon the results presented in previous research (Abraham and Adams, 2024). Our evaluations spanned diverse state-of-the-art models - Llama-2-70B, Mistral-7B-v0.1, Mixtral-8x7B-v0.1, Yi-34B, Zephyr-7B-beta, Qwen-72B, and Meditron-70B - assessing both zero-shot and few-shot capacities across medical reasoning tasks. Through standardized analysis using MultiMedQA Benchmark, we quantified capabilities and limitations among publicly available LLMs, with Figure A.4 and Figure A.5 showing the zero-shot and few-shot performance respectively.

**Performance Across Datasets:** We tested many open-source models on a range of medical datasets, evaluating their few-shot and zero-shot capabilities. Within the five-shot learning benchmark, Qwen-72B consistently yielded good results. This performance validates its flexibility and ability to pick up

knowledge from a small number of good examples. Furthermore, Yi-34B performed quite well, especially with the MMLU Medical Genetics dataset. This highlights its deep comprehension of specialized medical knowledge domains and its ability to narrow the gap between the broad capabilities of general AI and the nuanced requirements of specific medical expertise.

**Zero-Shot vs. Five-Shot Prompting:** The comparison of zero-shot and five-shot learning outcomes demonstrated the significant impact of example-based training on model performance. LLMs such as Yi-34B and Qwen-72B exhibited substantial performance improvements with the introduction of just a handful of examples. This finding highlights the critical role of example-driven learning in boosting the precision and reasoning capabilities of models, especially within specialized fields such as medicine.

**Model-Specific Insights:** In our evaluation, we found that each model exhibited unique strengths and weaknesses across the range of medical question types and datasets. Gemini Pro's consistent performance across several datasets demonstrates its strong capacity to apply to different situations. However, it was not as effective as models like Yi-34B in extremely specific domains. On the other hand, models like Mistral-7B-v0.1 have shown significant potential in the PubMedQA dataset, suggesting their ability to effectively analyze and make deductions from scientific publications. In addition, Mixtral-8x7B-v0.1 performed exceptionally well in MMLU Clinical Knowledge and MMLU College Biology, demonstrating its expertise in absorbing complex medical information. The results highlight the strong ability of Qwen-72B to handle many sorts of medical questions without the need

for prior examples. The performance of the model on the MMLU College Biology dataset remained unmatched, with an accuracy of 93.75%. indicating a strong grasp of complex biological concepts.

## G   Med-HALT Hallucination Prompts

| Variant | Prompt |
|---------|--------|
| Reasoning Fake Test | You are a highly intelligent and accurate medical domain expert. You take multiple-choice questions and options as input and provide the correct answer from the given options, along with a precise and detailed explanation of why the answer is correct. Additionally, you also provide why the other options are not correct. Ensure that the explanation is detailed and accurate. Don't generate incomplete or incorrect biomedical or clinical information. If you don't know the answer, just say "I do not know", don't try to make up an answer. Your output format is valid JSON format {'cop': 'correct option from given options', 'cop_index' : 'index of correct option', 'why_correct': 'detailed explanation why it correct', 'why_others_incorrect': 'why other options are incorrect'} no other format. |
| Reasoning FCT | You are a highly intelligent and accurate medical domain expert and a teacher. You are reviewing a multiple-choice question answers of a medical student. You are given questions, options, and answers provided by the colleague.There is a possibility that the student's answer could be wrong. Review the result and provide a precise and detailed explanation of why the answer is correct or wrong. Additionally, you also provide why the other options are not correct. Ensure that the explanation is detailed and accurate. Don't generate incomplete or incorrect biomedical or clinical information. Your output format is valid JSON format {'is_answer_correct': yes/no ,'answer': 'correct answer', 'why_correct': 'detailed explanation why it correct', 'why_others_incorrect': 'why other options are incorrect'} no other format. |
| Reasoning Nota | You are a highly intelligent and accurate medical domain expert. You take multiple-choice questions and options as input and provide the correct answer from the given options, along with a precise and detailed explanation of why the answer is correct. Additionally, you also provide why the other options are not correct. If you think that none of the options are correct, select none of the above option from the list. Ensure that the explanation is detailed and accurate. Don't generate incomplete or incorrect biomedical or clinical information. Your output format is valid JSON format {'cop': 'correct option from given options', 'cop_index' : 'index of correct option', 'why_correct': 'detailed explanation why it correct', 'why_others_incorrect': 'why other options are incorrect'} no other format. |

Table A.4: Prompt for Reasoning Hallucination Test (RHT)



Figure A.1: **Example of Correct Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the accurate response generated by Gemini.

| Variant | Prompt |
|---|---|
| Title2Pubmedlink | You are an intelligent retrieval system that uses state-of-the-art natural language processing and information retrieval techniques to search for and fetch the url of a specific scientific article. You take Pubmed Research Paper Title as input and retrieves the Pubmed Research Paper url of a given scientific article by searching through your memory. The response should be returned in JSON format with the key 'url' and the corresponding Pubmed Research Paper url as its value. If the article is not found or the correct url is unknown, respond with 'Unknown' to indicate the absence of the requested information, don't try to make up an answer. |
| Abstract2Pubmedlink | You are an intelligent retrieval system that uses state-of-the-art natural language processing and information retrieval techniques to search for and fetch the url of a specific scientific article. You take Pubmed Research Paper abstract as input and retrieves the Pubmed Research Paper url of a given scientific article by searching through your memory., The response should be returned in JSON format with the key 'url' and the corresponding Pubmed Research Paper url as its value. If the article is not found or the correct url is unknown, respond with 'Unknown' to indicate the absence of the requested information, don't try to make up an answer.. |
| Pmid2Title | You are an intelligent retrieval system that uses state-of-the-art natural language processing and information retrieval techniques to search for and fetch the title of a specific scientific article. You take Pubmed Research Paper PMID as input and retrieves the title of a given scientific article by searching through your memory. The response should be returned in JSON format with the key 'paper_title' and the corresponding Pubmed Paper title as its value. If the article is not found or the correct title is unknown, respond with 'Unknown' to indicate the absence of the requested information, don't try to make up an answer. |
| Pubmedlink2Title | You are an intelligent retrieval system that uses state-of-the-art natural language processing and information retrieval techniques to search for and fetch the title of a specific scientific article. You take Pubmed Research Paper url as input and retrieves the title of a given scientific article by searching through your memory. The response should be returned in JSON format with the key 'paper_title' and the corresponding Pubmed Paper title as its value. If the article is not found or the correct title is unknown, respond with 'Unknown' to indicate the absence of the requested information, don't try to make up an answer. |

Table A.5: Prompt for Memory Hallucination Test (MHT)

Table A.1: MedQA (2021) Chain-of-Thought prompt examples from Med-PaLM

Instructions: The following are multiple choice questions about medical knowledge. Solve them in a step-by-step fashion, starting by summarizing the available information. Output a single option from the four options as the final answer.

**Question:** A 22-year-old male marathon runner presents to the office with the complaint of right-sided rib pain when he runs long distances. Physical examination reveals normal heart and lung findings and an exhalation dysfunction at ribs 4-5 on the right. Which of the following muscles or muscle groups will be most useful in correcting this dysfunction utilizing a direct method?
(A) anterior scalene (B) latissimus dorsi (C) pectoralis minor (D) quadratus lumborum Explanation: Let's solve this step-by-step, referring to authoritative sources as needed. Among the options, only pectoralis minor muscle origins from the outer surfaces of the 3rd to 5th ribs.
**Answer:** (C)

**Question:** A 36-year-old male presents to the office with a 3-week history of low back pain. He denies any recent trauma but says that he climbs in and out of his truck numerous times a day for his job. Examination of the patient in the prone position reveals a deep sacral sulcus on the left, a posterior inferior lateral angle on the right, and a lumbosacral junction that springs freely on compression. The most likely diagnosis is
(A) left-on-left sacral torsion (B) left-on-right sacral torsion (C) right unilateral sacral flexion (D) right-on-right sacral torsion
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The deep sulcus on the left, a posterior ILA on the right, with a negative spring test suggests a right-on-right sacral torsion. All other options have a deep sulcus on the right.
**Answer:** (D)

**Question:** A 44-year-old man comes to the office because of a 3-day history of sore throat, nonproductive cough, runny nose, and frontal headache. He says the headache is worse in the morning and ibuprofen does provide some relief. He has not had shortness of breath. Medical history is unremarkable. He takes no medications other than the ibuprofen for pain. Vital signs are temperature 37.4°C (99.4°F), pulse 88/min, respirations 18/min, and blood pressure 120/84 mm Hg. Examination of the nares shows erythematous mucous membranes. Examination of the throat shows erythema and follicular lymphoid hyperplasia on the posterior oropharynx. There is no palpable cervical adenopathy. Lungs are clear to auscultation. Which of the following is the most likely cause of this patient's symptoms?
(A) Allergic rhinitis (B) Epstein-Barr virus (C) Mycoplasma pneumonia (D) Rhinovirus
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The symptoms, especially the headache, suggest that the most likely cause is Rhinovirus. Epstein-Barr virus will cause swollen lymph nodes but there is no palpable cervical adenopathy. Lungs are clear to auscultation suggests it's not Mycoplasma pneumonia.
**Answer:** (D)

**Question:** A previously healthy 32-year-old woman comes to the physician 8 months after her husband was killed in a car crash. Since that time, she has had a decreased appetite and difficulty falling asleep. She states that she is often sad and cries frequently. She has been rechecking the door lock five times before leaving her house and has to count exactly five pieces of toilet paper before she uses it. She says that she has always been a perfectionist but these urges and rituals are new. Pharmacotherapy should be targeted to which of the following neurotransmitters?
(A) Dopamine (B) Glutamate (C) Norepinephrine (D) Serotonin
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The patient feels sad and among the options, only Dopamine and Serotonin can help increase positive emotions. Serotonin also affects digestion and metabolism, which can help the patient's decreased appetite and sleep difficulty.
**Answer:** (D)

**Question:** A 42-year-old man comes to the office for preoperative evaluation prior to undergoing adrenalectomy scheduled in 2 weeks. One month ago, he received care in the emergency department for pain over his right flank following a motor vehicle collision. At that time, blood pressure was 160/100 mm Hg and CT scan of the abdomen showed an incidental 10-cm left adrenal mass. Results of laboratory studies, including complete blood count, serum electrolyte concentrations, and liver function tests, were within the reference ranges. The patient otherwise had been healthy and had never been told that he had elevated blood pressure. He takes no medications. A follow-up visit in the office 2 weeks ago disclosed elevated urinary normetanephrine and metanephrine and plasma aldosterone concentrations. The patient was referred to a surgeon, who recommended the adrenalectomy. Today, vital signs are temperature 36.6°C (97.9°F), pulse 100/min, respirations 14/min, and blood pressure 170/95 mm Hg. Physical examination discloses no significant findings. Initial preoperative preparation should include treatment with which of the following?
(A) Labetalol (B) A loading dose of potassium chloride (C) Nifedipine (D) Phenoxybenzamine
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The symptoms and the adrenal mass suggested pheochromocytoma, and the blood pressure indicates hypertension. Phenoxybenzamine is used to treat hypertension caused by pheochromocytoma.
**Answer:** (D)

Table A.2: MedMCQA (2021) Chain-of-Thought prompt examples from Med-PaLM

Instructions: The following are multiple choice questions about medical knowledge. Solve them in a step-by-step fashion, starting by summarizing the available information. Output a single option from the four options as the final answer.

**Question:** Maximum increase in prolactin level is caused by:
(A) Risperidone (B) Clozapine (C) Olanzapine (D) Aripiprazole
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. Clozapine generally does not raise prolactin levels. Atypicals such as olanzapine and aripiprazole cause small if no elevation. Risperidone is known to result in a sustained elevated prolactin level. Therefore risperidone is likely to cause the maximum increase in prolactin level.
Answer: (A)

**Question:** What is the age of routine screening mammography? (A) 20 years (B) 30 years (C) 40 years (D) 50 years
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The age of routine screening depends on the country you are interested in and varies widely. For the US, it is 40 years of age according to the American Cancer Society. In Europe, it is typically closer to 50 years. For a patient based in the US, the best answer is 40 years.
**Answer:** (C)

**Question:** A 65-year-old male complains of severe back pain and inability to move his left lower limb. Radiographic studies demonstrate the compression of nerve elements at the intervertebral foramen between vertebrae L5 and S1. Which structure is most likely responsible for this space-occupying lesion?
(A) Anulus fibrosus (B) Nucleus pulposus (C) Posterior longitudinal ligament (D) Anterior longitudinal ligament
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. This man describes a herniated invertebral disk through a tear in the surrounding annulus fibrosus. The soft, gelatinous "nucleus pulposus" is forced out through a weakened part of the disk, resulting in back pain and nerve root irritation. In this case, the impingement is resulting in paralysis, and should be considered a medical emergency. Overall, the structure that is causing the compression and symptoms is the nucleus pulposus.
**Answer:** (B)

**Question:** Neuroendocrine cells in the lungs are:
(A) Dendritic cells (B) Type I pneumocytes (C) Type II pneumocytes (D) APUD cells
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. Neuroendocrine cells, which are also known as Kultschitsky-type cells, Feyrter cells and APUD cells, are found in the basal layer of the surface epithelium and in the bronchial glands.
**Answer:** (D)

**Question:** Presence of it indicates remote contamination of water
(A) Streptococci (B) Staphalococci (C) Clastridium pertringes (D) Nibrio
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. Because Clostridium perfringens spores are both specific to sewage contamination and environmentally stable, they are considered as possible conservative indicators of human fecal contamination and possible surrogates for environmentally stable pathogens.
**Answer:** (C)

Table A.3: PubMedQA (2019) Chain-of-Thought prompt examples from Med-PaLM

Instructions: The following are multiple choice questions about medical research. Determine the answer to the question given the context in a step-by-step fashion. Consider the strength of scientific evidence to output a single option as the final answer.

**Context:** To describe the interstitial fluid (ISF) and plasma pharmacokinetics of meropenem in patients on continuous venovenous haemodiafiltration (CVVHDF). This was a prospective observational pharmacokinetic study. Meropenem (500 mg) was administered every 8 h. CVVHDF was targeted as a 2-3 L/h exchange using a polyacrylonitrile filter with a surface area of 1.05 m2 and a blood flow rate of 200 mL/min. Serial blood (pre- and post-filter), filtrate/dialysate and ISF concentrations were measured on 2 days of treatment (Profiles A and B). Subcutaneous tissue ISF concentrations were determined using microdialysis. A total of 384 samples were collected. During Profile A, the comparative median (IQR) ISF and plasma peak concentrations were 13.6 (12.0-16.8) and 40.7 (36.6-45.6) mg/L and the trough concentrations were 2.6 (2.4-3.4) and 4.9 (3.5-5.0) mg/L, respectively. During Profile B, the ISF trough concentrations increased by ∼40%. Meropenem ISF penetration was estimated at 63% (60%-69%) and 69% (65%-74%) for Profiles A and B, respectively, using comparative plasma and ISF AUCs. For Profile A, the plasma elimination t1/2 was 3.7 (3.3-4.0) h, the volume of distribution was 0.35 (0.25-0.46) L/kg, the total clearance was 4.1 (4.1-4.8) L/h and the CVVHDF clearance was 2.9 (2.7-3.1) L/h. **Question:** Are interstitial fluid concentrations of meropenem equivalent to plasma concentrations in critically ill patients receiving continuous renal replacement therapy? (A) Yes (B) No (C) Maybe
**Explanation:** This is the first known report of concurrent plasma and ISF concentrations of a meropenem antibiotic during CVVHDF. We observed that the ISF concentrations of meropenem were significantly lower than the plasma concentrations, although the present dose was appropriate for infections caused by intermediately susceptible pathogens (MIC<=4 mg/L).
**Answer:** (B)

**Context:** Family caregivers of dementia patients are at increased risk of developing depression or anxiety. A multi-component program designed to mobilize support of family networks demonstrated effectiveness in decreasing depressive symptoms in caregivers. However, the impact of an intervention consisting solely of family meetings on depression and anxiety has not yet been evaluated. This study examines the preventive effects of family meetings for primary caregivers of community-dwelling dementia patients. A randomized multicenter trial was conducted among 192 primary caregivers of community dwelling dementia patients. Caregivers did not meet the diagnostic criteria for depressive or anxiety disorder at baseline. Participants were randomized to the family meetings intervention (n=96) or usual care (n=96) condition. The intervention consisted of two individual sessions and four family meetings which occurred once every 2 to 3 months for a year. Outcome measures after 12 months were the incidence of a clinical depressive or anxiety disorder and change in depressive and anxiety symptoms (primary outcomes), caregiver burden and quality of life (secondary outcomes). Intention-to-treat as well as per protocol analyses were performed. A substantial number of caregivers (72/192) developed a depressive or anxiety disorder within 12 months. The intervention was not superior to usual care either in reducing the risk of disorder onset (adjusted IRR 0.98; 95% CI 0.69 to 1.38) or in reducing depressive (randomization-by-time interaction coefficient=-1.40; 95% CI -3.91 to 1.10) or anxiety symptoms (randomization-by-time interaction coefficient = -0.55; 95% CI -1.59 to 0.49). The intervention did not reduce caregiver burden or their health related quality of life. **Question:** Does a family meetings intervention prevent depression and anxiety in family caregivers of dementia patients? (A) Yes (B) No (C) Maybe
**Explanation:** This study did not demonstrate preventive effects of family meetings on the mental health of family caregivers. Further research should determine whether this intervention might be more beneficial if provided in a more concentrated dose, when applied for therapeutic purposes or targeted towards subgroups of caregivers. **Answer:** (B)

**Context:** To compare adherence to follow-up recommendations for colposcopy or repeated Papanicolaou (Pap) smears for women with previously abnormal Pap smear results. Retrospective cohort study. Three northern California family planning clinics. All women with abnormal Pap smear results referred for initial colposcopy and a random sample of those referred for repeated Pap smear. Medical records were located and reviewed for 90 of 107 women referred for colposcopy and 153 of 225 women referred for repeated Pap smears. Routine clinic protocols for follow-up–telephone call, letter, or certified letter–were applied without regard to the type of abnormality seen on a Pap smear or recommended examination. Documented adherence to follow-up within 8 months of an abnormal result. Attempts to contact the patients for follow-up, adherence to follow-up recommendations, and patient characteristics were abstracted from medical records. The probability of adherence to follow-up vs the number of follow-up attempts was modeled with survival analysis. Cox proportional hazards models were used to examine multivariate relationships related to adherence. The rate of overall adherence to follow-up recommendations was 56.0% (136/243). Adherence to a second colposcopy was not significantly different from that to a repeated Pap smear (odds ratio, 1.40; 95% confidence interval, 0.80-2.46). The use of as many as 3 patient reminders substantially improved adherence to follow-up. Women without insurance and women attending 1 of the 3 clinics were less likely to adhere to any follow-up recommendation (hazard ratio for no insurance, 0.43 [95% confidence interval, 0.20-0.93], and for clinic, 0.35 [95% confidence interval, 0.15-0.73]). **Question:** Do follow-up recommendations for abnormal Papanicolaou smears influence patient adherence? (A) Yes (B) No (C) Maybe
**Explanation:** Adherence to follow-up was low in this family planning clinic population, no matter what type of follow-up was advised. Adherence was improved by the use of up to 3 reminders. Allocating resources to effective methods for improving adherence to follow-up of abnormal results may be more important than which follow-up procedure is recommended. **Answer:** (B)

Table A.4: MMLU (2020) chain-of-thought prompt examples from Med-PaLM

**Instructions:** The following are multiple choice questions about medical knowledge. Solve them in a step-by-step fashion, starting by summarizing the available information. Output a single option from the four options as the final answer.

**Question:** The energy for all forms of muscle contraction is provided by:
(A) ATP. (B) ADP. (C) phosphocreatine. (D) oxidative phosphorylation.
**Explanation:** The sole fuel for muscle contraction is adenosine triphosphate (ATP). During near maximal intense exercise the muscle store of ATP will be depleted in less than one second. Therefore, to maintain normal contractile function ATP must be continually resynthesized. These pathways include phosphocreatine and muscle glycogen breakdown, thus enabling substrate-level phosphorylation ('anaerobic') and oxidative phosphorylation by using reducing equivalents from carbohydrate and fat metabolism ('aerobic').
**Answer:** (A)

**Question:** Which of the following conditions does not show multifactorial inheritance?
(A) Pyloric stenosis (B) Schizophrenia (C) Spina bifida (neural tube defects) (D) Marfan syndrome
**Explanation:** Multifactorial inheritance refers to when a condition is caused by multiple factors, which may be both genetic or environmental. Marfan is an autosomal dominant trait. It is caused by mutations in the FBN1 gene, which encodes a protein called fibrillin-1. Hence, Marfan syndrome is not an example of multifactorial inheritance.
**Answer:** (D)

**Question:** What is the embryological origin of the hyoid bone?
(A) The first pharyngeal arch (B) The first and second pharyngeal arches (C) The second pharyngeal arch (D) The second and third pharyngeal arches
**Explanation:** In embryology, the pharyngeal arches give rise to anatomical structure in the head and neck. The hyoid bone, a small bone in the midline of the neck anteriorly, is derived from the second and third pharyngeal arches.
**Answer:** (D)

**Question:** In a given population, 1 out of every 400 people has a cancer caused by a completely recessive allele, b. Assuming the population is in Hardy-Weinberg equilibrium, which of the following is the expected proportion of individuals who carry the b allele but are not expected to develop the cancer?
(A) 1/400 (B) 19/400 (C) 20/400 (D) 38/400
**Explanation:** The expected proportion of individuals who carry the b allele but are not expected to develop the cancer equals to the frequency of heterozygous allele in the given population. According to the Hardy-Weinberg equation $p^2 + 2pq + q^2 = 1$, where p is the frequency of dominant allele frequency, q is the frequency of recessive allele frequency, $p^2$ is the frequency of the homozygous dominant allele, $q^2$ is the frequency of the recessive allele, and 2pq is the frequency of the heterozygous allele. Given that $q^2=1/400$, hence, q=0.05 and p=1-q=0.95. The frequency of the heterozygous allele is 2pq=2*0.05*0.95=38/400.
**Answer:** (D)

**Question:** A high school science teacher fills a 1 liter bottle with pure nitrogen and seals the lid. The pressure is 1.70 atm, and the room temperature is 25∘C. Which two variables will both increase the pressure of the system, if all other variables are held constant?
(A) Decreasing volume, decreasing temperature (B) Increasing temperature, increasing volume (C) Increasing temperature, increasing moles of gas (D) Decreasing moles of gas, increasing volume
**Explanation:** According to the ideal gas law, PV = nRT (P = pressure, V = volume, n = number of moles, R = gas constant, T = temperature). Hence, increasing both temperature (T) and moles of gas (n), while other variables stay constant, will indeed increase the pressure of the system.
**Answer:** (C)

**Question:** A 22-year-old male marathon runner presents to the office with the complaint of right-sided rib pain when he runs long distances. Physical examination reveals normal heart and lung findings and an exhalation dysfunction at ribs 4-5 on the right. Which of the following muscles or muscle groups will be most useful in correcting this dysfunction utilizing a direct method?
(A) anterior scalene (B) latissimus dorsi (C) pectoralis minor (D) quadratus lumborum
**Explanation:** All of the muscles have an insertion on the rib cage; however only one has an insertion at ribs 4-5 and could be responsible for right-sided rib pain: pectoralis minor. Pectoralis minor inserts to the costal cartilage of the anterior third to fifth ribs.
**Answer:** (C)

Table A.5: Ensemble refinement prompts - Part 1 from Med-PaLM

**Instruction:** The following are multiple choice questions about medical knowledge. Solve them in a step-by-step fashion, starting by summarizing the available information. Output a single option from the four options as the final answer. We provide several student reasonings for the last question. Some of them may be correct and some incorrect. You can use the best correct arguments from these reasonings. Beware of wrong reasoning and do not repeat wrong reasoning.

**Question:** A 22-year-old male marathon runner presents to the office with the complaint of right-sided rib pain when he runs long distances. Physical examination reveals normal heart and lung findings and an exhalation dysfunction at ribs 4-5 on the right. Which of the following muscles or muscle groups will be most useful in correcting this dysfunction utilizing a direct method?
(A) anterior scalene (B) latissimus dorsi (C) pectoralis minor (D) quadratus lumborum
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. Among the options, only pectoralis minor muscle origins from the outer surfaces of the 3rd to 5th ribs.
**Answer:** (C)

**Question:** A 36-year-old male presents to the office with a 3-week history of low back pain. He denies any recent trauma but says that he climbs in and out of his truck numerous times a day for his job. Examination of the patient in the prone position reveals a deep sacral sulcus on the left, a posterior inferior lateral angle on the right, and a lumbosacral junction that springs freely on compression. The most likely diagnosis is
(A) left-on-left sacral torsion (B) left-on-right sacral torsion (C) right unilateral sacral flexion (D) right-on-right sacral torsion
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The deep sulcus on the left, a posterior ILA on the right, with a negative spring test suggests a right-on-right sacral torsion. All other options have a deep sulcus on the right.
**Answer:** (D)

**Question:** A 44-year-old man comes to the office because of a 3-day history of sore throat, nonproductive cough, runny nose, and frontal headache. He says the headache is worse in the morning and ibuprofen does provide some relief. He has not had shortness of breath. Medical history is unremarkable. He takes no medications other than the ibuprofen for pain. Vital signs are temperature 37.4°C (99.4°F), pulse 88/min, respirations 18/min, and blood pressure 120/84 mm Hg. Examination of the nares shows erythematous mucous membranes. Examination of the throat shows erythema and follicular lymphoid hyperplasia on the posterior oropharynx. There is no palpable cervical adenopathy. Lungs are clear to auscultation. Which of the following is the most likely cause of this patient's symptoms?
(A) Allergic rhinitis (B) Epstein-Barr virus (C) Mycoplasma pneumonia (D) Rhinovirus
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The symptoms, especially the headache, suggest that the most likely cause is Rhinovirus. Epstein-Barr virus will cause swollen lymph nodes but there is no palpable cervical adenopathy. Lungs are clear to auscultation suggests it's not Mycoplasma pneumonia.
**Answer:** (D)

**Question:** A previously healthy 32-year-old woman comes to the physician 8 months after her husband was killed in a car crash. Since that time, she has had a decreased appetite and difficulty falling asleep. She states that she is often sad and cries frequently. She has been rechecking the door lock five times before leaving her house and has to count exactly five pieces of toilet paper before she uses it. She says that she has always been a perfectionist but these urges and rituals are new. Pharmacotherapy should be targeted to which of the following neurotransmitters?
(A) Dopamine (B) Glutamate (C) Norepinephrine (D) Serotonin
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The patient feels sad and among the options, only Dopamine and Serotonin can help increase positive emotions. Serotonin also affects digestion and metabolism, which can help the patient's decreased appetite and sleep difficulty.
**Answer:** (D)

**Question:** A 42-year-old man comes to the office for preoperative evaluation prior to undergoing adrenalectomy scheduled in 2 weeks. One month ago, he received care in the emergency department for pain over his right flank following a motor vehicle collision. At that time, blood pressure was 160/100 mm Hg and CT scan of the abdomen showed an incidental 10-cm left adrenal mass. Results of laboratory studies, including complete blood count, serum electrolyte concentrations, and liver function tests, were within the reference ranges. The patient otherwise had been healthy and had never been told that he had elevated blood pressure. He takes no medications. A follow-up visit in the office 2 weeks ago disclosed elevated urinary normetanephrine and metanephrine and plasma aldosterone concentrations. The patient was referred to a surgeon, who recommended the adrenalectomy. Today, vital signs are temperature 36.6°C (97.9°F), pulse 100/min, respirations 14/min, and blood pressure 170/95 mm Hg. Physical examination discloses no significant findings. Initial preoperative preparation should include treatment with which of the following?
(A) Labetalol (B) A loading dose of potassium chloride (C) Nifedipine (D) Phenoxybenzamine
**Explanation:** Let's solve this step-by-step, referring to authoritative sources as needed. The symptoms and the adrenal mass suggested pheochromocytoma, and the blood pressure indicates hypertension. Phenoxybenzamine is used to treat hypertension caused by pheochromocytoma.
**Answer:** (D)

**Question**: A 67-year-old man with transitional cell carcinoma of the bladder comes to the physician because of a 2-day history of ringing sensation in his ear. He received this first course of neoadjuvant chemotherapy 1 week ago. Pure tone audiometry shows a sensorineural hearing loss of 45 dB. The expected beneficial effect of the drug that caused this patient's symptoms is most likely due to which of the following actions?

(A) Inhibition of proteasome (B) Hyperstabilization of microtubules (C) Generation of free radicals (D) Cross-linking of DNA

**Students' reasonings:**

**1 reasoning**: Cisplatin is a chemotherapy drug that causes sensorineural hearing loss. Cisplatin is an alkylating agent that causes cross-linking of DNA.

**2 reasoning**: Cisplatin is a platinum-based chemotherapy drug used to treat several types of cancer, including bladder cancer. Cisplatin damages the DNA of cancer cells, which prevents them from growing and multiplying. Cisplatin can cause several side effects, including hearing loss. Ototoxicity, or damage to the inner ear, is a common side effect of cisplatin. Cisplatin can cause both temporary and permanent hearing loss. Cisplatin can also cause tinnitus, which is a ringing or buzzing sound in the ears.

**3 reasoning**: Cisplatin is the most common cause of ototoxicity in cancer patients. Cisplatin is an alkylating agent that acts by cross-linking DNA.

**4 reasoning**: Cisplatin is a chemotherapeutic agent that can cause ototoxicity. Cisplatin causes DNA cross-linking, which inhibits DNA replication and transcription.

**5 reasoning**: Cisplatin is a platinum-based chemotherapy drug that is used to treat a variety of cancers, including bladder cancer. Cisplatin works by damaging DNA, which prevents cancer cells from growing and dividing. Cisplatin can cause a number of side effects, including hearing loss. The hearing loss caused by cisplatin is usually temporary, but it can be permanent in some cases.

**6 reasoning**: Cisplatin is a platinum-containing compound that is used to treat various types of cancer, including transitional cell carcinoma of the bladder. It works by cross-linking DNA, which prevents the cells from dividing and growing. Cisplatin can cause ototoxicity, which is damage to the inner ear that can cause hearing loss and tinnitus (ringing in the ears).

**7 reasoning**: Cisplatin is a chemotherapeutic agent which is used in the treatment of transitional cell carcinoma of the bladder. Cisplatin causes ototoxicity by generation of free radicals.

**8 reasoning**: Cisplatin is a platinum-containing compound that cross-links DNA and inhibits DNA replication. It is used in the treatment of many solid tumors, including transitional cell carcinoma of the bladder. Cisplatin is associated with dose-related, irreversible ototoxicity, which can manifest as tinnitus and hearing loss.

**9 reasoning**: Cisplatin is the drug that most commonly causes ototoxicity. Cisplatin is a platinum-containing compound that forms inter- and intra-strand DNA cross-links, thereby inhibiting DNA replication and transcription and producing cytotoxicity.

**10 reasoning**: Cisplatin is a platinum-containing compound that can cause ototoxicity. Cisplatin causes DNA cross-linking and cell death.

**11 reasoning**: Cisplatin is a platinum-based chemotherapy agent that is used to treat a variety of cancers, including bladder cancer. Cisplatin works by damaging the DNA of cancer cells, which prevents them from growing and dividing. Cisplatin can cause a number of side effects, including hearing loss, ringing in the ears (tinnitus), and kidney damage. Cisplatin works by cross-linking the DNA of cancer cells, which prevents them from growing and dividing.

**Explanation**:

Table A.7: PubMedQA (2019) few-shot prompt examples from Med-PaLM

**INSTRUCTIONS:** This is a multiple choice question about medical research. Determine the answer to the question based on the strength of the scientific evidence provided in the context. Valid answers are yes, no or maybe. Answer yes or no if the evidence in the context supports a definitive answer. Answer maybe if the evidence in the context does not support a definitive answer, such as when the context discusses both conditions where the answer is yes and conditions where the answer is no.

FEW_SHOT_TEMPLATE:
Instructions: {INSTRUCTIONS}
Context: {TRAIN_CONTEXT_1}
Question:{TRAIN_QUESTION_1}
Answer: The answer to the question given the context is {TRAIN_ANSWER_1}.

Instructions: {INSTRUCTIONS}
Context: {TRAIN_CONTEXT_2}
Question:{TRAIN_QUESTION_2}
Answer: The answer to the question given the context is {TRAIN_ANSWER_2}.

Instructions: {INSTRUCTIONS}
Context: {TRAIN_CONTEXT_3}
Question:{TRAIN_QUESTION_3}
Answer: The answer to the question given the context is {TRAIN_ANSWER_3}.

Instructions: {INSTRUCTIONS}
Context: {EVAL_CONTEXT}
Question:{EVAL_QUESTION}

Figure A.2: **Example of Correct Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the accurate response generated by Gemini.



Figure A.3: **Example of Correct Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the accurate response generated by Gemini.

Figure A.4: **Example of Correct Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the accurate response generated by Gemini.



Figure A.5: **Example of incorrect Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the incorrect response generated by Gemini.

44

Figure A.6: **Example of incorrect Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the incorrect response generated by Gemini.



Figure A.7: **Example of incorrect Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the incorrect response generated by Gemini.

Figure A.8: **Example of incorrect Gemini Output on Visual Question Answering Benchmark** This figure provides a randomly selected sample question from the VQA benchmark alongside the incorrect response generated by Gemini.

# Retrieval augmented text-to-SQL generation for epidemiological question answering using electronic health records

**Angelo Ziletti***
Bayer AG
angelo.ziletti@bayer.com

**Leonardo D'Ambrosi**
Bayer AG

## Abstract

Electronic health records (EHR) and claims data are rich sources of real-world data that reflect patient health status and healthcare utilization. Querying these databases to answer epidemiological questions is challenging due to the intricacy of medical terminology and the need for complex SQL queries. Here, we introduce an end-to-end methodology that combines text-to-SQL generation with retrieval augmented generation (RAG) to answer epidemiological questions using EHR and claims data. We show that our approach, which integrates a medical coding step into the text-to-SQL process, significantly improves the performance over simple prompting. Our findings indicate that although current language models are not yet sufficiently accurate for unsupervised use, RAG offers a promising direction for improving their capabilities, as shown in a realistic industry setting.

## 1 Introduction

Real-world data (RWD) are data routinely gathered from various sources that capture aspects of patient health status and the provision of health care. This encompasses electronic health records (EHR), medical claims data, disease registries, and emerging sources like digital health technologies. By investigating epidemiological quantities like patients' counts and demographics, disease incidence and prevalence, natural history of diseases, and treatment patterns in real-world clinical practice, researchers and healthcare organizations can identify for example target patient populations with unmet needs, discover unknown benefits of available drugs, evaluate potential for market entry, and estimate the potential enrolment of clinical trials.

**Problem Statement.** Addressing epidemiological questions using RWD databases is complex, as it requires not only an understanding of the data's characteristics, including biases, confounders, and

limitations, but also involves interpreting medical terminology across various ontologies, formulating precise SQL queries, executing these queries, and accurately synthesizing the results.

**Contributions.** With this paper, we present a straightforward and effective end-to-end approach to answer epidemiological questions based on data queried from EHR/Claims databases.

- We release a dataset of manually annotated question-SQL pairs designed for epidemiological research, and adhering to the widely-adopted Observational Medical Outcomes Partnership Common Data Model (OMOP-CDM) (OMOP-CDM, 2023).

- We integrate a medical coding step into the text-to-SQL process, enhancing data retrieval and clinical context comprehension.

- We show that retrieval augmented generation (RAG) significantly improves performance compared with static instruction prompting, as confirmed by extensive benchmarking with top-tier large language models (LLMs).

- We share our dataset, code, and prompts[1] to foster reproducibility and catalyse a community-driven effort towards advancing this research area.

The presented approach is currently deployed at Bayer in experimental mode. Epidemiologists and data analysts are using the system to explore and evaluate its capabilities, ensuring that its use is carefully monitored and supervised.

## 2 The Dataset

Our dataset was created through a manual curation process, engaging specialists in epidemiological

---

[1] https://github.com/Bayer-Group/text-to-sql-epi-ehr-naacl2024

| Quantity | Value |
|---|---|
| # of question/SQL pairs (all) | 306 |
| # of different tables used (all) | 13 |
| # of different columns used (all) | 44 |
| # logical conditions/query | 6.4 (6.7) |
| # nesting levels/query | 1.5 (1.1) |
| # tables/query | 2.7 (0.9) |
| # columns/query | 6.3 (4.7) |
| # medical entities/query | 2.0 (4.1) |
| Question length [char]/query | 91.7 (81.2) |
| SQL query length [char]/query | 796.4 (448.5) |

Table 1: Summary statistics of the dataset. For sample statistics, average and standard deviation (in brackets) are reported.

studies to contribute typical questions from their work. Despite its modest size, the dataset offers a realistic selection of epidemiological questions within industry practice, and exhibits a high degree of complexity. 53 samples require more than two level of nesting, and 19 more than three levels. Correctly answering questions often require multiple logical steps: selection of population(s) of interest, relationship between events within a specific time frame, aggregation statistics, and basic mathematical operations (e.g., ratios). The dataset features questions in their natural, free-form language and it is augmented with two paraphrased versions per question-SQL pair, increasing volume while also offering validated labels for retrieval algorithms. Statistics on the dataset are shown in Table 1. Due to budget limits, we will use one version per question for subsequent evaluations.

**Applicability across RWD databases**. To address the challenge of data retrieval variability across databases with differing data models, we leverage the OMOP-CDM. This model, underpinned by standardized vocabularies (Reich et al., 2024), harmonizes observational healthcare data and it is widely recognized as the standard for RWD analysis, with data from over 2.1 billion patient records across 34 countries (Voss et al., 2023; Reich et al., 2024).

## 3 Methods

Our methodology, outlined in Fig. 1, employs LLM prompting to translate natural language questions into SQL queries. It advances EHR text-to-SQL methods beyond the constraints of exact or string-based matching to fully encompass the semantic complexities of clinical terminology (Wang et al., 2020; Lee et al., 2022). To achieve this, we introduce a step where an LLM generates SQL

with placeholders for medical entities (e.g., [condition@disphagia] in Fig. 1d), which are then mapped to precise clinical ontology terms (Sec. 4.1, Fig. 1d-e). This yields executable queries that accurately retrieve database information. Building on the success of RAG in enhancing LLMs for complex NLP tasks (Lewis et al., 2020), we use our dataset (Sec. 2) as an external knowledge base. Relevant question-SQL pairs are extracted and incorporated into the prompt, refining SQL generation. The completed SQL queries, embedded with medical codes, are run on an OMOP CDM-compliant database (Fig. 1f) to facilitate data retrieval. If needed, an answer can be articulated from the retrieved data through further LLM prompting (Fig. 1g).

## 4 Evaluation

### 4.1 Experimental setup

**Large language models.** We employ several leading LLMs as of February 2024: OpenAI's GPT-3.5 Turbo (Brown et al., 2020) and GPT-4 Turbo (OpenAI, 2023), Google's GeminiPro 1.0 (Gemini Team, 2023), Anthropic's Claude 2.1 (Anthropic AI, 2023), and Mistral AI's Mixtral 8x7B and Mixtral Medium (Mistral AI, 2023), with Mixtral 8x7B being the only open-source model (Jiang et al., 2024). We use one simple and one advanced prompt. The simple prompt provides essential instructions for creating queries that adhere to the conventions of the pipeline (Fig. 1). The advanced prompt adds detailed directives on concept IDs, race analysis, geographical analysis, date filters, column naming, patient count, age calculation, and additional instructions on SQL query validity review. Following Pourreza and Rafiei (2023), we allow LLMs up to three attempts to self-correct non-executable SQL queries using the compiler's error feedback.

**Retrieval augmented generation.** For similarity computation in RAG, we apply entity masking to substitute medical entities with generic labels (e.g., <DRUG>). We utilize the BGE-LARGE-EN-V1.5 embedding model from Hugging Face (Wolf et al., 2020), which has been fine-tuned for retrieval augmentation of LLMs (Zhang et al., 2023). We opt for masked question selection rather than utilizing the query because it eliminates the need for an initial LLM call to generate SQL for retrieval, while maintaining a comparable accuracy (Gao et al., 2023).

Figure 1: From a question in natural language to an answer in natural language using electronic health record or claims databases: end-to-end workflow.

**Medical coding.** LLMs extract medical entities and integrate them into SQL as placeholders (Fig. 1d), effectively recasting the medical coding task into medical entity normalization (Portelli et al., 2022; Ziletti et al., 2022; Zhang et al., 2022; Limsopatham and Collier, 2016). To perform entity normalization, we first compute the cosine similarity of each entity's SapBERT embeddings (Liu et al., 2021) with SNOMED ontology terms, and select the top-50 matches. Then, similarly to Yang et al. (2022), we prompt GPT-4 Turbo to verify whether a given code should be assigned to the input entity, refining the list.

**Database and evaluation.** The evaluation data reported are obtained querying the DE-SynPUF dataset (SynPUF, 2010), which is a synthetic dataset that emulates the structure of actual claims data. It includes 6.8 million beneficiary records, 112 million claims records, and 111 million prescription drug events records (Gonzales et al., 2023). The same analysis could be applied to any database conforming to the OMOP-CDM, thus potentially allowing access to 2.1 billion patient records (Reich et al., 2024). For evaluation, we manually developed a dataset of question-SQL pairs, as detailed in Sec. 2. These are then executed against the DE-SynPUF dataset, and the retrieved data from both reference and generated queries are compared to assess performance. This process reflects the practical use of SQL queries on healthcare databases. To ensure a realistic eval-

uation setup, the actual question being evaluated is removed from the RAG procedure. A generated query is marked as correct if it retrieves data enabling an answer that aligns with the reference query's answer (within a 10% tolerance), and incorrect otherwise. The tolerance compensates for variations from GPT-4 Turbo-based medical coding, maintaining the focus on text-to-SQL evaluation accuracy.

### 4.2 Experimental results

Results are shown in Table 2, and outlined below.
**Enhanced performance with detailed prompting.**
Advanced prompting typically increases execution scores across models (except GPT-3.5 Turbo), but its impact on accuracy varies: Claude 2.1, Mistral-m, and GPT-4 Turbo show marked accuracy improvements with the advanced prompt, whereas Mixtral, GeminiPro, and GPT-3.5 Turbo see no such gains, suggesting that the additional details in the prompt may not benefit smaller or less sophisticated models. Overall performance is quite poor with either prompting methods.
**Performance gains with contextual information.**
The inclusion of relevant examples via RAG significantly and consistently improves performance (Table 2, cf. RAG-top1/2/5 vs Prompt(advanced)). Notably, Mistral-m and GPT-4 Turbo exhibit marked improvements, suggesting they may possess a more advanced few-shot learning ability relative to the other models. Models outperform zero-shot

| | Mixtral | | GeminiPro | | Claude 2.1 | | Mistral-m | | GPT-3.5 Turbo | | GPT-4 Turbo | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc | Exec | Acc | Exec | Acc | Exec | Acc | Exec | Acc | Exec | Acc | Exec |
| Prompt (simple) | 2.0 | 7.8 | 6.9 | 29.4 | 20.6 | 53.5 | 8.8 | 32.4 | 20.2 | 67.0 | **28.4** | **77.5** |
| Prompt (advanced) | 2.9 | 18.6 | 6.9 | 34.7 | 25.5 | 78.4 | 17.6 | 44.1 | 15.8 | 63.4 | **38.2** | **91.2** |
| RAG-random1 | 19.6 | 46.1 | 11.8 | 35.3 | 33.3 | 76.5 | 38.2 | 68.3 | 29.0 | 84.0 | **50.0** | **97.1** |
| RAG-top1 | 33.3 | 52.0 | 38.2 | 59.8 | 29.4 | 73.5 | 50.0 | 69.6 | 59.8 | 90.2 | **72.5** | **97.1** |
| RAG-top2 | 20.6 | 40.2 | 37.3 | 56.9 | 38.6 | 75.2 | 46.1 | 73.5 | 61.8 | 94.1 | **77.5** | **98.0** |
| RAG-top5 | 22.5 | 44.1 | 35.0 | 62.0 | 34.3 | 71.6 | 51.0 | 73.5 | 52.0 | 95.1 | **77.5** | **97.1** |
| RAG-top1-oracle | 52.0 | 62.7 | 67.6 | 73.5 | 58.8 | 83.3 | 56.9 | 74.5 | **91.1** | **99.0** | 82.8 | 95.0 |

Table 2: Comparative evaluation of LLMs' performance on text-to-SQL generation for epidemiological question answering. Accuracy (Acc) and executability (Exec) percentages are presented across different models and prompting conditions. Best results are in bold, while second best are underlined. RAG-top1/2/5 indicates the use of the top 1, 2, or 5 most similar questions to augment generation. RAG-random1 and RAG-top1-oracle scenarios provide models with a random dataset sample and the correct SQL query, respectively, for context.

prompting also when given a random dataset sample (RAG-random1), indicating that exposure to dataset structure and domain-specific language is helpful, even without query-specific context.

**Diminishing returns with increased context.** Providing a single example (RAG-top-1) leads to substantial improvements in performance, but adding more top results (RAG-top2 and RAG-top5) does not result in a similar increase. Some models exhibit a performance peak or a minor decline with additional context, indicating a limit to the beneficial amount of context.

**Superiority of GPT-4 Turbo.** GPT-4 Turbo is the best model overall by a large margin, followed by GPT-3.5 Turbo. Mistral-m outperforms both Claude 2.1 and GeminiPro. The open-source Mixtral model lags behind proprietary models in both accuracy and executability across all scenarios.

**Model-specific approach to oracle context.** In the RAG-top1-oracle scenario, where the prompt includes the correct SQL query, GPT-3.5 Turbo unexpectedly surpasses GPT-4 Turbo by closely mirroring the provided context, favouring direct replication. In contrast, GPT-4 Turbo and other models take a "deliberative" approach, often modifying the input, which, while useful for complex reasoning, hinders tasks that require exact copying.

## 5 Related Work

**Text-to-SQL datasets for EHRs.** The MIMIC-SQL dataset (Wang et al., 2020) comprises 10 000 template-generated questions for the MIMIC-III (Johnson et al., 2016) database. It contains both question designed to retrieve patient-specific information, and questions on patients counts with logical and basic mathematical operations. Tarbell et al. (2023) noted limited diversity in MIMIC-

SQL's queries, possibly affecting its utility for testing text-to-SQL model generalizability. emrKBQA (Raghavan et al., 2021) contains 1 million patient-specific questions, also based on MIMIC-III. EHRSQL(Lee et al., 2022) is a dataset created by extracting templates from clinical questions posed by hospital staff, which are then used to generate a comprehensive set of queries for MIMIC-III and eICU (Pollard et al., 2018). It relies on an earlier, less performing text-to-text model for query generation (Raffel et al., 2020). All these datasets do not adhere with OMOP-CDM, and they opt for direct string matching for concept retrieval. The closest dataset to ours is the OMOP query library (OHDSI, 2019; OMOP-CDM-Query-Library, 2019), which is a collection of queries in OMOP-CDM. We adapted and included fifteen SQL queries from this library pertinent to epidemiological research into our dataset. Park et al. (2023) use rule-based methods and GPT-4 to translate clinical trial eligibility criteria into SQL queries for OMOP-CDM.

**Text-to-SQL with LLMs and in-domain demonstrations.** Prompting LLMs has proven effective, often outperforming specialized fine-tuned models in text-to-SQL task (Pourreza and Rafiei, 2023). Both in-domain (Chang and Fosler-Lussier, 2023a) and out-of-domain (Chang and Fosler-Lussier, 2023b) demonstrations improve LLMs' performance. Gao et al. (2023) explores retrieval scenarios for in-domain demonstration selection. To the best of our knowledge, the exploration of these text-to-SQL methods within EHR (or biomedical) research has not yet extended to small datasets that are critical for industry applications.

# 6 Conclusion

In this work, we presented the task of answering epidemiological questions using RWD. We demonstrated that RAG is effective in improving performance on all tested scenarios. Our study extends the demonstrated efficacy of RAG from general text-to-SQL benchmarks (Gao et al., 2023; Chang and Fosler-Lussier, 2023b) to include to small, domain-specific biomedical datasets, underlining its utility in data-scarce industry settings. The primary limitation is the dataset's limited size and specialized focus on epidemiological questions, suggesting further research should broaden its scope and scale.

# 7 Acknowledgments

# References

Anthropic AI. 2023. Model card and evaluations for claude models. `https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf`. Accessed: February 15, 2024.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Shuaichen Chang and Eric Fosler-Lussier. 2023a. How to prompt llms for text-to-sql: A study in zero-shot, single-domain, and cross-domain settings.

Shuaichen Chang and Eric Fosler-Lussier. 2023b. Selective demonstrations for cross-domain text-to-SQL. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14174–14189, Singapore. Association for Computational Linguistics.

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.

Gemini Team. 2023. Gemini: A family of highly capable multimodal models. Technical report, Google. Accessed: February 15, 2024.

Aldren Gonzales, Guruprabha Guruswamy, and Scott R. Smith. 2023. Synthetic data in health care: A narrative review. *PLOS Digital Health*, 2(1):1–16.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts.

Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Liwei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3(1):160035.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. In *Advances in Neural Information Processing Systems*, volume 35, pages 15589–15601. Curran Associates, Inc.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA. Curran Associates Inc.

Nut Limsopatham and Nigel Collier. 2016. Normalising medical concepts in social media texts by learning semantic representation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1014–1023, Berlin, Germany. Association for Computational Linguistics.

Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. 2021. Self-alignment pretraining for biomedical entity representations. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4228–4238, Online. Association for Computational Linguistics.

Mistral AI. 2023. Generative endopoints of mistral ai. https://docs.mistral.ai/platform/endpoints/. Accessed: February 15, 2024.

OHDSI. 2019. *The Book of OHDSI*, 1 edition. Observational Health Data Sciences and Informatics, Seoul, Korea. Accessed: 2021-03-30.

OMOP-CDM. 2023. Omop cdm common data model. https://ohdsi.github.io/CommonDataModel/. Accessed: January 23, 2024.

OMOP-CDM-Query-Library. 2019. Omop cdm query library. https://github.com/OHDSI/QueryLibrary. Accessed: January 23, 2024.

OpenAI. 2023. Gpt-4 technical report.

Jimyung Park, Yilu Fang, and Chunhua Weng. 2023. Criteria2Query 3.0 Powered by Generative Large Language Models. Observational Health Data Sciences and Informatics (OHDSI). https://www.ohdsi.org/wp-content/uploads/2023/10/423-Park-BriefReport.pdf.

Tom J. Pollard, Alistair E. W. Johnson, Jesse D. Raffa, Leo A. Celi, Roger G. Mark, and Omar Badawi. 2018. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific Data*, 5(1):180178.

Beatrice Portelli, Simone Scaboro, Enrico Santus, Hooman Sedghamiz, Emmanuele Chersoni, and Giuseppe Serra. 2022. Generalizing over long tail concepts for medical term normalization. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8580–8591, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

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*.

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. *J. Mach. Learn. Res.*, 21(1).

Preethi Raghavan, Jennifer J Liang, Diwakar Mahajan, Rachita Chandra, and Peter Szolovits. 2021. emrKBQA: A clinical knowledge-base question answering dataset. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 64–73, Online. Association for Computational Linguistics.

Christian Reich, Anna Ostropolets, Patrick Ryan, Peter Rijnbeek, Martijn Schuemie, Alexander Davydov, Dmitry Dymshyts, and George Hripcsak. 2024. OHDSI Standardized Vocabularies—a large-scale centralized reference ontology for international data harmonization. *Journal of the American Medical Informatics Association*, page ocad247.

SynPUF. 2010. Medicare claims synthetic public use files (synpufs). https://www.cms.gov/data-research/statistics-trends-and-reports/medicare-claims-synthetic-public-use-files. Accessed: January 23, 2024.

Richard Tarbell, Kim-Kwang Raymond Choo, Glenn Dietrich, and Anthony Rios. 2023. Towards understanding the generalization of medical text-to-sql models and datasets. *AMIA ... Annual Symposium proceedings. AMIA Symposium*, 2023:669–678.

Erica A Voss, Clair Blacketer, Sebastiaan van Sandijk, Maxim Moinat, Michael Kallfelz, Michel van Speybroeck, Daniel Prieto-Alhambra, Martijn J Schuemie, and Peter R Rijnbeek. 2023. European health data and evidence network—learnings from building out a standardized international health data network. *Journal of the American Medical Informatics Association*, 31(1):209–219.

Ping Wang, Tian Shi, and Chandan K. Reddy. 2020. Text-to-sql generation for question answering on electronic medical records. In *Proceedings of The Web Conference 2020*, WWW '20, page 350–361, New York, NY, USA. Association for Computing Machinery.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Zhichao Yang, Sunjae Kwon, Zonghai Yao, and Hongfeng Yu. 2022. Multi-label few-shot icd coding as autoregressive generation with prompt. *Proceedings of the AAAI Conference on Artificial Intelligence. AAAI Conference on Artificial Intelligence*, 37 4:5366–5374.

Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. 2023. Retrieve anything to augment large language models.

Sheng Zhang, Hao Cheng, Shikhar Vashishth, Cliff Wong, Jinfeng Xiao, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2022. Knowledge-rich self-supervision for biomedical entity linking. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 868–880, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Angelo Ziletti, Alan Akbik, Christoph Berns, Thomas Herold, Marion Legler, and Martina Viell. 2022. Medical coding with biomedical transformer ensembles and zero/few-shot learning. In *Proceedings of*

*the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, pages 176–187, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

# ClinicalMamba: A Generative Clinical Language Model on Longitudinal Clinical Notes

**Zhichao Yang[1], Avijit Mitra[1], Sunjae Kwon[1], Hong Yu[1,2]**

[1] College of Information and Computer Sciences, University of Massachusetts Amherst
[2] Department of Computer Science, University of Massachusetts Lowell
{zhichaoyang,avijitmitra,sunjaekwon}@umass.edu  hong_yu@uml.edu

## Abstract

The advancement of natural language processing (NLP) systems in healthcare hinges on language models' ability to interpret the intricate information contained within clinical notes. This process often requires integrating information from various time points in a patient's medical history. However, most earlier clinical language models were pretrained with a context length limited to roughly one clinical document. In this study, We introduce ClinicalMamba, a specialized version of the Mamba language model, pretrained on a vast corpus of longitudinal clinical notes to address the unique linguistic characteristics and information processing needs of the medical domain. ClinicalMamba models, with 130 million and 2.8 billion parameters, demonstrate superior performance in modeling clinical language across extended text lengths compared to Mamba and other clinical models based on longformer and Llama. With few-shot learning, ClinicalMamba achieves notable benchmarks in speed and performance, outperforming existing clinical language models and large language models like GPT-4 in longitudinal clinical tasks.

## 1 Introduction

Clinical narratives, such as patient histories, consultation notes, and discharge summaries, contain detailed and complex information that extends over long text sequences (Wu et al., 2019). To fully understand a patient's condition, treatments, and outcomes, NLP systems need to integrate information from various parts of these narratives, which often requires understanding the context provided in those long form text (Blumenthal, 2010).

Understanding the sequence of health events is crucial for diagnoses, treatment plans, and patient monitoring (Wang et al., 2024; Yang et al., 2023; Kraljevic et al., 2023; Eva, 2005). This often involves putting together information from different time points within a patient's health history (Gao et al., 2024). Long context enables NLP systems to perform temporal reasoning by tracking events over time longitudinally, which is essential for tasks like predicting disease progression or extracting medical relation (Chen et al., 2023; Jia et al., 2019; Wiegreffe et al., 2019).

It becomes imperative to design models for the need for processing longer texts (Parmar et al., 2023; Tay et al., 2020). Prior studies have introduced Mamba (Gu and Dao, 2023), a selective state space model, that selects and compresses all necessary information into latent space from context, and achieves linear-time efficiency with context length. While these advancements have been primarily directed towards processing general domain text, the unique linguistic features of clinical narratives differ significantly from general domain (Lehman et al., 2023), motivating us to develop specialized Mamba models in the clinical domain.

In this work, we build and publicly release ClinicalMamba - a Mamba model pretrained on longitudinal clinical notes. Furthermore, we demonstrate that ClinicalMamba outperforms multiple language models on longitudinal clinical NLP tasks. In particular, our contributions are as follows:

- We publicly release ClinicalMamba with 130m and 2.8b parameters trained on MIMIC-III (Johnson et al., 2016). [1]

- Through distributed training, ClinicalMamba-2.8b model was pretrained in under 60 hours on 4 A100 GPUs and it is the first clinical autoregressive language model with a 16k maximum token length.

- Through few-shot prompt-based finetuning, we demonstrate both ClinicalMamba outperforms original Mamba, GPT4, and other existing clinical long context language models

---

[1] https://github.com/whaleloops/ClinicalMamba

Figure 1: Perplexity of different generative language models on MIMIC-III when evaluated at various preceding context lengths (1k, 4k, and 16k tokens). The X-axis is in the log scale. The subfigure is a zoom-out plot with perplexity ranges 0-100. Experiment settings and detailed results are in section 5.

on well-established long context clinical information extractions tasks: cohort selection for clinical trial and international classification of diseases (ICD) coding.

## 2 Related Work

### 2.1 Pretraining clinical narratives

The rapid expansion of the utilization of electronic health records (EHRs) into the healthcare landscape underscores an urgent need for a clinical language model (Kang et al., 2019). Previous work, such as Alsentzer et al. (2019) on Clinical BERT embeddings and Huang et al. (2019) with ClinicalBERT and Lewis et al. (2020) with ClinicalRoberta, adapted general-purpose language models to the clinical domain to enhance performance on clinical tasks. These models have been pivotal in demonstrating the effectiveness of adapting general-purpose NLP tools to the intricacies of clinical text. Similarly, the creation of GatorTron (Yang et al., 2022a) scales up clinical language models to billions of parameters, while NYUTron (Jiang et al., 2023) harness billions of unstructured data found in EHRs. Both underscores the potential of domain-adapted language models to advance clinical NLP by improving performance across various tasks such as concept extraction and outcome prediction (Yang et al., 2022a; Jiang et al., 2023).

To handle complex and nuanced tasks, recent studies investigated training generative models with prompt (Kweon et al., 2023; Peng et al., 2023;

Wang et al., 2023; Lu et al., 2022; Wang and Sun, 2022). These models not only excel in classification but also in generating clinically relevant text that can be indistinguishable from human-written notes. Most previous methods focus on pretraining transformer models with a context window less than 2k tokens. However, we pretrained a selective state space model with a context window of 16k tokens, which includes more than 98% of the visits in MIMIC-III.

### 2.2 Clinical information extraction on long document

Handling long texts in clinical NLP has always been challenging. Traditional methods of information extraction tackle this by marking specific locations within the sentence, but such labeling is not always available, and hiring annotators can be costly (Fu et al., 2020; Mitra et al., 2023). Recent advancements in document information extraction involve pairing labels with documents (Kwon et al., 2022; Deshpande et al., 2024). However, BERT and Roberta struggle with processing these lengthy documents directly.

To address this, prior research introduced *Hierarchical-ClinicalRoberta*, which involves breaking down long documents into shorter segments of 512 tokens, applying ClinicalRoberta to each segment to obtain embeddings, and then using additional layers to leverage these embeddings for label classification (Huang et al., 2022; Zhang and Jankowski, 2022). However, this method combines information from each segment only at the final layer, which can hinder performance when training data is limited.

To mitigate this issue, *ClinicalLongformer* is designed to efficiently process longer context length by employing a self-attention mechanism across all layers, which is key to its proficiency in managing dense information exchanges within a specified contextual range (Li et al., 2022; Ji et al., 2023). This mechanism, while powerful, is limited by its focus on a predetermined window of text, restricting its scope to what falls within this window.

To overcome these limitations, the Mamba model emerges as a revolutionary approach. Mamba employs a selective state space model strategy to meticulously choose critical data for incorporation into its state (Gu and Dao, 2023), thereby, enhancing its capability to manage information beyond the conventional self-attention window. In

general domain language modeling, Mamba surpasses Transformers of equivalent size in task performance and speed.

## 3 Methods

### 3.1 Pretraining

We gather 82,178 hospital visits along with their deidentified free-text clinical notes (2,083,180) from 46,520 patients in MIMIC-III (Johnson et al., 2016). Rather than breaking down the notes into chunks of 512 tokens to act as individual data instances, we aggregate all notes related to a visit longitudinally. The distribution of token counts per data instance is detailed in Table A.1. For information on our text pre-processing methods, please refer to section A.1.

Following previous works (Li et al., 2019; Lee et al., 2019), we continue to pretrain Mamba using MIMIC-III clinical notes with the causal language modeling objective. This pretraining process utilizes 4 Nvidia A100-80GB GPUs. It's important to note that some of our downstream evaluation tasks utilize a small subset (6,049) of hospital visits from MIMIC-III, so we exclude them from the pretaining data. A comprehensive training recipe is available in section A.2.

### 3.2 Prompt based fine-tuning

We leverage the inherent capabilities of pre-trained language models by introducing a novel fine-tuning strategy that aligns with the specific demands of few-shot learning in clinical NLP. Recognizing the limitations of traditional fine-tuning methods when applied to clinical NLP tasks with limited labeled data, we adapted a prompt-based fine-tuning mechanism following previous works (Gao et al., 2021; Yang et al., 2022b; Taylor et al., 2023). Specifically, we first identify a set of representative prompts that encapsulate key aspects of the clinical tasks, such as the patient's alcohol consumption. These prompts are then appended after each input clinical note and incorporated into the fine-tuning phase, where the language model learns to associate them with label tokens (Yes/No) based on a limited dataset. The generated label tokens are then mapped to label space.

As shown in Figure 2, we transfer the downstream information extraction task into label token generation, which is similar to next token prediction during pretraining. We prompt based fine-tuned mamba and other baseline models unless



Figure 2: Illustration of Prompt-based fine-tuning.

otherwise specified on the following tasks.

### 3.3 Fine-tuning tasks

Cohort selection for clinical trial addresses the challenge of interpreting unstructured clinical narratives to streamline the patients selection process (Wornow et al., 2024; Jin et al., 2023; Wong et al., 2023). It aims to classify patients based on whether they meet 13 specific eligibility criteria, such as the usage of aspirin to prevent myocardial infarction, excessive alcohol consumption, and HbA1c values between 6.5 and 9.5%, among others. The input contains multiple clinical notes with a total length of 4924 tokens on average. This dataset was released as part of n2c2 challenge (track 1) in 2018 (Stubbs et al., 2019).

ICD coding interprets complex clinical narratives, translating them into standardized codes that facilitate accurate billing, statistical analysis, and healthcare management. It aims to extract patient's disease and procedure codes from clinical text. We followed general instructions from Mullenbach et al. (2018) in building this task, but instead of using a single discharge summary as input, we used all previous discharge summaries and assigned ICD code descriptions from previous visits. We further filtered 50 infrequent codes as Code-rare and 50 frequent codes as Code-common following Yang et al. (2022b). The average length is 4,223 and 7,062 tokens respectively. Detailed dataset statistics are shown in Table A.1.

## 4 Experiments

For each fine-tuning task, we measured the micro precision, micro recall, micro F1 scores, and micro receiver operating characteristic/area under the curve (ROCAUC). We compared our model with the following baselines:

**GPT-4** is a large language model designed to understand and generate human-like text based on the input it receives. We applied zero-shot prompting

to each downstream task, using ACAN and original prompt introduced in Wornow et al. (2024). GPT-4 (version 2023-12-01-preview) was accessed securely through the Azure OpenAI API. We opt out of human review of the data by signing the Additional Use Case Form[2]. We set the sampling temperature for decoding to 0.1.

**Asclepius-R** (Kweon et al., 2023) is a clinical generative language model trained on MIMIC-III discharge summaries and corresponding instruction-answer pairs. It has 7 billion parameters with a maximum input of 4096 tokens.

**ClinicalLlama2** (CLlama2) is similar to Asclepius-R, but it was trained on all types of MIMIC-III note with the same computation budget as ClinicalMamba-2.8b (60 hours on 4 A100). It has 7 billion parameters with 4096 max context length.

**ClinicalLongformer** (CLongformer) (Li et al., 2022) is a clinical knowledge enriched version of Longformer that was further pretrained using MIMIC-III clinical notes. It has 149 million parameters with a maximum input of 4096 tokens. We only used local attention and does not apply global attention for computation efficiency.

**Hierachical ClinicalRoberta** (Hi-CRoberta) (Huang et al., 2022), utilizes multiple embedding from clinical Roberta (Lewis et al., 2020). It first segment clinical notes into chunks of 512 tokens to obtain their embeddings, embeddings are then pooled by concatenation and finally a linear classification head during downstream task. It has 110 million parameters with a max of 16384 tokens. We did not apply prompt based fine-tuning as this is not a generative model.

**MultiResCNN** (Li and Yu, 2020) encode free text with Multi-Filter ResidualCNN, and applied label code attention mechanism to enable each ICD code to attend different parts of the document.

## 5 Results & Discussions

In this section, we will first compare the model's language modeling ability on MIMIC-III clinical notes. We will then describe the evaluation on different clinical information extraction tasks. Finally, we will share a case study, which illustrates that our ClinicalMamba can recall patient information with long history.

ClinicalMamba stands as the sole model capable of handling clinical notes of up to 16k tokens.

| Model | Prec | Recall | F1 | AUC |
|---|---|---|---|---|
| **C**Llama2 | 70.0 | 79.1 | 77.7 | 84.3 |
| Hi-**C**Roberta | 72.4 | 82.6 | 79.2 | 88.1 |
| **C**Longformer | 69.7 | 78.6 | 76.1 | 83.5 |
| GPT-4 | 88.1 | 79.9 | 84.8 | - |
| Mamba-130m | 75.4 | 80.2 | 77.7 | 85.7 |
| **C**Mamba-130m | 79.0 | 86.2 | 82.2 | 91.8 |
| **C**Mamba-2.8b | **88.6** | **89.5** | **88.8** | **95.7** |

Table 1: Results on cohort selection task, where **C** is model pretrained in clinical domain.

As demonstrated in Figure 1, the perplexity for ClinicalMamba-2.8b decreased from 3.11 to 2.61 as the context length expanded from 1k to 16k tokens during inference. This is in contrast to the performance of prior clinical autoregressive language models, where perplexity levels rose with increased context lengths. For instance, with ClinicalLlama-7b, perplexity escalated from 2.82 to 94.02 as the context length grew from 4k to 16k. This limitation arises because these models were trained on contexts not exceeding 4k, impairing their accuracy for next token prediction when given previous contexts beyond 4k.

In the domain of extracting information from longitudinal clinical records, ClinicalMamba demonstrates superior performance compared to Mamba. ClinicalMamba achieved ROCAUC scores of 91.8, 42.3, and 94.2 on Cohort selection, Code-rare, and Code-common, while Mamba obtained ROCAUC scores of 85.7, 37.8, and 92.8 respectively. ClinicalMamba also outperformed previous long-range clinical language models with similar number of parameters. ClinicalMamba significantly outperformed Hierachical-ClinicalRoberta and ClinicalLongformer by relatively 52.7% and 19.1% on ROCAUC respectively. This is particularly notable in the Code-rare task with limited training data (5 shots), where ClinicalMamba attained an AUC of 91.1, compared to 77.1 of Hierachical-ClinicalRoberta and 80.5 of ClinicalLongformer.

Surprisingly, ClinicalMamba-2.8b also outperformed zero-shot GPT-4, achieving F1 scores of 88.8, 56.6, and 73.6 on Cohort selection, Code-rare, and Code-common tasks, whereas GPT-4 obtained a F1 score of 84.8, 33.3, and 68.2 respectively.

We also present a case with long history attached in the supplementary material. Both ClinicalMamba and ClinicalLongformer models identi-

| Model | Code-rare | | | | Code-common | | | |
|---|---|---|---|---|---|---|---|---|
| | Prec | Recall | F1 | AUC | Prec | Recall | F1 | AUC |
| MultiResCNN | 20.34 | 2.07 | 5.19 | 47.20 | 70.50 | 60.78 | 66.24 | 92.04 |
| Hi-**C**Roberta | 46.19 | 10.96 | 16.74 | 77.11 | 73.76 | 65.01 | 69.23 | 93.14 |
| **C**Longformer | 50.27 | 17.81 | 28.69 | 80.52 | **78.42** | 64.97 | 71.14 | 94.24 |
| GPT-4 | 30.91 | 36.12 | 33.29 | - | 72.48 | 62.28 | 68.19 | - |
| Mamba-130m | 57.75 | 28.08 | 37.79 | 84.80 | 73.71 | 62.87 | 68.94 | 92.75 |
| **C**Mamba-130m | 70.97 | 30.14 | 42.31 | 91.08 | 76.82 | 68.03 | **74.34** | 94.23 |
| **C**Mamba-2.8b | **75.28** | **45.89** | **56.51** | **92.75** | 75.53 | **72.12** | 73.64 | **94.54** |

Table 2: Results on ICD coding task, where **C** indicates model pretrained in clinical domain.

fied that this case met the criteria for history of advanced cardiovascular disease, and were tasked to interpret the prediction using SHAP (Lundberg and Lee, 2017). ClinicalMamba identified sentence "FINAL DIAGNOSIS: Acute MI" with highest SHAP value. This sentence locates at the first note among the patient history. In contrast, ClinicalLongformer identified a sentence, the patient "will be admitted in stable condition for further evaluation and rule out for myocardial infarction", at the last note. ClinicalLongformer also misinterpreted the negation in this sentence. In this case, ClinicalMamba remembered the correct long-distant sentence that leads to the correct answer, while ClinicalLongformer preferred wrong short-distant sentence to support its prediction.

## 6 Conclusion

In this study, we developed and released Mamba models pretrained on a large collection of clinical notes. Our findings demonstrate the superior performance of our ClinicalMamba in extracting information from long text documents compared to other models. We strongly believe that clinical NLP researchers can benefit from such long-context generative language models that alleviates the need of a substantial computational power, without any performance trade-off. Building on the groundwork laid by this study, future endeavors can further refine and expand the capabilities of ClinicalMamba.

## Limitations

This work has several notable limitations. First, we do not experiment with more recent parameter-efficient fine-tuning strategies such as soft prompting (Lester et al., 2021) and Low-Rank Adaptation (LoRa) (Hu et al., 2021). This potentially

undermined ClinicalMamba on downstream tasks. Second, our adaptation of the Mamba framework was restricted solely to textual data documented during visits. EHRs are rich with multifaceted information, including but not limited to radiology images taken at different times and Electrocardiogram waveforms that span various periods. Future research could develop a multimodal Mamba framework to leverage all other modalities. Third, the MIMIC-III dataset, which serves as the foundation of our study, only includes notes from the intensive care unit of a single hospital within the United States. This limits the generalizability of our findings, as care practices vary significantly across different institutions and countries. We did not pretrain on MIMIC-IV because it only has a limited number of notes (and also limited type: discharge summary and radiology report) per visit. Lastly, the linguistic scope of the MIMIC dataset is limited to English, which presents a barrier to understanding and applying our findings in non-English speaking contexts. Addressing these limitations could substantially broaden the applicability and relevance of our work in future endeavors.

## Ethics Statement

In this research, we gained authorized access to the MIMIC and N2C2 dataset and used de-identified clinical notes following their license agreement and HIPAA regulations. When language models are trained on extensive clinical text, they can inherit biases within the data. For instance, they might prefer inquiries concerning smoking habits or link specific medical conditions to certain demographic groups. These biases could be mitigated by enhancing model alignment with each patient's background.

# References

Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

David Blumenthal. 2010. Launching hitech. *The New England journal of medicine*, 362 5:382–5.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Zekai Chen, Mariann Micsinai Balan, and Kevin Brown. 2023. Language models are few-shot learners for prognostic prediction. *ArXiv*, abs/2302.12692.

Vijeta Deshpande, Minhwa Lee, Zonghai Yao, Zihao Zhang, Jason Brian Gibbons, and Hong Yu. 2024. Localtweets to localfealth: A mental health surveillance framework based on twitter data.

K. Eva. 2005. What every teacher needs to know about clinical reasoning. *Medical Education*, 39.

Sunyang Fu, David C. Chen, Huan He, Sijia Liu, Sungrim Moon, Kevin J. Peterson, Feichen Shen, Liwei Wang, Yanshan Wang, Andrew Wen, Yiqing Zhao, Sunghwan Sohn, and Hongfang Liu. 2020. Clinical concept extraction: A methodology review. *Journal of biomedical informatics*, page 103526.

Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and Marinka Zitnik. 2024. Units: Building a unified time series model.

Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3816–3830, Online. Association for Computational Linguistics.

Albert Gu and Tri Dao. 2023. Mamba: Linear-time sequence modeling with selective state spaces. *ArXiv*, abs/2312.00752.

J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685.

Chao-Wei Huang, Shang-Chi Tsai, and Yun-Nung Chen. 2022. PLM-ICD: Automatic ICD coding with pretrained language models. In *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 10–20, Seattle, WA. Association for Computational Linguistics.

Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. 2019. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *ArXiv*, abs/1904.05342.

Shaoxiong Ji, Tianlin Zhang, Kailai Yang, Sophia Ananiadou, E. Cambria, and Jörg Tiedemann. 2023. Domain-specific continued pretraining of language models for capturing long context in mental health. *ArXiv*, abs/2304.10447.

Robin Jia, Cliff Wong, and Hoifung Poon. 2019. Document-level n-ary relation extraction with multiscale representation learning. 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 3693–3704, Minneapolis, Minnesota. Association for Computational Linguistics.

Lavender Yao Jiang, Xujin C. Liu, Nima Pour Nejatian, Mustafa Nasir-Moin, Duo Wang, Anas Abidin, Kevin Eaton, Howard A. Riina, Ilya Laufer, Paawan Punjabi, Madeline Miceli, Nora C. Kim, Cordelia Orillac, Zane Schnurman, Christopher Livia, Hannah Weiss, David Kurland, Sean N. Neifert, Yosef Dastagirzada, Douglas Kondziolka, Alexander T. M. Cheung, Grace Yang, Mingzi Cao, Mona G. Flores, Anthony B Costa, Yindalon Aphinyanaphongs, Kyunghyun Cho, and Eric Karl Oermann. 2023. Health system-scale language models are all-purpose prediction engines. *Nature*, 619:357 – 362.

Qiao Jin, Zifeng Wang, Charalampos S Floudas, Jimeng Sun, and Zhiyong Lu. 2023. Matching patients to clinical trials with large language models. *ArXiv*.

Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Li wei H. Lehman, Mengling Feng, Mohammad Mahdi Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3.

Tian Kang, Shirui Zou, and Chunhua Weng. 2019. Pretraining to recognize pico elements from randomized controlled trial literature. *Studies in health technology and informatics*, 264:188 – 192.

Zeljko Kraljevic, Dan Bean, Anthony Shek, Rebecca Bendayan, Harry Hemingway, Joshua Au Yeung, Alexander Deng, Alfie Baston, Jack Ross, Esther Idowu, James T Teo, and Richard J Dobson. 2023. Foresight – generative pretrained transformer (gpt) for modelling of patient timelines using ehrs.

Sunjun Kweon, Junu Kim, Jiyoun Kim, Sujeong Im, Eunbyeol Cho, Seongsu Bae, Jungwoo Oh, Gyubok Lee, Jong Hak Moon, Seng Chan You, Seungjin Baek, Chang Hoon Han, Yoon Bin Jung, Yohan Jo, and E. Choi. 2023. Publicly shareable clinical large language model built on synthetic clinical notes. *ArXiv*, abs/2309.00237.

Sunjae Kwon, Zonghai Yao, Harmon Jordan, David Levy, Brian Corner, and Hong Yu. 2022. MedJEx: A medical jargon extraction model with Wiki's hyperlink span and contextualized masked language model score. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11733–11751, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36:1234 – 1240.

Eric P. Lehman, Evan Hernandez, Diwakar Mahajan, Jonas Wulff, Micah J. Smith, Zachary M. Ziegler, Daniel Nadler, Peter Szolovits, Alistair E. W. Johnson, and Emily Alsentzer. 2023. Do we still need clinical language models? *ArXiv*, abs/2302.08091.

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.

Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoyanov. 2020. Pretrained language models for biomedical and clinical tasks: Understanding and extending the state-of-the-art. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 146–157, Online. Association for Computational Linguistics.

Fei Li, Yonghao Jin, Weisong Liu, Bhanu Pratap Singh Rawat, Pengshan Cai, and Hong Yu. 2019. Fine-tuning bidirectional encoder representations from transformers (bert)–based models on large-scale electronic health record notes: An empirical study. *JMIR Med Inform*, 7(3):e14830.

Fei Li and Hong Yu. 2020. Icd coding from clinical text using multi-filter residual convolutional neural network. *Proceedings of the AAAI Conference on Artificial Intelligence. AAAI Conference on Artificial Intelligence*, 34 5:8180–8187.

Yikuan Li, Ramsey M. Wehbe, Faraz S. Ahmad, Hanyin Wang, and Yuan Luo. 2022. Clinical-longformer and clinical-bigbird: Transformers for long clinical sequences. *ArXiv*, abs/2201.11838.

Qiuhao Lu, Dejing Dou, and Thien Nguyen. 2022. ClinicalT5: A generative language model for clinical text. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5436–5443, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Neural Information Processing Systems*.

Avijit Mitra, Richeek Pradhan, Rachel D. Melamed, Kun Chen, David C. Hoaglin, Katherine L. Tucker, Joel I. Reisman, Zhichao Yang, Weisong Liu, Jack Tsai, and Hong Yu. 2023. Associations Between Natural Language Processing–Enriched Social Determinants of Health and Suicide Death Among US Veterans. *JAMA Network Open*, 6(3):e233079–e233079.

Niklas Muennighoff, Alexander M. Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. 2023. Scaling data-constrained language models. *ArXiv*, abs/2305.16264.

James Mullenbach, Sarah Wiegreffe, Jon Duke, Jimeng Sun, and Jacob Eisenstein. 2018. Explainable prediction of medical codes from clinical text. 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 1101–1111, New Orleans, Louisiana. Association for Computational Linguistics.

Mihir Parmar, Aakanksha Naik, Himanshu Gupta, Disha Agrawal, and Chitta Baral. 2023. Longbox: Evaluating transformers on long-sequence clinical tasks. *ArXiv*, abs/2311.09564.

C.A.I. Peng, Xi Yang, Aokun Chen, Kaleb E. Smith, Nima M. Pournejatian, Anthony B Costa, Cheryl Martin, Mona G. Flores, Ying Zhang, Tanja Magoc, Gloria P. Lipori, Duane A. Mitchell, Naykky Maruquel Singh Ospina, Mustafa Mamon Ahmed, William R. Hogan, Elizabeth A. Shenkman, Yi Guo, Jiang Bian, and Yonghui Wu. 2023. A study of generative large language model for medical research and healthcare. *NPJ Digital Medicine*, 6.

Amber Stubbs, Michele Filannino, Ergin Soysal, Sam Henry, and Özlem Uzuner. 2019. Cohort selection for clinical trials: n2c2 2018 shared task track 1. *Journal of the American Medical Informatics Association : JAMIA*.

Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2020. Long range arena: A benchmark for efficient transformers. *ArXiv*, abs/2011.04006.

Niall Taylor, Yi Zhang, Dan W. Joyce, Ziming Gao, Andrey Kormilitzin, and Alejo Nevado-Holgado. 2023. Clinical prompt learning with frozen language models. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–11.

Jiaqi Wang, Junyu Luo, Muchao Ye, Xiaochen Wang, Yuan Zhong, Aofei Chang, Guanjie Huang, Ziyi Yin, Cao Xiao, Jimeng Sun, and Fenglong Ma. 2024. Recent advances in predictive modeling with electronic health records.

Yuqing Wang, Yun Zhao, and Linda Petzold. 2023. Are large language models ready for healthcare? a comparative study on clinical language understanding.

Zifeng Wang and Jimeng Sun. 2022. PromptEHR: Conditional electronic healthcare records generation with prompt learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2873–2885, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Sarah Wiegreffe, Edward Choi, Sherry Yan, Jimeng Sun, and Jacob Eisenstein. 2019. Clinical concept extraction for document-level coding. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 261–272, Florence, Italy. Association for Computational Linguistics.

Cliff Wong, Sheng Zhang, Yu Gu, Christine Moung, Jacob Abel, Naoto Usuyama, Roshanthi Weerasinghe, Brian Piening, Tristan Naumann, Carlo Bifulco, and Hoifung Poon. 2023. Scaling clinical trial matching using large language models: A case study in oncology.

Michael Wornow, Alejandro Lozano, Dev Dash, Jenelle A. Jindal, Kenneth W. Mahaffey, and Nigam H. Shah. 2024. Zero-shot clinical trial patient matching with llms. *ArXiv*, abs/2402.05125.

Stephen T Wu, Kirk Roberts, Surabhi Datta, Jingcheng Du, Zongcheng Ji, Yuqi Si, Sarvesh Soni, Qiong Wang, Qiang Wei, Yang Xiang, Bo Zhao, and Hua Xu. 2019. Deep learning in clinical natural language processing: a methodical review. *Journal of the American Medical Informatics Association : JAMIA*.

Xi Yang, Aokun Chen, Nima M. Pournejatian, Hoo-Chang Shin, Kaleb E. Smith, Christopher Parisien, Colin B. Compas, Cheryl Martin, Anthony B Costa, Mona G. Flores, Ying Zhang, Tanja Magoc, Christopher A. Harle, Gloria P. Lipori, Duane A. Mitchell, William R. Hogan, Elizabeth A. Shenkman, Jiang Bian, and Yonghui Wu. 2022a. A large language model for electronic health records. *NPJ Digital Medicine*, 5.

Zhichao Yang, Avijit Mitra, Weisong Liu, Dan Berlowitz, and Hong Yu. 2023. Transformehr: transformer-based encoder-decoder generative model to enhance prediction of disease outcomes using electronic health records. *Nature Communications*, 14.

Zhichao Yang, Shufan Wang, Bhanu Pratap Singh Rawat, Avijit Mitra, and Hong Yu. 2022b. Knowledge injected prompt based fine-tuning for multi-label few-shot ICD coding. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1767–1781, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Ning Zhang and Maciej Jankowski. 2022. Hierarchical bert for medical document understanding. *ArXiv*, abs/2204.09600.

# A Appendix

## A.1 Text preprocessing

We followed Huang et al. (2019) to format notes during text preprocessing. But We did not convert text to lowercase because Mamba tokenizer is able to process both upper and lower cases. For notes on each patient's hospital visit, we sorted notes by their charted date and concatenated notes into one string. We used string "- - {NoteType} note - -" to separate the notes. Table A.2 shows comprehensive values of {NoteType}.

For pretraining data, we truncate notes with more than 16k tokens, however, this is only less than 2%, a length distribution is provided in Figure A.1. We exclude a small subset (6,049) of hospital visits due to the evaluation of MIMIC ICD coding and MIMIC hospital readmission prediction, The visit ids (hadm_id) are documented in the github.

## A.2 Pretraining recipe

ClinicalMamba-2.8b is a selective state space model designed using replication of the Mamba architecture (Gu and Dao, 2023). ClinicalMamba refers to the class of models, while 2.8b represents the number of parameters of this particular pretrained model. We also pretrained ClinicalMamba-130m using the same pretraining data from the previous section. The specific values of hyperparameters are shown in Table A.3. These models were trained for 763 million English tokens over 7000 steps (3 epochs) (Muennighoff et al., 2023). It was trained as an autoregressive language model, using cross-entropy loss (Brown et al., 2020). For learning rate scheduling, we followed Mamba and chose linear learning rate warmup with cosine decay to $1e - 5$. We found this important setting to avoid loss overflow. It took under 60 hours to pretrain ClinicalMamba-2.8b in on 4 Nvidia Tesla A100-80GB GPUs.



Figure A.1: Long tail distribution of number of tokens per each visit. Y-axis is the density (sum to 1.0).

|       |        | Cohort selection | Code-rare | Code-common |
|-------|--------|-----------------:|----------:|------------:|
| shots | mean   | 89 | 5 | 918 |
| tokens | mean  | 4924 | 4223 | 7062 |
|        | median | 4632 | 3236 | 5177 |
|        | 99%    | 10781 | 14345 | 13356 |
|        | max    | 13989 | 18480 | 14773 |

Table A.1: Number of instances per label (shots) and number of tokens per input.

| Category | Count | % | Len |
|----------|------:|--:|----:|
| Nursing | 506,528 | 73 | 241 |
| Radiology | 338,834 | 83.3 | 449 |
| ECG | 123,042 | 61.3 | 43 |
| Physician | 92,426 | 18.2 | 1369 |
| Discharge summary | 47,572 | 96.7 | 2195 |
| Echo | 34,064 | 45.8 | 464 |
| Respiratory | 32,798 | 8.1 | 205 |
| Nutrition | 7,971 | 6.4 | 602 |
| General | 7,710 | 6.4 | 290 |
| Rehab Services | 5,321 | 4.6 | 622 |
| Social Work | 2,294 | 2.8 | 446 |
| Case Management | 939 | 1.3 | 260 |
| Pharmacy | 97 | 0.1 | 512 |
| Consult | 78 | 0.1 | 1206 |

Table A.2: Statistic of note events documented in MIMIC-III dataset. Each column represents a) the number of notes, b) proportion of visits, c) average number of words for each note type.

| Hyperparameter | Value |
|----------------|------:|
| num param | 130m/2.8b |
| num layer | 24/64 |
| dim model | 768/2560 |
| context len | 16k |
| num vocab | 50277 |
| position emb | None |
| optimizer | Adam |
| beta1 | 0.9 |
| beta2 | 0.95 |
| epsilon | 1e-5 |
| batch size | 32 |
| weight decay | 0.1 |
| gradient clipping | 1.0 |
| peak learning rate | 1e-3/6e-4 |

Table A.3: Hyperparameters used to train ClinicalMamba.

# Working Alliance Transformer for Psychotherapy Dialogue Classification

**Baihan Lin[1], Guillermo Cecchi[2], Djallel Bouneffouf[2]**
[1]Icahn School of Medicine at Mount Sinai, New York, NY
[2]IBM TJ Watson Research Center, Yorktown Heights, NY
baihan.lin@mssm.edu, {gcecchi@us., djallel.bouneffouf@}ibm.com

## Abstract

As a predictive measure of the treatment outcome in psychotherapy, the working alliance measures the agreement of the patient and the therapist in terms of their bond, task and goal. Long been a clinical quantity estimated by the patients' and therapists' self-evaluative reports, we believe that the working alliance can be better characterized using natural language processing technique directly in the dialogue transcribed in each therapy session. In this work, we propose the Working Alliance Transformer (WAT), a Transformer-based classification model that has a psychological state encoder which infers the working alliance scores by projecting the embedding of the dialogues turns onto the embedding space of the clinical inventory for working alliance. We evaluate our method in a real-world dataset with over 950 therapy sessions with anxiety, depression, schizophrenia and suicidal patients and demonstrate an empirical advantage of using information about therapeutic states in the sequence classification task of psychotherapy dialogues.

## 1 Introduction

The working alliance between the therapist and the patient is an important measure of the clinical outcome and a qualitative predictor of therapeutic effectiveness in psychotherapy (Wampold, 2015; Bordin, 1979). The alliance entails a number of cognitive and emotional aspects of the interaction between these two agents, such as their shared understanding of the objectives to be attained and the tasks to be completed, as well as the bond, trust, and respect that will develop during the course of the therapy. While traditional methods to quantify the alliance depend on self-evaluative reports with point-scales valuation by patients and therapists about whole sessions (Horvath, 1981), the digital era of mental health can enable new research fronts utilizing real-time transcripts of the dialogues between the patients and therapists. By analyzing the psychotherapy dialogues, we are interested in studying the usage of natural language processing technique to extract out turn-level features of the working alliance and see if it can help better inform us of the clinical condition of the patient.

Here we present Working Alliance Transformer (WAT), a transformer-based classification model to classify the psychotherapy sessions into different psychiatric conditions. Our methods consists of a psychological state encoder that quantifies the degree of patient-therapist alliance by projecting each turn in a therapeutic session onto the representation of clinically established working alliance inventories, using language modeling to encode both turns and inventories, which was originally proposed in (Lin et al., 2022) as an analytical tool. This allows us not only to quantify the overall degree of alliance but also to identify granular patterns its dynamics over shorter and longer time scales.

We collated and preprocessed the Alex Street Counseling and Psychotherapy Transcripts dataset (Street, 2023), which consists of transcribed recordings of over 950 therapy sessions between multiple anonymized therapists and patients that belong to four types of psychiatric conditions: anxiety, depression, schizophrenia and suicidal. (The data publisher mentions that they have more clinical conditions other than the analyzed 4 classes, but due to the licensing and access limitations, we can only obtain the 4 classes we presented.) This multi-part collection includes speech-translated transcripts of the recordings from real therapy sessions, 40,000 pages of client narratives, and 25,000 pages of reference works. As open science and data sharing initiatives in the psychiatry domains become more prominent, we believe our methodologies can be adapted in a responsible way to a broader spectrum of clinical conditions. On this dataset, we evaluate quantitatively the effectiveness of this inference method in improving the classification / diagnosis capability of deep learning models to linguistically

predict psychiatric conditions from therapy transcripts. Lastly, we discuss how our approach may be used as a companion tool to provide feedback to the therapist and to augment learning opportunities for training therapists.

## 2 Methods

We describe our pipeline in Figure 1. Given the transcripts of a therapy session and the medical records of the patient. The dialogue are separated into pairs of turns as the timestamps. We can either choose to only use the turns by the patients, or by the therapists, or use both, as a paired input. Empirically, the patients' turns are usually more narrative, as they are describing themselves, while the therapists' turns are usually more declarative, as they are usually confirming the patients, or leading conversations to certain topic.

Each patient response turn $S_i^p$ followed by a therapist response turn $S_i^t$ is treated as a dialogue pair. In total, these materials include over 200,000 turns together for the patient and therapist and provide access to the broadest range of clients for our linguistic analysis of the therapeutic process of psychotherapy. On the other hand, we have access to the Working Alliance Inventory (WAI), the clinical instrument. The modern WAI consists of 36 statements in a self-report questionnaire which measures the therapeutic bond, task agreement, and goal agreement (Horvath, 1981; Tracey and Kokotovic, 1989; Martin et al., 2000), where the Since the original 12-item version (Tracey and Kokotovic, 1989), the inventory has used parallel versions for clients and therapist with good psychometric properties and helped establish the importance of therapeutic alliance in predicting treatment outcomes. The modern version of the inventory consists of 36 questions, where the rater (i.e. the patient or the therapist) is asked to rate each statement on a 7-point scale (1=never, 7=always)(Martin et al., 2000). This inventory is disorder-agnoistic, meaning that it measures the alliance factors across all types of therapies, and provides a record of the mapping from the alliance measurement and the corresponding cognitive constructs underlying the measurement under a unified theory of therapeutic change (Horvath and Greenberg, 1994).

The inference goal is to compute a score that characterizes the working alliance given the clinical inventory, with for instance, a feature vector of 36 dimension that correspond to the 36 alliance

---

**Algorithm 1** Working Alliance Transformer (WAT)

1: **Input**: a session with $T$ turns
2: **Output**: a label for psychiatric condition
3: **for** i = 1,2,···, T **do**
4:     Transcribe dialogue turn pairs $(S_i^p, S_i^t)$
5:     **for** $(I_j^p, I_j^t) \in$ inventories $(I^p, I^t)$ **do**
6:         $W_j^{p_i} = \mathrm{similarity}(Emb(I_j^p), Emb(S_i^p))$
7:         $W_j^{t_i} = \mathrm{similarity}(Emb(I_j^t), Emb(S_i^t))$
8:     **end for**
9:     (Patient) $x_c = concat(Emb(S_i^t), W^{p_i})$
10:     (Therapist) $x_t = concat(Emb(S_i^t), W^{t_i})$
11:     (Dyad) $x = concat(x_t, x_c)$
12:     Aggregated feature $X.append(x)$
13: **end for**
14: obtain prediction $y = Transformer(X)$

---

measure of interests in the inventory. Operationally, the goal is to derive from these 36 items three alliance scales: the task scale, the bond scale and the goal scale. They measures the three major themes of psychotherapy outcomes: (1) the collaborative nature of the patient-therapist relationship; (2) the affective bond between therapist and patient, and (3) the therapist's and patient's capabilities to agree on treatment-related short-term tasks and long-term goals. The score corresponding to the three scales comes from a key table which specifies the positivity or the sign weight to be applied on the questionnaire answer when summing in the end. The full scale is simply the sum of the scores of the three scales. The key table is like a weighting matrix that specifies the directionalities of the scales. After computing the information regarding the predicted clinical outcome with our inferred working alliance scores, this feature vector highlights a bias towards what the clinicians would care about in the psychotherapy given the metrics provided by the working alliance inventory. We would then able to further use this information to potentially inform us of the psychiatric condition of a given patient. As such, we propose the Working Alliance Transformer (WAT), a classification model that utilizes an inference module that informs the downstream classifier where the current state is with respect to the therapeutic trajectory or landscape in the psychotherapy treatment of this patient. Is this patients approaching a breakthrough? Or is he or she susceptible to a rupture of trust? These therapeutic information about alliance can vary across clinical conditions, and thus, potentially beneficial to the

Figure 1: Architecture of working alliance transformer for psychiatric condition classification using the psychological state encoder from working alliance

diagnosis and monitoring of psychiatric disorders.

Algorithm 1 outlines the classification process. During the session, the dialogue between the patient and therapist are transcribed into pairs of turns. We denote the patient turn as $S_i^p$ followed by the therapist turn $S_i^t$, as a dialogue pair. Similarly, the inventories of working alliance questionnaires come in pairs ($I^p$ for the patient, and $I^t$ for the therapist, each with 36 statements). We compute the distributed representations of both the dialogue turns and the inventories with the sentence embeddings. The working alliance scores can then be computed as the cosine similarity between the embedding vectors of the turn and its corresponding inventory vectors. Following (Lin et al., 2022), we use SentenceBERT (Reimers and Gurevych, 2019) and Doc2Vec embedding (Le and Mikolov, 2014) as our sentence embeddings for the working alliance inference. With that, for each turn (either by patient or by therapist), we obtain a 36-dimension working alliance score. For the classification, we concatenate the 36-dimension working alliance scores computed from the current turn in the dialogue, along with the sentence embedding of the current turn, as our feature vector to fed into our Transformer sequence classifier.

The analytical features enabled by the working alliance inference are not only useful for the classification we investigate in this study but also other downstream tasks, such as predictive modeling and real-time analytics. In our case, the turns in a dialogue or monologue are fed into the sentence embedding sequentially as individual entries. And then, given the sentence embedding, we feed them each into the psychological state encoder that in-

fer the psychological or therapeutic state of the dialogue at this turn. The encoder will generate a vector that characterizes the state, such as the 36-dimension working alliance scores, corresponding to the 36 working alliance inventory items. Then, the model aggregate both the sentence embedding feature vector and the psychological state vector. In this case, we concatenate them together as a first step. Since we feed our input sentence by sentence (or turn by turn), we have a sequence of combined feature vector, which is then fed into a sequence classifier. We use the transformer (Vaswani et al., 2017) as our classifier for its effectiveness in various sequence-based learning tasks, and potential interpretability from its attention weights. The output of this classification model is the predicted clinical condition of this sequence. The sequence of turns we feed to generate a label is a trimmed segment of the session of psychotherapy transcript.

## 3 Results

Here we present the transcript classification results.

**Experimental setting.** The psychotherapy dataset we evaluate is highly *imbalanced* across the four clinical conditions (495 anxiety sessions, 373 depression sessions, 71 schizophrenia sessions, and 12 suicidal sessions). If we directly train our models on this dataset, the classifier is likely to be highly biased towards the majority class. To correct for this imbalanceness issue, we are using the sampling technique. Instead of going through the entire training data in epochs, we train the models in sampling iterations. In each iteration we randomly choose a class and then randomly sample one session from the class pool. Before we sample

66

Table 1: Classification accuracy (%) of psychotherapy sessions

| | SentenceBERT | | | Doc2Vec | | |
|---|---|---|---|---|---|---|
| | Patient turns | Therapist turns | Both turns | Patient turns | Therapist turns | Both turns |
| WAT (working alliance embedding) | **27.6** | **27.0** | **26.0** | **34.1** | 25.7 | **31.9** |
| WAT (working alliance score) | 26.1 | 23.4 | 25.5 | 28.9 | 23.7 | **31.9** |
| Embedding Transformer | 24.8 | 24.0 | 25.5 | 31.8 | **26.2** | 29.9 |
| WA-LSTM (working alliance embedding) | **35.0** | **36.9** | 23.3 | **46.0** | 27.7 | 29.6 |
| WA-LSTM (working alliance score) | 24.5 | 34.2 | 22.6 | 30.2 | 24.7 | **43.4** |
| Embedding LSTM | 23.0 | 36.0 | 22.9 | 44.3 | **31.1** | 31.1 |

the sessions, we split the dataset into 20/80 as our test set and training set. Then during the training or the test phase, we perform the sampling technique for each iteration only in the fully separated training and test sets. Then, for each sampled session, we feed into the classification model the first 50 dialogue turns of our transcript, turn by turn, and the sequence classifier will output a label predicting which psychiatric condition this session belongs to.

**Model architecture.** We evaluate two classifier backbones. The first one is the classical transformer model. For the multi-head attention module, we set the number of heads to be 4 and the dimension of the hidden layer to be 64. The dropout rates for the positional encoding layer and the transformer blocks are both set to be 0.5. The second backbone is a 64-neuron Long Short-Term Memory (LSTM) network (Hochreiter and Schmidhuber, 1997).

**Ablation and baseline models.** For each of the two classifiers, we compare three types of features as the input we feed into the sequence classifier component. The first one, the working alliance embedding, is the concatenated feature vector of both the sentence embedding vector and the psychological state vector (which in our case, is the 36-dimension inferred working alliance scores). The second type of feature, the working alliance score, is an ablation model which only uses the state vector (the working alliance score vector). The third type of feature, the embedding, is the baseline which only uses the sentence embedding vector directly. In other words, The working alliance score introduces the bias for WAI. The sentence embedding doesn't. The working alliance embedding is the feature that combines both with concatenation. And since we have two sentence embeddings to choose from (the sentence BERT and Doc2Vec), they each have 9 models in the evaluation pool. Other than the classifier types (Transformer or LSTM), the embedding types (Sentence-BERT or Doc2Vec) and the feature types (working alliance embedding, working alliance scores, or

simply sentence embedding), we also compared using only the dialogue turns from the patients, from the therapists, and from both the patients and the therapists. In the case where we use the turns from both the patients and the therapists, we consider them as a pair, and concatenate them together as a combined feature. This is as opposed to treating them as subsequent sequences, because we believe that the therapist's response are loosely semantic labels for the patient's statements, and thus, serve different semantic contexts that should be considered side by side, instead of sequentially, which would assume a homogeneity between time steps.

**Training procedure.** For all 12 models, we train them for over 50,000 iterations using the stochastic gradient descent with a learning rate of 0.001 and a momentum of 0.9. Since the training set is relatively small for our neural network models, we observe some of the models exhibit overfitting at early stages before we finish the training. As a result, we report the performance of their checkpoints where they converge and have a plateau performance. Then in the testing phase, we randomly sample class-balanced 1,000 samples.

**Empirical results.** We report the classification accuracy as our evaluation metrics. Since we have four classes, and the evaluation is corrected for imbalanceness with the sampling technique.

Overall, we observe a benefit of using the working alliance embedding as our features in Transformer and LSTM-based model architectures. Among all the models, the WA-LSTM model with working alliance embedding using only the patient turns obtains the best classification result (46%), followed by the WA-LSTM model using only the working alliance score using both turns from the patients and therapists (43.4%). This suggest the advantage of taking into account the predicted clinical outcomes in characterizing these sessions given their clinical conditions. We also notice that the inference of the therapeutic working alliance with Doc2Vec appears to be more beneficial in model-

ing the patient turns than the therapist turns, while the SentenceBERT-based inference appears to be advantageous in both therapist and patient features.

Comparing the two sequential learners, the Transformer, due to the additional attention mechanism, yields a more stable learning phase. When using SentenceBERT as its embedding, we observe a modest benefit when training on only the patient turns, which might suggest an interference of features between the therapists' and patients' working alliance information. The Transformers using the working alliance embedding, i.e. both the sentence embedding and their therapeutic states (i.e. the inferred working alliance score vector) are the best performing ones. When using Doc2Vec as the embedding, the best performing models are both the Transformers using some of the working alliance information from our inference module as features.

## 4 Discussion

Our analytic approach reveals insightful features of therapeutic relationship and their usefulness in terms of clinical diagnosis merely based on the patient-doctor conversations. In our prior work, we observe systematic differences in the mean inferred alliance scores between patients and therapists, and also across disorders (Lin et al., 2022). However the in-session evolution of the inferred scores provide a much more interesting perspective, as shown in our dialogue sequence classification results. In particular, while all conditions show a systematic misalignment of scores between patients and therapists, this is significantly starker for suicidality, something observable in the mean as well as in the time trace for full and sub-scales, which can be useful for early detection of suicidal thoughts.

As more and more successful applications of AI are deployed in clinical domains, there are many ethical considerations we practitioners of machine learning should be aware of and take into considerations. When dealing with patient data, the privacy and security is a top priority. Following the suggestion of best practices from (Matthews et al., 2017), all examples in this paper as well as the dataset we analyzed are properly anonymized with pre- and post-processing techniques. In addition, the dataset itself was sourced with proper license from ProQuest's Alexander Street platform. We remove all personally identifiable information (meta data, user name, identifiers, doctors' name) from the dataset.

Since the clinical domain of this work is men-

tal health and psychological well-being, there are additional ethical considerations. Emerging techniques in wearable devices, digital health records, brain imaging measurements, smartphone applications and social media are gradually transforming the landscape of the monitoring and treatment of mental health illness. However, most of these attempts are proof of concept as identified by this review (Graham et al., 2019), and requires extensive caution to prevent from the pitfall of over-interpreting preliminary results. The limitations of these prior studies, including our work here, reside in the difficulty of a systematic clinical validation and a uncertain future expectation of the technological readiness for patient care and therapeutic decision making approved by authorities. For instance, it was recently shown that despite the high predictability of statistical learning-based methods in analyzing large datasets in support of clinical decisions in psychiatry, existing machine learning solutions is highly susceptible to overfitting in realistic tasks which has usually a small sample sizes in the data, missing data points for some subjects, and highly correlated variables (Iniesta et al., 2016). These properties in real-world applications limits the out-of-sample generalizability of the results.

## 5 Conclusions

In this work, we present a Transformer-based classification model that characterizes the sequence of therapeutic states as beneficial feature to improve the classification of psychological dialogues into different psychiatric conditions. It combines the domain expertise from clinically validated psychiatry inventories with the distributed deep representations of language modeling provide a turn-level encoding of working alliance at a turn-level resolution. We demonstrate on a real-world psychotherapy dialogue dataset that using this additional granular representation of the interaction dynamics between patients and therapists is beneficial both for interpretable post-session insights and linguistically diagnosing the patients.

Our results suggest that the inferred scores of therapeutic or psychological states of patient-doctor alignment can be useful in downstream tasks, such as diagnosis. Although not a main focus in this work, future work would include a more systematic investigation of such downstream tasks, and exploiting the attention mechanism of the transformer blocks for interpretations.

# References

Edward S Bordin. 1979. The generalizability of the psychoanalytic concept of the working alliance. *Psychotherapy: Theory, research & practice*, 16(3):252.

Sarah Graham, Colin Depp, Ellen E Lee, Camille Nebeker, Xin Tu, Ho-Cheol Kim, and Dilip V Jeste. 2019. Artificial intelligence for mental health and mental illnesses: an overview. *Current psychiatry reports*, 21(11):1–18.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Adam O Horvath. 1981. *An exploratory study of the working alliance: Its measurement and relationship to therapy outcome*. Ph.D. thesis, University of British Columbia.

Adam O Horvath and Leslie S Greenberg. 1994. *The working alliance: Theory, research, and practice*, volume 173. John Wiley & Sons.

Raquel Iniesta, D Stahl, and Peter McGuffin. 2016. Machine learning, statistical learning and the future of biological research in psychiatry. *Psychological medicine*, 46(12):2455–2465.

Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196. PMLR.

Baihan Lin, Guillermo Cecchi, and Djallel Bouneffouf. 2022. Deep annotation of therapeutic working alliance in psychotherapy. *arXiv preprint arXiv:2204.05522*.

Daniel J Martin, John P Garske, and M Katherine Davis. 2000. Relation of the therapeutic alliance with outcome and other variables: a meta-analytic review. *Journal of consulting and clinical psychology*, 68(3):438.

Tara Matthews, Kathleen O'Leary, Anna Turner, Manya Sleeper, Jill Palzkill Woelfer, Martin Shelton, Cori Manthorne, Elizabeth F Churchill, and Sunny Consolvo. 2017. Stories from survivors: Privacy & security practices when coping with intimate partner abuse. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 2189–2201.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.

Alexander Street. 2023. counseling and psychotherapy transcripts series.

Terence J Tracey and Anna M Kokotovic. 1989. Factor structure of the working alliance inventory. *Psychological Assessment: A journal of consulting and clinical psychology*, 1(3):207.

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.

Bruce E Wampold. 2015. How important are the common factors in psychotherapy? an update. *World Psychiatry*, 14(3):270–277.

# Building A German Clinical Named Entity Recognition System without In-domain Training Data

**Siting Liang[1], Hans-Jürgen Profitlich[1], Maximilian Klass[2], Niko Möller-Grell[2],**
**Celine-Fabienne Bergmann[2], Simon Heim[2], Christian Niklas[2], Daniel Sonntag[1,*]**

[1]German Research Center for Artificial Intelligence, Germany
[2]Heidelberg University Hospital, Germany
[*]University of Oldenburg, Germany
siting.liang|hans-juergen.profitlich|daniel.sonntag@dfki.de
maximilian.klass|christian.niklas@med.uni-heidelberg.de

## Abstract

Clinical Named Entity Recognition (NER) is essential for extracting important medical insights from clinical narratives. Given the challenges in obtaining expert training datasets for real-world clinical applications related to data protection regulations and the lack of standardised entity types, this work represents a collaborative initiative aimed at building a German clinical NER system with a focus on addressing these obstacles effectively. In response to the challenge of training data scarcity, we propose a **Conditional Relevance Learning (CRL)** approach in low-resource transfer learning scenarios. **CRL** effectively leverages a pre-trained language model and domain-specific open resources, enabling the acquisition of a robust base model tailored for clinical NER tasks, particularly in the face of changing label sets. This flexibility empowers the implementation of a **Multilayered Semantic Annotation (MSA)** schema in our NER system, capable of organizing a diverse array of entity types, thus significantly boosting the NER system's adaptability and utility across various clinical domains. In the case study, we demonstrate how our NER system can be applied to overcome resource constraints and comply with data privacy regulations. Lacking prior training on in-domain data, feedback from expert users in respective domains is essential in identifying areas for system refinement. Future work will focus on the integration of expert feedback to improve system performance in specific clinical contexts.

## 1 Introduction

Clinical Named Entity Recognition (NER) plays a central role in extracting valuable information from medical texts as essential features for developing clinical decision support systems. In this work, we concentrate on the German language and its application within clinics in Germany. Clinical documents can originate from a variety of sources. Each source has its unique characteristics, making it challenging to develop a one-size-fits-all NER system (Sonntag et al., 2016; Sonntag and Profitlich, 2019; Profitlich and Sonntag, 2021; Borchert et al., 2022; Roller et al., 2022). In developing German clinical NER systems, challenges arise when strict privacy rules are applied to data sources and the complexities associated with expert annotation of training data (Kittner et al., 2021; Roller et al., 2022). Previous related research efforts in German language have tackled these challenges using a range of techniques, from rule-based approaches to transfer learning methods in low resources scenarios (Frei and Kramer, 2021; Schäfer et al., 2022; Liang et al., 2023b,a). In this paper, we describe the development and assessment of an adaptive German clinical NER system without prior training on in-domain data. This goal is motivated by the principles of Interactive Machine Learning (IML) (Fails and Olsen Jr, 2003; Dudley and Kristensson, 2018), particularly when dealing with the challenges of annotating complex medical texts.

In this work, we investigate innovative transfer learning techniques and support the evolution of dynamic annotation schemas by engaging expert users from the medical field. We aim to address the constraints posed by limited data resources. The system has been developed as part of a cooperative project involving a machine learning lab and a university institute for medical informatics. Our system is equipped with a dedicated web-based User Interface (web-UI) for correcting system-generated annotations, which is instrumental in our case study involving **cardiology**. Qualified experts who are granted access to review the specific documents sourced from the hospital's internal database in medical informatics can utilize this tool to interact with system-generated outputs via a standalone website. The main objective of the case study is to conduct a comprehensive analysis of the performance of the NER system when applied to a non-

distributed dataset without violating privacy regulations. This analysis helps identify areas where adjustments are needed to refine the annotation schema and offers valuable insights to guide further fine-tuning of the model's performance to maintain its relevance to domain-specific nuances, context, and entity variations. Figure 1 displays our collaborative research environment. All system-related modules and model checkpoints are deployed at the hospital endpoint to ensure strict data security.



Figure 1: Overview of our collaborative research between experts in the field of interactive machine learning and medical informatics.

The main contributions of our work in the field of German clinical NER are as follows:

- Firstly, we leverage cross-domain transfer learning methods inspired by Liang et al. (2023a), using pre-trained language models and domain-relevant open-source Germain language datasets. We refer the approach to **Conditional Relevance Learning (CRL)** with an architecture that extends from a BERT-based encoder (Kenton and Toutanova, 2019) and incorporates a token-level binary classifier (see subsection 3.1). **CRL** has the potential to reduce the need of domain-specific training data compared to data-specific classifiers in low-resource scenarios. It also offers significant flexibility in adapting to changing label sets across different clinical domains.

- To enhance the adaptability and utility of our NER system across a range of medical texts (Widdows et al., 2002; Roller et al., 2022), we establish a comprehensive set of entity types categorised into six distinct semantic groups, drawing from the semantic ontology from the Unified Medical Language System (**UMLS**) Metathesaurus[1] (Bodenreider, 2004) and domain-specific annotations provided by Roller et al. (2022). This extensive annotation schema is referred to as **Multilayered Semantic Annotation (MSA)** (see subsection 3.3).

[1] http://umls.nlm.nih.gov

This forms the basis for a dynamic and expandable semantic annotation schema as new clinical use cases emerge over time. As the NER system is deployed in specific clinical contexts, we can refine the annotation schema to align with the unique requirements of each use case it serves, ensuring the system's effectiveness and relevance across diverse clinical scenarios with minimal modification needed.

- Moreover, our **case study** plays a critical role in our broader efforts to improve the adaptability and utility of our clinical NER system. It offers essential insights into the performance of the NER system lacking in-domain training data and highlights areas where improvements are necessary.

Our collaborative effort is achieving NER system's adaptability across clinical domains and improving the robustness of the NER system's performance through the engagement of domain experts in applications. Ultimately, we aim to address the impact of strict privacy rules on data accessibility and annotation quality and contribute to research on the development of clinical information extraction systems in the German healthcare sector. Codes and demonstration are available in GitHub repository https://github.com/sitingGZ/bert-sner-cardio.

## 2 Related Work

Efforts to address the scarcity of in-domain training data in German NER training have led to the exploration of two main strategies. One strategy involves translating annotated English corpora, such as the n2c2 dataset (Henry et al., 2020) and DDI dataset (Segura-Bedmar et al., 2013), to synthesize domain-related German training datasets (Frei and Kramer, 2021; Schäfer et al., 2022). However, the accuracy and feasibility of the resulting NER models remain limited by nuances and contextual differences between languages and clinical domains, which affect their applicability to real German clinical texts. Another line of work involves the manual annotation efforts of curating generic medical datasets (Widdows et al., 2002; Borchert et al., 2022) by annotating on open-source corpora derived from medical journals. Furthermore, there are ongoing endeavours to develop domain-specific German clinical data sets that adhere to the English reference guidelines (Kittner et al., 2021; Richter-Pechanski et al.,

2023).

Task-specific datasets, like Roller et al. (2022) are rare but valuable for enhancing NER model performance by capturing domain-specific nuances. However, their adoption is hindered by the resource-intensive manual annotation workload. Llorca et al. (2023); Liang et al. (2023a) attempted to generalise various entity labels from different annotated German datasets to achieve a reasonable amount of cross-domain training data. While Llorca et al. (2023) introduced harmonised versions of German medical corpora through the Big-BIO framework (Fries et al., 2022), contributing to the creation of common metadata sets for improved NER model performance across different medical text sources. However, the study found limited generalisation of NER models across different datasets. It remains challenging to effectively leverage diverse medical corpora for entity recognition in German texts. Liang et al. (2023a) augmented training data by mapping the original entity labels from source datasets to semantic types based on the ULMS ontologies and utilized a German BERT encoder with a binary token classifier to efficiently recognize medical entities by prompting with various semantic types followed by the same medical text. The novel training framework has shown effective cross-domain knowledge transfer and enhanced performance in low-resource German NER tasks.

While previous approaches offer potential solutions for the data scarcity issue, none have been able to develop a cross-domain adaptive NER system. Our work represents a step towards a more efficient and flexible long-term solution, as we combine the NER training framework from Liang et al. (2023a) with a multilayered semantic annotation schema, specifically targeting NER challenges within clinical information extraction and adapting to evolving clinical use cases.

## 3 Approach

### 3.1 Conditional Relevance Learning (CRL)

Developing in-domain training data from scratch requires substantial effort and resources for data collection and annotation, which is time-consuming and costly (Kittner et al., 2021; Roller et al., 2022; Richter-Pechanski et al., 2023). While transfer learning using pre-trained models can be beneficial, achieving acceptable performance still necessitates an effective transfer learning framework that leverages domain-specific fine-tuning techniques rather than relying solely on a large amount of labelled data in the source tasks (Llorca et al., 2023).

Liang et al. (2023a) shows that the performance of NER models firstly trained with domain-related corpus (Widdows et al., 2002) through a set of harmonized entity labels and novel training objective can be effectively generalised to clinical target tasks (Kittner et al., 2021; Roller et al., 2022) with much less fine-tuning data. While a harmonized label set proves highly advantageous in aggregating different relevant datasets to obtain a reasonable amount of training data, the limitation lies in the intricacies of the carefully designed matching process used to convert the entity labels from the target tasks to a unified label set. In this work, we only adopt the training framework from Liang et al. (2023a) which leverages pre-trained BERT-based encoder[2] and a token-level binary classifier on top of the BERT-based encoder predicts the contextual relevance score for individual tokens in a medical text input conditioned on the preceding label words. Table 1 presents two training examples of using the novel training objective from Liang et al. (2023a). The semantic type and the medical term phrases to be extracted are annotated as class 1. The remaining part of the input is marked as class 0. The training data also include negative samples, where no entity phrases can be extracted for the preceding semantic types.

| Input | Target |
|---|---|
| [CLS] **Clinical Drug** [SEP] **Zofran** 4mg for nausea | [0, 1, 1, 0, 1, 0, 0, 0] |
| [CLS] Diagnostic Procedure [SEP] Zofran 4mg for nausea | [0, 0, 0, 0, 0, 0, 0, 0] |

Table 1: Training example in line with the idea of conditional relevance learning. The model learns to recognize how different tokens in the input text should be associated with specific entity types.

Depart from the approach of Liang et al. (2023a), we preserve the original labels from the training sources, rather than their transformation into a unified label set. We aim to promote a more nuanced understanding of the diverse range of entities present in medical texts, ultimately improving the adaptability and effectiveness of the NER system. Figure 2 shows the format of the training data utilized in line with the approach of **CRL**. More information about the utilized datasets can be found

---

[2]https://www.deepset.ai/german-bert

Figure 2: Our training examples for CRL approach. During inference, users can select relevant labels or introduce new ones. User feedback ensures system adaptability.

in Subsection 3.2.

**CRL** allows the NER system to adapt to new entity types from different label sets. During inference, users have the flexibility to choose from the entire entity label set seen during training, select only the relevant ones, or even introduce new, unseen semantic types as needed. Furthermore, user feedback regarding the applicability of the system to their specific use cases is essential. Following this adaptation, we apply the trained model on an unseen clinical dataset from the cardiovascular domain in our case study and make NER classification based on the **MSA** schema which is explained in Subsection 3.3.

### 3.2 Training Data

**MUCHMORE**[3] is a bilingual labelled corpus collecting English and German abstracts from 41 medical journals and containing the semantic annotations mapped to UMLS medical concepts (Widdows et al., 2002). **Ex4CDS**[4] is a corpus containing entity annotations related to the outcomes of kidney transplantation in nephrology clinics (Roller et al., 2022). Both are readily available open-source datasets. In this work, we train and adapt the NER system for clinical applications using these datasets with **CRL**. The entity types present in the training

data from **MUCHMORE** and **Ex4CDS** are shown in Table 6 and Table 5 respectively (see Appendix A).

### 3.3 Multilayered Semantic Annotation (MSA)

During training, our goal is to encourage the model's comprehension of semantic diversity through a broad array of semantic types as entity labels. In practice, a medical phrase can represent an entity type of *Diagnostic Procedure* and contain the name of a *Medical Device* applied in the procedure. Multiple entity types can be predicted to the same text span, which is particularly valuable for nested and discontinuous NER tasks where entities are embedded within others (Yan et al., 2021). Hence, having a clear semantic annotation schema is beneficial for both application and performance analysis purposes.

The **UMLS** Metathesaurus integrating millions of medical concepts are widely applied knowledge sources for mining medical terms tasks (Aronson, 2001, 2006; Savova et al., 2010; Widdows et al., 2002; Borchert et al., 2022; Llorca et al., 2023). We design a **MSA** schema that aims to identify entities in multiple semantic dimensions based on **UMLS** ontology, thus the NER system can extract a broader spectrum of clinically relevant information, leading to the discovery of advanced medical knowledge. In addition to the incorporation of stan-

---

[3]https://muchmore.dfki.de/resources1.htm
[4]https://github.com/DFKI-NLP/Ex4CDS

dardised semantic types from **UMLS**, we add two semantic aspects, e.g. **Health State** and **Factuality** from Roller et al. (2022) in the **MSA** schema, illustrated in Table 2.

| Semantic Group | Entity Types |
|---|---|
| Physical Object | Anatomical Structure, Clinical Drug, Medical Device |
| Conceptual Entity | Clinical Attribute, Quantitative Concept, Laboratory or Test Result, Temporal Concept |
| Procedure | Laboratory Procedure, Diagnostic Procedure, Therapeutic or Preventive Procedure |
| Phenomenon or Process | Injury or Poisoning, Disease Physiologic Function, Pathologic Function |
| Health State | Healthy Condition, Deteriorated Condition |
| Factuality | Negated, Minor, Speculated |

Table 2: **MSA** encompasses six semantic groups. Each semantic group contains multiple entity types to facilitate fine-grained disambiguation. Most entity types are from the **UMLS** semantic types, except for groups **Health State** and **Factuality** (Roller et al., 2022).

## 3.4 Automatic Entity Annotation

**CRL** transforms the NER task into a token-level binary classification task, it predicts a relevance score for each token based on the preceding entity labels. We employ a threshold-based approach to transform the prediction scores made by the models into entity recognition results during the inference phase. Tokens with scores above the predefined threshold are considered part of entities and are assigned the corresponding entity type. In this process, we consider the variations in prediction scores, domain shifts, entity type unbalance, and the organization of entity types into semantic groups within **MSA**. They are critical aspects of ensuring the adaptability and effectiveness of the NER system in cross-domain transfer scenarios. As a result, the key steps include: **(1)** Prediction scores are generated for each token in the new dataset. **(2)** For each specific entity type, the maximum prediction score among the tokens associated with that type is identified, which serves as the maximum confidence for the entity type in the dataset. **(3)** The prediction scores for each token associated with an entity type are normalized by dividing the maximum confidence score to ensure a common range for comparison. **(4)** A set of thresholds is applied to determine the entity type assignment

for the token. Tokens with normalized prediction scores above the assigned threshold are labelled with the corresponding entity type. **(5)** In cases where entity types within the same semantic group may be assigned to the same span, we assign priorities to entity types based on their normalized prediction scores to determine the most appropriate entity type for that span. Figure 3 displays the annotation results of the selected document based on the **MSA**.



Figure 3: UI snippet of displaying MSA results for sentences of the selected document (above) and the selected semantic groups (left). The definitions of the entity types (coloured by groups) can be checked by clicking on the **[Check type definitions]** button.

**CRL** and **MSA** facilitate the integration of new semantic types and evolving clinical NER tasks without modification to the system architecture. However, building the NER system without in-domain training data remains challenging in understanding the domain specialities. Hence, we seek to evaluate the performance of our NER system through a detailed case study.

## 4 Case Study

In the absence of an annotated test dataset that closely resembles most real-world scenarios, evaluating the system's performance necessitates user evaluation. In our case study, the NER system is deployed on a German clinical corpus from the cardiovascular domain, where doctor letters undergo anonymization and time-shifting to comply with privacy regulations (Richter-Pechanski et al., 2023). Expert annotators can review and modify the system's output through the web-UI (see Figure 3 and 4), facilitating ongoing refinement based on valuable user feedback within the applied domain. This user-centric evaluation approach ensures that the system is continuously optimized to meet the specific needs and requirements of the clinical context, thereby enhancing its practical utility and effective-

| Physical Object | Conceptual Entity | Procedure | Phenomenon or Process | Health State | Factuality |
|---|---|---|---|---|---|

| T059: Laboratory Procedure | T060: Diagnostic Procedure | T061: Therapeutic or Preventive Procedure | O: null |
|---|---|---|---|

| Semantic Groups | ☑Arterieller | ☑CW-Doppler | ☐-LRB- | ☐Verschlussdruecke | ☐in |
|---|---|---|---|---|---|
| Physical Object | 0 | T074 | 0 | 0 | 0 |
| Conceptual Entity | T034 | T034 | 0 | T081 | T081 |
| Procedure | T060 | T060 | 0 | 0 | 0 |
| Phenomenon or Process | 0 | 0 | 0 | T039 | T039 |
| Health State | 0 | 0 | 0 | 0 | 0 |
| Factuality | 0 | 0 | 0 | 0 | 0 |

‹ 1 2 3 4 5 6 7 **8** 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 ›

Figure 4: UI snippet of annotation revision. Annotators are requested to revise the system-generated annotations by semantic groups. Tokens from one sentence (the 8th sentence in this example) are placed in the header of the revision table, one token per column. The annotators can make changes on the entity types labelled to a single token (selecting one column) or to multiple tokens (selecting multiple columns) by clicking on the most appropriate entity type (*Diagnostic Procedure* in the second row) under the selected semantic group (*Procedure* in the example). Label *O: null* indicates that the token is not recognised as an entity according to the current annotation schema.

ness.

## 4.1 Evaluation Setup

In the human evaluation, we restrict the evaluation scope to the medical texts from two typical sections in the clinic routine, e.g. **Findings** and **Diagnosis**. They indicate different types of medical texts. Table 3 presents the scope of the evaluation. Table 4 shows the most frequent words in different sections of medical texts. These words represent the specialised content of each section. The comparison between these two types of medical texts is presented in each metric.

Three senior medical students, experienced in clinical annotation projects, are the expert annotators in our case study. The annotators work on the same documents. They are instructed to correct system-generated annotations through the standalone user interface shown in Figure 4. The system-generated annotations, which are utilized for revision, are generated using a threshold of 0.5. This threshold is set to avoid many false positives. Since no gold standard test dataset is established, these revisions serve as a form of ground truth to measure the NER system performance on the applied domain.

**Target Data in Case Study.** The most frequent words in texts from different clinic sections (**Findings** and **Diagnosis**) are listed in Table 4. They provide insights into the specialities of each

|  | #Docs | #Sents | #Words |
|---|---|---|---|
| Findings | 8 | 136 | 1562 |
| Diagnosis | 11 | 155 | 1831 |

Table 3: Amount of data, of two different sections, at document-level (column 1), sentence-level (column 2) and word-level (column 3) respectively.

section. The words from **Findings** section are mostly related to the examination and lab results and those from **Diagnosis** section indicate the assessments and patient conditions.

**Metrics.** We measure the inter-rater reliability for each semantic group presented in Figure 5 (Cohen's Kappa[5] and Fleiss's Kappa[6]) to display the degree of agreement among multiple annotators. These scores play a crucial role in how to measure and assess the system's performance since we use the revisions of different annotators as a form of ground truth. In semantic groups with higher agreement, the evaluation metrics, such as Precision, Recall and F-scores presented in Figure 6, are more reliable indicators of the system's NER capabilities. Conversely, discrepancies in F-scores signal challenges in reaching a consensus among annotators of a given semantic group.

[5]https://scikit-learn.org/stable/modules/generated/sklearn.metrics.cohen_kappa_score.

75

| | top 20 frequent words |
|---|---|
| Findings | Kein (*no*), min, ms, QTc, QRS, Sinusrhythmus (*Sinus rhythm*), Ruhe-EKG (*ECG at rest*), Befund (*finding*) , PQ, Nachweis (*proof*), Normofrequenter (*normal frequent*), signifikanten (*significate*), R-Progression, Regelrechte (*regular*), Kammerendteilveraenderungen (*chamber end part changes*), Herzfrequenz (*heart rate*), Linkstyp (*left-side type*), rechts (*right*), regelmaessig (*regularly*), Beurteilung (*assessment*) |
| Diagnosis | Diagnosen (*diagnosis*), rechts (*right*), PTCA, links (*left*), Koronare (*Coronary*), Stenose (*stenosis*), RCA, Pumpfunktion (*pump function*), ED,TTE, LCX, 3-Gefaesserkrankung (*three-vessel disease*), DE-Stentimplantation (*drug eluting coronary implantation*), ohne (*without*), guter (*good*), linksventrikulaerer (*left ventricular*), Erfolgreiche (*successful*), Therapie (*therapy*), Vorhofflimmern (*Atrial fibrillation*), Rekanalisation (*Revascularization*) |

Table 4: The 20 most frequent words (the *English translations*) in texts from **Findings** and **Diagnosis** sections across the applied doctor letters, excluding the stop words (articles, prepositions and numbers).

## 4.2 Analysis

The metrics show that our NER system misses some entities but is relatively confident in the correctness of the labelled entities (higher precision and lower recall scores). Our NER system maintains a consistent performance, as indicated by the average F-score around 0.5 across a variety of the semantic groups with moderate agreement, e.g. **Physical Objects, Conceptual Entity, Procedure, Phenomenon or Process**. Compared to the results for section **Findings**, better system performance is observed for section **Diagnosis** based on the metrics in general. These results indicate that our NER system can provide a certain degree of

Figure 5: Inter-annotator agreement scores of pair-wise Cohen' Kappa and overall agreement of annotators based on multiple Fleiss's Kappa scores across different semantic groups and two types of medical texts.

reliable results in a zero-shot setting.

For the **Health State** group, higher F-scores suggest that the NER system effectively recognizes entities in this group, despite lower inter-annotator agreement scores. It suggests that annotators may agree less on which specific tokens represent the entities within this semantic group, but agree more on the broader context, leading to good performance in terms of F-scores.

The **Factuality** semantic group containing entities such as *Negated, Minor* and *Speculated*, presents a set of specific challenges. The user feedback highlights a notable ambiguity in the annotation process. While the system tends to miss many entities associated with *Minor* and *Speculated*, it is generally effective at capturing instances related to *Negated*. Furthermore, exemplified ambiguities rise in cases like *"no proof for a specific disease"* when deciding whether to annotate the terms *"no"*, *"no proof"* or the entire span with the type *Negated*. This highlights the need for further refinement and specificity in the annotation schema to ensure consistent and unambiguous annotations within the **Factuality** semantic group, as well as collecting

Figure 6: Precision, Recall and F-scores for NER results for medical text from **Findings** (up) and **Diagnosis** (down) section across different semantic groups.

more annotations to fine-tune the model's performance within this group.

### 4.3 System Usability Feedback

We administer a questionnaire to collect comprehensive feedback on the system's usability from the annotators regarding the system usability for future refinement. The answers indicate a System Usability Score (SUS) of 68, reflecting a moderately positive perception of the system's usability. Annotators generally found the annotation revision process straightforward and did not require technical support. However, a common concern was the time-consuming nature of reviewing annotations for every token within each possible semantic group. Annotators also noted that the NER system tends to overlook specific entities and there were some ambiguities related to specific entity types that require further clarifications in the annotation guidelines. These observations align with the evaluation metrics. This invaluable feedback suggests opportunities to enhance the efficiency of the NER system and improve the user experience in further use cases.

### 5 Conclusion

In summary, our work aims to overcome the challenges involved in developing a NER system

for German clinical NLP applications without in-domain training data. We propose using advanced transfer learning methods and focusing on direct adaptability to new datasets with **CRL** and **MSA**. Our case study, which involves close collaboration with domain experts in specific clinical applications, yields invaluable insights that contribute to the overall improvement of the system. These insights are essential for tailoring the system to meet the specific information extraction requirements in target domains. Our work represents a significant advancement in clinical information extraction, alleviating limitations associated with data scarcity and cross-domain transferability. In future research, we will concentrate on incorporating expert feedback into the adaptation pipeline and models' fine-tuning, ultimately creating a continuous learning ecosystem tailored to the distinct clinical context. This effort aligns with our primary objective of advancing NER technology to effectively address challenges related to data scarcity, medical text diversity, and ever-changing label sets.

### Ethical Statement

The annotators involved in the case study are co-authors of this paper and were not compensated for their research contributions.

### Limitations

Importantly, our NER system's adaptability to new datasets to address data scarcity limitations is a key achievement of our research. However, the absence of interactive annotation rounds for resolving disagreements among annotators has prevented the creation of a more refined standard test dataset. In future endeavours, we plan to overcome these limitations by focusing on long-term data collection initiatives aimed at fine-tuning the system and adjusting model weights to better suit specific domains. Due to resource constraints, our ability to conduct comprehensive evaluations across a wider spectrum of clinical domains is restricted. Additionally, we recognize the necessity of exploring additional use cases to expand our understanding of medical text diversity, introduce new entity labels, and enhance the overall robustness of our cross-domain NER system. By addressing these limitations and pursuing a more extensive and diverse set of clinical data, we aim to further elevate the adaptability and utility of our NER system.

## References

Alan R Aronson. 2001. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program. page 17. American Medical Informatics Association.

Alan R Aronson. 2006. Metamap: Mapping text to the umls metathesaurus. *Bethesda, MD: NLM, NIH, DHHS*, 1:26.

Olivier Bodenreider. 2004. The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270.

Florian Borchert, Christina Lohr, Luise Modersohn, Jonas Witt, Thomas Langer, Markus Follmann, Matthias Gietzelt, Bert Arnrich, Udo Hahn, and Matthieu-P. Schapranow. 2022. GGPONC 2.0 - the German clinical guideline corpus for oncology: Curation workflow, annotation policy, baseline NER taggers. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3650–3660, Marseille, France. European Language Resources Association.

John J Dudley and Per Ola Kristensson. 2018. A review of user interface design for interactive machine learning. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 8(2):1–37.

Jerry Alan Fails and Dan R Olsen Jr. 2003. Interactive machine learning. In *Proceedings of the 8th international conference on Intelligent user interfaces*, pages 39–45.

Johann Frei and Frank Kramer. 2021. Gernermed – an open german medical ner model.

Jason Fries, Leon Weber, Natasha Seelam, Gabriel Altay, Debajyoti Datta, Samuele Garda, Sunny Kang, Rosaline Su, Wojciech Kusa, Samuel Cahyawijaya, et al. 2022. Bigbio: A framework for data-centric biomedical natural language processing. *Advances in Neural Information Processing Systems*, 35:25792–25806.

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. *Journal of the American Medical Informatics Association: JAMIA*, 27(1):3.

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.

Madeleine Kittner, Mario Lamping, Damian T Rieke, Julian Götze, Bariya Bajwa, Ivan Jelas, Gina Rüter, Hanjo Hautow, Mario Sänger, Maryam Habibi, et al. 2021. Annotation and initial evaluation of a large annotated german oncological corpus. *JAMIA open*, 4(2):ooab025.

Siting Liang, Mareike Hartmann, and Daniel Sonntag. 2023a. Cross-domain German medical named entity recognition using a pre-trained language model and unified medical semantic types. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 259–271, Toronto, Canada. Association for Computational Linguistics.

Siting Liang, Mareike Hartmann, and Daniel Sonntag. 2023b. Cross-lingual German Biomedical Information Extraction: from Zero-shot to Human-in-the-Loop. *arXiv e-prints*, pages arXiv–2301.

Ignacio Llorca, Florian Borchert, and Matthieu-P. Schapranow. 2023. A meta-dataset of German medical corpora: Harmonization of annotations and cross-corpus NER evaluation. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 171–181, Toronto, Canada. Association for Computational Linguistics.

Hans-Jürgen Profitlich and Daniel Sonntag. 2021. A case study on pros and cons of regular expression detection and dependency parsing for negation extraction from german medical documents. technical report. *CoRR*, abs/2105.09702.

Phillip Richter-Pechanski, Philipp Wiesenbach, Dominic M Schwab, Christina Kiriakou, Mingyang He, Michael M Allers, Anna S Tiefenbacher, Nicola Kunz, Anna Martynova, Noemie Spiller, et al. 2023. A distributable german clinical corpus containing cardiovascular clinical routine doctor's letters. *Scientific Data*, 10(1):207.

Roland Roller, Aljoscha Burchardt, Nils Feldhus, Laura Seiffe, Klemens Budde, Simon Ronicke, and Bilgin Osmanodja. 2022. An annotated corpus of textual explanations for clinical decision support. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2317–2326, Marseille, France. European Language Resources Association.

Guergana K Savova, James J Masanz, Philip V Ogren, Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-Schuler, and Christopher G Chute. 2010. Mayo clinical text analysis and knowledge extraction system (ctakes): architecture, component evaluation and applications. *Journal of the American Medical Informatics Association*, 17(5):507–513.

Henning Schäfer, Ahmad Idrissi-Yaghir, Peter Horn, and Christoph Friedrich. 2022. Cross-language transfer of high-quality annotations: Combining neural

machine translation with cross-linguistic span alignment to apply NER to clinical texts in a low-resource language. In *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 53–62, Seattle, WA. Association for Computational Linguistics.

Isabel Segura-Bedmar, Paloma Martínez, and María Herrero-Zazo. 2013. Semeval-2013 task 9: Extraction of drug-drug interactions from biomedical texts (ddiextraction 2013). In *Second Joint Conference on Lexical and Computational Semantics (\* SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 341–350.

Daniel Sonntag and Hans-Jürgen Profitlich. 2019. An architecture of open-source tools to combine textual information extraction, faceted search and information visualisation. *Artificial intelligence in medicine*, 93:13–28.

Daniel Sonntag, Volker Tresp, Sonja Zillner, Alexander Cavallaro, Matthias Hammon, André Reis, Peter A Fasching, Martin Sedlmayr, Thomas Ganslandt, Hans-Ulrich Prokosch, et al. 2016. The clinical data intelligence project. *Informatik-Spektrum*, 39(4):290–300.

Dominic Widdows, Beate Dorow, and Chiu-Ki Chan. 2002. Using parallel corpora to enrich multilingual lexical resources. pages 240–245.

Hang Yan, Tao Gui, Junqi Dai, Qipeng Guo, Zheng Zhang, and Xipeng Qiu. 2021. A unified generative framework for various NER subtasks. 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 5808–5822, Online. Association for Computational Linguistics.

## A   Entity Types in Training Data

Entity types derived from **Ex4CDS** are presented in Table 5. 134 semantic types from **UMLS** semantic network ontology in 2001 are annotated in **MUCHMORE** corpora. However, the number of annotations of each semantic type is extremely imbalanced ranging from less than 10 terms to at most 8202, see Table 6.

| Entity Type | Description | Corresponding Type Names |
|---|---|---|
| Condition | A pathological medical condition of a patient can describe for instance a symptom or a disease. | Sign or Symptom; Disease or Syndrome; Finding |
| DiagLab | Particular diagnostic procedures have been carried out. | Laboratory Procedure; Diagnostic Procedure |
| LabValues | Mentions of lab values. | Clinical Attribute |
| Measure | Mostly numeric values, often in the context of medications or lab values, but can also be a description if a value changes, e.g. raises. | Quantitative Concept |
| Medication | A medication. | Pharmacologic Substance |
| Process | Describes particular process, such as blood pressure, or heart rate, often related to vital parameters. | Physiologic Function |
| TimeInfo | Describes temporal information, such as 2 weeks ago or January. | Temporal Concept |
| Health State* | A positive condition of the patient. | Healthy Condition |
| Factuality* | Factuality regarding symtoms and diseases (present or not, present but in a lower amount, kind of speculation). | Negated, Minor, Speculated |

Table 5: Entity types, descriptions in Ex4CDS and the corresponding type names (* Type names are matched to the UMLS semantic types except for *HealthState, Factuality*, where no proper semantic type is found and retained the natural words of the entity types).

| ID | Semantic Type | Description | Amount |
|---|---|---|---|
| T101 | Patient or Disabled Group | An individual or individuals classified according to a disability, disease, condition or treatment. | 8202 |
| T047 | Disease or Syndrome | A condition which alters or interferes with a normal process, state, or activity of an organism. It is usually characterized by the abnormal functioning of one or more of the host's systems, parts, or organs. Included here is a complex of symptoms descriptive of a disorder. | 7636 |
| T023 | Body Part, Organ, or Organ Component | A collection of cells and tissues which are localized to a specific area or combine and carry out one or more specialized functions of an organism. This ranges from gross structures to small components of complex organs. These structures are relatively localized in comparison to tissues. | 7070 |
| T169 | Functional Concept | A concept which is of interest because it pertains to the carrying out of a process or activity. | 5569 |
| T061 | Therapeutic or Preventive Procedure | A procedure, method, or technique designed to prevent a disease or a disorder, or to improve physical function, or used in the process of treating a disease or injury. | 5542 |
| T046 | Pathologic Function | A disordered process, activity, or state of the organism as a whole, of a body system or systems, or of multiple organs or tissues. Included here are normal responses to a negative stimulus as well as pathologic conditions or states that are less specific than a disease. Pathologic functions frequently have systemic effects. | 3974 |
| T191 | Neoplastic Process | A new and abnormal growth of tissue in which the growth is uncontrolled and progressive. The growths may be malignant or benign. | 3806 |
| T170 | Intellectual Product | A conceptual entity resulting from human endeavor. Concepts assigned to this type generally refer to information created by humans for some purpose. | 3266 |
| T081 | Quantitative Concept | A concept which involves the dimensions, quantity or capacity of something using some unit of measure, or which involves the quantitative comparison of entities. | 3049 |
| T033 | Finding | That which is discovered by direct observation or measurement of an organism attribute or condition, including the clinical history of the patient. The history of the presence of a disease is a 'Finding' and is distinguished from the disease itself. | 2621 |
| T060 | Diagnostic Procedure | A procedure, method, or technique used to determine the nature or identity of a disease or disorder. This excludes procedures which are primarily carried out on specimens in a laboratory. | 2621 |
| T184 | Sign or Symptom | An observable manifestation of a disease or condition based on clinical judgment, or a manifestation of a disease or condition which is experienced by the patient and reported as a subjective observation. | 2547 |
| T024 | Tissue | An aggregation of similarly specialized cells and the associated intercellular substance. Tissues are relatively non-localized in comparison to body parts, organs or organ components. | 2533 |
| T121 | Pharmacologic Substance | A substance used in the treatment or prevention of pathologic disorders. This includes substances that occur naturally in the body and are administered therapeutically. | 2403 |
| T037 | Injury or Poisoning | A traumatic wound, injury, or poisoning caused by an external agent or force. | 2080 |
| T029 | Body Location or Region | An area, subdivision, or region of the body demarcated for the purpose of topographical description. | 1865 |
| T040 | Organism Function | A physiologic function of the organism as a whole, of multiple organ systems, or of multiple organs or tissues. | 1540 |
| T041 | Mental Process | A physiologic function involving the mind or cognitive processing. | 1429 |
| T078 | Idea or Concept | An abstract concept, such as a social, religious or philosophical concept. | 1309 |
| T032 | Organism Attribute | A property of the organism or its major parts. | 1281 |
| T073 | Manufactured Object | A physical object made by human beings. | 1226 |
| T091 | Biomedical Occupation or Discipline | A vocation, academic discipline, or field of study related to biomedicine. | 1213 |
| T123 | Biologically Active Substance | A generally endogenous substance produced or required by an organism, of primary interest because of its role in the biologic functioning of the organism that produces it. | 1187 |
| T100 | Age Group | An individual or individuals classified according to their age. | 1149 |
| T062 | Research Activity | An activity carried out as part of research or experimentation. | 1148 |
| T079 | Temporal Concept | A concept which pertains to time or duration. | 1124 |

Table 6: Most frequent **UMLS** semantic types annotated in the MUCHMORE corpus. The numbers in the third column are the amount of annotated terms appear in the training data.

# DAIC-WOZ: On the Validity of Using the *Therapist's prompts* in Automatic Depression Detection from Clinical Interviews

**Sergio Burdisso**[*,1], **Ernesto A. Reyes-Ramírez**[2], **Esaú Villatoro-Tello**[*,1],
**Fernando Sánchez-Vega**[2,3], **A. Pastor López-Monroy**[2] and **Petr Motlicek**[1,4]

[1]Idiap Research Institute, Martigny, Switzerland
[2]Mathematics Research Center (CIMAT), Gto, Mexico
[3]Consejo Nacional de Humanidades, Ciencias y Tecnologías (CONAHCYT), México
[4]Brno University of Technology, Brno, Czech Republic

## Abstract

Automatic depression detection from conversational data has gained significant interest in recent years. The DAIC-WOZ dataset, interviews conducted by a human-controlled virtual agent, has been widely used for this task. Recent studies have reported enhanced performance when incorporating interviewer's prompts into the model. In this work, we hypothesize that this improvement might be mainly due to a bias present in these prompts, rather than the proposed architectures and methods. Through ablation experiments and qualitative analysis, we discover that models using interviewer's prompts learn to focus on a specific region of the interviews, where questions about past experiences with mental health issues are asked, and use them as discriminative shortcuts to detect depressed participants. In contrast, models using participant responses gather evidence from across the entire interview. Finally, to highlight the magnitude of this bias, we achieve a 0.90 F1 score by intentionally exploiting it, the highest result reported to date on this dataset using only textual information. Our findings underline the need for caution when incorporating interviewers' prompts into models, as they may inadvertently learn to exploit targeted prompts, rather than learning to characterize the language and behavior that are genuinely indicative of the patient's mental health condition.

## 1 Introduction

Recent advances in Artificial Intelligence (AI) have increased the existing enthusiasm among medical professionals and clinicians when considering the potential for AI-based solutions to make mental healthcare more accessible and to reduce the burden of psychiatric institutions (Passos et al., 2023). This possibility has led some psychiatrists to argue

that the use of AI might result in more standardized and objective measures of mental health (Pendse et al., 2022).

Consequently, the automatic analysis of clinical interviews has been recognized as a promising direction for the development of automatic solutions that will help to improve the diagnostic consistency of depression detection (Tao et al., 2023; Zou et al., 2022; Burdisso et al., 2019; Valstar et al., 2016). The Distress Analysis Interview Corpus - Wizard of Oz (DAIC-WOZ) dataset (Gratch et al., 2014) stands out as the most representative multimodal resource which has been commonly used for training and validating depression classification models within a clinical setup. Most existing studies leverage the participant answers for depressive assessment, varying from single-modality methods, i.e., text transcripts, speech (Burdisso et al., 2023; Villatoro-Tello et al., 2021a; Xezonaki et al., 2020; Mallol-Ragolta et al., 2019), to multi-modal approaches (text + speech + video) (Zhuang et al., 2024; Fang et al., 2023; Shen et al., 2022; Yoon et al., 2022; Villatoro-Tello et al., 2021b). However, recent studies that incorporate therapist's prompts during training, argue that such information works as supplementary context to better extract salient cues from participant answers (Zhuang et al., 2024; Shen et al., 2022; Niu et al., 2021; Dai et al., 2021), reporting high classification performances.

In this paper, we investigate the validity of using the interviewer's prompts from the DAIC-WOZ dataset in automatic depression detection scenarios. We hypothesize that the reported results using both interviewer and participant information may be artificially inflated by a bias induced by the interviewer, failing to generalize to real-world scenarios where such biases may not exist. The impact of over-reporting performance in the DAIC-WOZ dataset has been already pointed by (Bailey and Plumbley, 2021) due to the presence of gender bias. Nevertheless, and to the best of our knowledge, this

---
*Corresponding authors.
{sergio.burdisso, esau.villatoro}@idiap.ch

is the first work to report the existence of a strong bias in the interviewer's prompts and to show that models can effectively exploit it as discriminative shortcuts.

## 2 The DAIC-WOZ Dataset

The DAIC-WOZ dataset contains clinical interviews in North American English, performed by an animated virtual (human-controlled, i.e., Wizard of OZ) interviewer, called Ellie, designed to support the diagnosis of different psychological distress conditions. The DAIC-WOZ stands as a valuable resource frequently utilized by the NLP community, attributed to its rigorous data collection methods and the scarcity of newer data sources exploring comparable phenomena. DAIC-WOZ is a multi-modal corpus, composed by audio and video recordings, and transcribed text from the interviews. To the date, the DAIC-WOZ corpus represents a unique and valuable resource, accumulating over 1K citations since its release.[1]

Ellie conducts semi-structured interviews that are intended to create interactional situations favorable to the assessment of distress indicators correlated with depression, anxiety or post-traumatic stress disorder (PTSD). Theoretically, the advantage of Ellie over a human interviewer is the implicit replicability and consistency of the prompts and accompanying gestures. Thus, Ellie has a finite repertoire of 191 prompts, i.e., general questions (*what are you like when you don't get enough sleep?*), neutral backchannels (*uh huh*), positive empathy (*that's great*), negative empathy (*i'm sorry*), surprise responses (*wow!*), continuation prompts (*could you tell me more about that?*), and miscellaneous prompts ( *don't know; thank you*). Table 1 shows a few statistics from the dataset.[2]

## 3 Methodology

To assess the reliability of using Ellie's prompts for automatic depression detection on DAIC-WOZ, we first examine some of the highest results reported in the recent past using this dataset, summarized in Table 2. We can categorize published works into two primary groups: (a) those using solely the participant (P) responses and, (b) those incorporating Ellie's (E) prompts to the model. It seems that

| Speaker | Partition | Voc. size | Avg. #words | Avg. #tokens |
|---|---|---|---|---|
| Ellie (E) | *train* | 232 | 190.3 ($sd$=26.9) | 567.2 ($sd$=79.10) |
| | *eval* | 216 | 184.8 ($sd$=50.2) | 540.7 ($sd$=148.5) |
| Participant (P) | *train* | 5858 | 621.1 ($sd$=326.2) | 1606.2 ($sd$=893.9) |
| | *eval* | 3268 | 664.2 ($sd$=281.7) | 1756.3 ($sd$=814.7) |

Table 1: DAIC-WOZ contains 107 training files (77 control [C] and 30 depressed [D]), an evaluation set of 35 files (23 [C] and 12 [D]). Table shows the vocabulary size and the average interview length measure in words and *WordPiece* tokens, with its corresponding standard deviation (*sd*) values.

works from group (b) exhibit an overall superior performance compared to those of group (a). To investigate whether this improvement may stem from a bias in Ellie's prompts, before delving into a qualitative analysis, we proposed an initial ablation experiment. Concretely, we evaluated two versions of the same models: one employing only participant responses and another solely using Ellie's prompts. Subsequently, we assess the performance difference between these versions, aiming to quantify the challenge in identifying depressed subjects based on participant responses versus Ellie's prompts. Furthermore, we tested an ensemble approach to measure how complementary these two aspects are to each other.

In particular, we will conduct an ablation experiment using two models: a strong BERT-based baseline model and the Graph Convolutional Network (GCN) model described in Burdisso et al. (2023), which is the best-performing model that relies solely on the participant's text (see Table 2). The choice of these two models aims to compare the baselines against the best-performing model, as well as to analyze models with different natures, namely a bidirectional sequential model and a sequence-agnostic one. Moreover, as will be described below, the GCN model has an attractive interpretability property that we will use in Section 5 for the qualitative analysis. Thus, by analyzing the differences between these two models, we can determine whether the observed patterns hold independently of the model's nature. The models are described as follows:

• **LongBERT:** a BERT-based classification model. More precisely, we used a pre-trained BERT-based Longformer (Beltagy et al., 2020) model with a final linear layer added to classify the input using the encoding of the special *[CLS]* token, following common practice. The choice of using the

---

[1]Rough estimation based on the citation counts of (Gratch et al., 2014; DeVault et al., 2014) in Google scholar.

[2]Labels of the test set are not publicly available due to the AVEC competition (Valstar et al., 2016).

Longformer variant of BERT (Devlin et al., 2019), instead of the standard Transformer (Vaswani et al., 2017) version, stems from the fact that most interviews in DAIC-WOZ are long documents exceeding the 512 token limit (see Table 1).

• **GCN:** The two-layer Graph Convolutional Network (GCN) described in Burdisso et al. (2023) that uses two types of nodes to characterize the interviews: word nodes and participant nodes. In this graph, nodes are represented at three distinct levels: one-hot encoded vectors, embeddings in a latent space (after applying the first convolution), and in a two-dimensional "output space," (after the second convolution) where each dimension corresponds to the probability of belonging to the depression or the control group. Note that since the two type of nodes are represented in the same space, this last learned representation contains probabilities not only for the participants but also for *all the words*. This is an attractive quality of the model that allows us to track down Ellie's bias to particular subset of words and prompts (as described in Section 5).

## 4 Experiments and Results

We trained and evaluated two variants of the GCN: one exclusively using the participant's responses as in the original paper (Burdisso et al., 2023), denoted as *P*-GCN, and another one solely using Ellie's prompts, referred to as *E*-GCN. Similarly, we also fine-tuned and evaluated the same two versions of the Longformer BERT model, referred to as *P-longBERT* and *E-longBERT*, respectively.[3] Table 2 shows the obtained results. When using only the participant responses, *P*-GCN achieved a similarly high F1 score (0.85) to the score reported in the original paper (0.84), and *P*-longBERT a score (0.72) similar to other published works employing solely participant data (e.g. 0.69). On the other hand, when using Ellie, both *E*-GCN and *E-longBERT* achieve comparably higher F1 score. Notably, *E-longBERT*, by simply utilizing Ellie's prompts, managed to achieve the same score (0.84) as the original GCN paper, and the *E*-GCN outperformed all main previously published works that solely rely on textual input, with a score of 0.88. This suggests that when employing Ellie's prompts, the depression and control groups become more easily distinguishable. For instance, the F1

---

[3]Details are provided in Appendix A. Source code to replicate our study available at https://github.com/idiap/bias_in_daic-woz.

| Model | Source | | | $F_1$ score | | |
|-------|--------|---|---|------|---|---|
| | *P* | *E* | *M* | *Avg.* | *D* | *C* |
| Mallol-Ragolta et al. (2019) | ✓ | | | 0.60 | - | - |
| Xezonaki et al. (2020) | ✓ | | | 0.69 | - | - |
| Villatoro-Tello et al. (2021a) | ✓ | | | 0.64 | 0.52 | 0.77 |
| Burdisso et al. (2023) | ✓ | | | **0.84** | **0.80** | **0.89** |
| Williamson et al. (2016) | ✓ | ✓ | | 0.84 | - | - |
| Toto et al. (2021) | ✓ | ✓ | | **0.86** | - | - |
| Shen et al. (2022) | ✓ | ✓ | | 0.83 | - | - |
| Milintsevich et al. (2023) | ✓ | ✓ | | 0.80 | - | - |
| Agarwal and Dias (2024) | ✓ | ✓ | | 0.77 | - | - |
| Niu et al. (2021) | ✓ | ✓ | ✓ | 0.92 | - | - |
| Dai et al. (2021) | ✓ | ✓ | ✓ | **0.96** | - | - |
| Shen et al. (2022) | ✓ | ✓ | ✓ | 0.85 | - | - |
| Zhuang et al. (2024) | ✓ | ✓ | ✓ | 0.88 | 0.85 | 0.91 |
| *P-longBERT* | ✓ | | | 0.72 | 0.64 | 0.80 |
| *E-longBERT* | | ✓ | | **0.84** | **0.80** | **0.89** |
| *P-longBERT* ∧ *E-longBERT* | ✓ | ✓ | | 0.79 | 0.70 | 0.88 |
| *P*-GCN | ✓ | | | 0.85 | 0.81 | 0.88 |
| *E*-GCN | | ✓ | | **0.88** | **0.85** | **0.91** |
| *P*-GCN ∧ *E*-GCN | ✓ | ✓ | | **0.90** | **0.87** | **0.94** |

Table 2: Main previously published results on DAIC-WOZ evaluation set along with our obtained results. Performance is reported in terms of the $F_1$ score for both control (*C*) and depression (*D*) classes, as well as their macro average (*Avg.*). Results are marked with the source data used: (P) and (E) text from the participant and Ellie; (M) multimodal, e.g., speech and video. The global-best result among models using only textual content is **<u>underlined</u>**, while the best results in each group is highlighted in **bold**.

score of the *longBERTs* for the depression group (D) improves from 0.64 to 0.80 when using Ellie's prompts.

Finally, we performed a simple voting ensemble between the two variants of each model, denoted using the "and" symbol (∧). Participants are classified as positive (*i.e.*, in the depression group) only when both variants, Ellie *and* Participant, classify them as positive. As shown in Table 2, the ensemble approach enables the GCN-based model to achieve a remarkable F1 score of 0.90, the highest reported score to date among models exclusively utilizing textual content. These results suggest that the integration of both Ellie and participant content could be complementary for certain models, further exploiting Ellie's bias to make the depression and control groups even more easily distinguishable.

## 5 Analysis and Discussion

Overall, experimental results suggest that Ellie's prompts contain information that the models can exploit to more easily classify the participants. This

Figure 1: Heatmaps illustrating the distribution of learned keywords by each model across the progression of each interview. The x-axis represents individual interviews, while the y-axis denotes the percentage of the conversation from the beginning (0%) to the end (100%). The white vertical line in each plot indicates the training and evaluation splits respectively. Finally, in the *E*-GCN evaluation split region, the small red rectangle depicts the interview segment showed in Fig. 2.

is reasonable when considering that therapists adjust their questioning patterns based on the subjects' responses and may adapt their inquiries to delve deeper into specific aspects *when detecting potential depressive symptoms.*

To explore this possibility further, as mentioned in Section 3, we leveraged the GCN-based model's ability to learn a common representation for both participant and word nodes in the same output space. Firstly, we extracted the words that both GCN models learned to use to identify the depressed group, which we will refer to as keywords.[4] Subsequently, we analyzed the distribution of these keywords throughout the progression of each interview to contrast the depressed group against the control group, allowing us to visualize how easily distinguishable the two groups are from the perspectives of both Ellie (*E*-GCN) and the participant (*P*-GCN) models. Figure 1 illustrates the distributions obtained from our analysis, highlighting the contrasting behavior of the *E*-GCN and *P*-GCN models. The *P*-GCN distribution exhibits variability across interviews, with no distinct pat-

tern emerging from the distribution of keywords. In contrast, the *E*-GCN model displays a clear and consistent pattern, with concrete regions where keywords concentrate. That is, the participant model gathers evidence from various parts of the conversations, whereas Ellie's model focuses mainly on very specific segments, *i.e.* specific questions, to classify the participants. Furthermore, by contrasting the distributions for the depressed group against the control group, we observe that it is easier to distinguish between them using *E*-GCN than *P*-GCN. This suggests that Ellie's keywords are not only more localized but also possess greater discriminatory power. Note that for *E*-GCN, in contrast with the control group, almost all the interviews in the depressed group have colored regions, and they are mostly concentrated in a single segment that appears *after halfway the interviews.*[5] Interestingly, most of these segments correspond to a phase in the interview where Ellie begins to ask more personal questions about past experiences with mental

---

[4]Words $w$ such $P(depressed \mid w) > P(\neg depressed \mid w)$

[5]As shown in Table A2, to validate this observation further, we fine-tuned *E-longBERT* on the second half of interviews, achieving 0.84 F1 (same as full interviews). Using only the first half dropped F1 to 0.60, highlighting the importance of this latter portion.

Figure 2: Illustrative segment from interview "381" in the evaluation set, highlighted in Figure 1. Conversation turns are color-coded based on the proportion of keywords present, with keywords underlined for emphasis.

health issues. Figure 2 shows one such segment. Here, we see the segment containing the only four questions that Ellie's model used to classify the participant, disregarding everything else in the conversation, including the question "*Have you been diagnosed with depression?*" Note that such questions may be asked to different participants, but an affirmative answer triggered Ellie to delve deeper into specific questions, questions that models could easily learned to identify and exploit to correctly classify the participants.

### 5.1 Implications in Clinical Practice

In clinical practice the final psychiatric diagnosis is typically determined through a clinical interview, often semi-structured, where rating scales serve as additional sources of information to aid in diagnosis. However, these rating scales have limitations, as responses can be influenced by factors such as the patient's emotional state, comorbidities, relationship with the clinician, and patient self-bias (e.g., participants may be more likely to exaggerate their symptoms (Mao et al., 2023)).

Accordingly, the final goal of screening tools such as Ellie, is to contribute towards the replicability, consistency, standardization and the construction of objective measures that support the diagnosis of different mental disorders (Pendse et al., 2022).

As shown, the overall analysis described in this paper uncovers interesting biases in the data and shows how ostensibly good performance of NLP models can be deceiving and stress the importance of paying attention to the data and the rationales of the models rather than simply focusing on the superficial performance numbers. Thus, for automatic depression detection systems to be applicable in real-life clinical practice, systems must be able to provide practitioners whit interpretable and transparent insights to validate systems decisions. There are complex interactions happening during a clinical interview, and accurately modeling is still an open challenge, highlighting the need to develop robust and ethical AI systems for this important and sensitive application domain.

## 6 Conclusions

Our analysis reveals that the prompts posed by the interviewer, Ellie, contain biases that allow models to more easily distinguish between depressed and control participants in the DAIC-WOZ dataset. By analyzing the keywords learned by the models, we discover that Ellie's model tends to focus on highly localized segments of the interviews, primarily concentrated in the latter portion where more personal mental health questions are asked. In contrast, the model using participant responses alone does not exhibit such localization, instead gathering evidence from across the entire conversations. More broadly, our findings underline the need for caution when incorporating interviewers' prompts into mental health diagnostic models. Interviewers often strategically adapt their questioning to probe for potential symptoms. As a result, models may learn to exploit these targeted prompts as discriminative shortcuts, rather than learning to characterize the language and behavior that are truly indicative of mental health conditions.

## 7 Ethical Considerations

In this section, we elaborate on the potential ethical issues.

1. **Data privacy, participant demographics, and consent.** All the experiments reported in this paper were made on the publicly available DAIC-WOZ dataset, a valuable resource used for training and validating depression detection systems from clinical interviews.

This particular dataset was collected by the Institute for Creative Technologies at the University of Southern California. According to the original paper, the DAIC-WOZ dataset received approval from Institutional Ethics Board. All the participants, including the U.S. armed forces veterans and general public from the Greater Los Angeles metropolitan area, were informed that their interviews will be used for academic purposes. All personal details like names, ages, and professions are either removed or anonymized, eliminating any risk of personal information exposure. Original videos from the interviews are not provided, but instead vector features of facial actions and eye gaze are given, making it impossible to reconstruct the participants' appearance. In general, the information of participants was rigorously protected.

2. **The role of AI-based diagnosis.** Our performed experiments aimed at highlighting the importance of using interpretable AI-based solutions as an assistant tools. Thus, the goal is not to replace human experts (psychologists and psychiatrists) but to develop systems that should be used only as support tools. The principle of leaving the decision to the machine would imply major risks for decision making in the health field, a mistake that in high-stakes healthcare settings could prove detrimental or even dangerous. The experiments reported in this paper represent a step forward on the development of bias-aware models in the context of clinical interviews analysis.

## 8 Limitations

In this section we discuss the limitations of the study described in this paper.

1. **Task configuration.** In this paper we only focused on the task of depression detection from clinical interviews, i.e., a controlled scenario where a mental health expert (therapist) conducts an interview with the goal to identify different psychological distress conditions present in the interviewed participant. This setup is significantly different from the so called "wild setting", which refers to the analysis of daily messages, e.g., social media posts. Thus, the findings and claims made in this paper are limited to a clinical setup, and might

not be applicable to different setups. As part of our future work, we plan to validate the impact of prompts generated by a fully automatic therapist in similar setups, in particular in the E-DAIC (DeVault et al., 2014) corpus.

2. **Corpus and modality specific.** Our study is limited to textual modality present in the DAIC-WOZ corpus. Given that the acoustic modality contains also Ellie's interventions, we would like to confirm the presence of the same bias in the acoustic modality. Thus, as part of our future work, we plan to extend our analysis to the additional modalities present in the selected corpus. Similarly, our findings apply specifically to the DAIC-WOZ corpus, hence we cannot confirm the presence of the same type biases in similar corpora. As part of our immediate work, we will replicate our analysis with other datasets like E-DAIC (DeVault et al., 2014), EATD (Shen et al., 2022), or the recently released ANDROIDS (Tao et al., 2023) dataset.

## Acknowledgements

## References

Navneet Agarwal and Gaël Dias. 2024. Analysing Relevance of Discourse Structure for Improved Mental Health Estimation. In *9th Workshop on Computational Linguistics and Clinical Psychology (CLPSYCH) associated to 18th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, Saint Julian, Malta.

Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2623–2631.

Andrew Bailey and Mark D Plumbley. 2021. Gender bias in depression detection using audio features. In *2021 29th European Signal Processing Conference (EUSIPCO)*, pages 596–600. IEEE.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv:2004.05150*.

Sergey Brin and Lawrence Page. 1998. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7):107–117.

Sergio Burdisso, Marcelo Luis Errecalde, and Manuel Montes y Gómez. 2019. Towards measuring the severity of depression in social media via text classification. In *XXV CACIC*, pages 577–588.

Sergio Burdisso, Esaú Villatoro-Tello, Srikanth Madikeri, and Petr Motlicek. 2023. Node-weighted Graph Convolutional Network for Depression Detection in Transcribed Clinical Interviews. In *Proc. INTERSPEECH 2023*, pages 3617–3621.

Ilias Chalkidis, Xiang Dai, Manos Fergadiotis, Prodromos Malakasiotis, and Desmond Elliott. 2022. An exploration of hierarchical attention transformers for efficient long document classification. *arXiv preprint arXiv:2210.05529*.

Zhijun Dai, Heng Zhou, Qingfang Ba, Yang Zhou, Lifeng Wang, and Guochen Li. 2021. Improving depression prediction using a novel feature selection algorithm coupled with context-aware analysis. *Journal of affective disorders*, 295:1040–1048.

David DeVault, Ron Artstein, Grace Benn, Teresa Dey, Ed Fast, Alesia Gainer, Kallirroi Georgila, Jon Gratch, Arno Hartholt, Margaux Lhommet, et al. 2014. Simsensei kiosk: A virtual human interviewer for healthcare decision support. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*, pages 1061–1068.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Ming Fang, Siyu Peng, Yujia Liang, Chih-Cheng Hung, and Shuhua Liu. 2023. A multimodal fusion model with multi-level attention mechanism for depression detection. *Biomedical Signal Processing and Control*, 82:104561.

Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Stratou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-Philippe Morency. 2014. The distress analysis interview corpus of human and computer interviews. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3123–3128, Reykjavik, Iceland. European Language Resources Association (ELRA).

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Adria Mallol-Ragolta, Ziping Zhao, Lukas Stappen, Nicholas Cummins, and Björn W. Schuller. 2019. A Hierarchical Attention Network-Based Approach for Depression Detection from Transcribed Clinical Interviews. In *Proc. Interspeech 2019*, pages 221–225.

Kaining Mao, Yuqi Wu, and Jie Chen. 2023. A systematic review on automated clinical depression diagnosis. *npj Mental Health Research*, 2(1):20.

Kirill Milintsevich, Kairit Sirts, and Gael Dias. 2023. Towards automatic text-based estimation of depression through symptom prediction. *Brain Informatics*, 10.

Meng Niu, Kai Chen, Qingcai Chen, and Lufeng Yang. 2021. Hcag: A hierarchical context-aware graph attention model for depression detection. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4235–4239.

Ives Cavalcante Passos, Francisco Diego Rabelo-da Ponte, and Flavio Kapczinski. 2023. *Digital mental health: a practitioner's guide*. Springer.

Sachin R Pendse, Daniel Nkemelu, Nicola J Bidwell, Sushrut Jadhav, Soumitra Pathare, Munmun De Choudhury, and Neha Kumar. 2022. From treatment to healing: Envisioning a decolonial digital mental health. In *CHI Conference on Human Factors in Computing Systems*, pages 1–23.

Ying Shen, Huiyu Yang, and Lin Lin. 2022. Automatic depression detection: An emotional audio-textual corpus and a gru/bilstm-based model. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6247–6251. IEEE.

Fuxiang Tao, Anna Esposito, and Alessandro Vinciarelli. 2023. The Androids Corpus: A New Publicly Available Benchmark for Speech Based Depression Detection. In *Proc. INTERSPEECH 2023*, pages 4149–4153.

Ermal Toto, ML Tlachac, and Elke A. Rundensteiner. 2021. Audibert: A deep transfer learning multimodal classification framework for depression screening. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, CIKM '21, page 4145–4154, New York, NY, USA. Association for Computing Machinery.

Michel Valstar, Jonathan Gratch, Björn Schuller, Fabien Ringeval, Denis Lalanne, Mercedes Torres Torres, Stefan Scherer, Giota Stratou, Roddy Cowie, and Maja Pantic. 2016. Avec 2016: Depression, mood, and emotion recognition workshop and challenge. In *Proceedings of the 6th international workshop on audio/visual emotion challenge*, pages 3–10. ACM.

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.

Esaú Villatoro-Tello, Gabriela Ramírez-de-la Rosa, Daniel Gática-Pérez, Mathew Magimai.-Doss, and Héctor Jiménez-Salazar. 2021a. Approximating the mental lexicon from clinical interviews as a support tool for depression detection. In *Proc. ICMI'21*, page 557–566.

Esaú Villatoro-Tello, S. Pavankumar Dubagunta, Julian Fritsch, Gabriela Ramírez de-la Rosa, Petr Motlicek, and Mathew Magimai-Doss. 2021b. Late Fusion of the Available Lexicon and Raw Waveform-Based Acoustic Modeling for Depression and Dementia Recognition. In *Proc. Interspeech 2021*, pages 1927–1931.

James R Williamson, Elizabeth Godoy, Miriam Cha, Adrianne Schwarzentruber, Pooya Khorrami, Youngjune Gwon, Hsiang-Tsung Kung, Charlie Dagli, and Thomas F Quatieri. 2016. Detecting depression using vocal, facial and semantic communication cues. In *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, pages 11–18.

Danai Xezonaki, Georgios Paraskevopoulos, Alexandros Potamianos, and Shrikanth S. Narayanan. 2020. Affective conditioning on hierarchical attention networks applied to depression detection from transcribed clinical interviews. In *Interspeech*.

Jeewoo Yoon, Chaewon Kang, Seungbae Kim, and Jinyoung Han. 2022. D-vlog: Multimodal vlog dataset for depression detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12226–12234.

Chen Zhuang, Deng Jiawen, Zhou Jinfeng, Wu Jincenzi, Qian Tieyun, and Minlie Huang. 2024. Depression detection in clinical interviews with LLM-empowered structural element graph. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.

Bochao Zou, Jiali Han, Yingxue Wang, Rui Liu, Shenghui Zhao, Lei Feng, Xiangwen Lyu, and Huimin Ma. 2022. Semi-structural interview-based chinese multimodal depression corpus towards automatic preliminary screening of depressive disorders. *IEEE Transactions on Affective Computing*.

| Model | Learning Rate | Epoch | Features | Macro F$_1$ |
|-------|---------------|-------|----------|-------------|
| P-GCN | 1.022e-06 | 10 | *top-250* | 0.85 |
| E-GCN | 1.124e-06 | 10 | *auto* | 0.88 |

Table A1: Best hyperparameters obtained for the GCN models after optimization along with the obtained macro averaged F$_1$ score.

## A   Technical details

### A.1   Graph Convolutional Network

A Graph Convolutional Network (GCN) is a multi-layer neural network that operates directly on a graph and induces embedding vectors of nodes based on the properties of their neighbors. In this work we use the inductive two-layer GCN described in Burdisso et al. (2023). Let $A \in \mathcal{R}^{n \times n}$ be the weighted adjacency matrix of the graph connecting words and interviews of the DAIC-WOZ training set, the GCN is defined as:

$$H^{(1)} = \sigma(\tilde{A}H^{(0)}W^{(0)}) \tag{1}$$

$$Z = \text{softmax}(\tilde{A}H^{(1)}W^{(1)}) \tag{2}$$

where $\tilde{A} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ represents the normalized symmetric adjacency matrix, $W^{(0)}$ is the learned node embeddings lookup table, and $W^{(1)}$ represents the learned weight matrix in the second layer. Loss is computed by means of the cross-entropy between $Z_i$ and the one-hot encoded ground truth label $Y_i$ for all $i$-th interview in the training set. Following the original paper, we set $k = 64$ for the $k$-dimensional feature matrix $H^{(1)} \in \mathcal{R}^{n \times k}$. The adjacency matrix is defined as follows:

$$A_{ij} = \begin{cases} mi(i,j) & \text{if } i, j \text{ are words \& } mi(i,j) > 0 \\ pr(i,j) & \text{if } i, j \text{ are words \& } i = j \\ \text{tf-idf}_{i,j} & \text{if } i \text{ is interview \& } j \text{ is word} \\ 0 & \text{otherwise} \end{cases} \tag{3}$$

where *mi* is the point-wise mutual information and *pr* the *PageRank* (Brin and Page, 1998) score for node $i$.

Finally, in Section 5 we extracted all the words that the model learned to associate to the depressed category. To select these keywords we selected all words $i$ such that $P(depressed \mid word_i) > P(control \mid word_i)$, that is, $keywords = \{word_i \mid Z_{i,depressed} > 0.5\}$.

| Model | Learning Rate | Epoch | Macro $F_1$ |
|-------|--------------|-------|-------------|
| P-longBERT | 2.497e-03 | 10 | 0.72 |
| *first half* | 1.352e-03 | *10* | *0.67* |
| *second half* | 6.051e-03 | *10* | *0.73* |
| E-longBERT | 1.044e-03 | 6 | 0.84 |
| *first half* | 8.209e-04 | *9* | *0.60* |
| *second half* | 5.075e-04 | *7* | *0.84* |

Table A2: Best hyperparameters obtained for the long-BERT models after optimization along with the obtained macro averaged $F_1$ score.

## A.2 Longformer BERT

The Longformer (Beltagy et al., 2020) replaces the quadratic self-attention mechanism of Transformers (Vaswani et al., 2017) with a combination of global and local windowed attention, scaling linearly with sequence length. This modification enables efficient processing of documents with thousands of tokens, consistently outperforming Transformer-based models on long document tasks. In particular, we used the version of Longformer described in Chalkidis et al. (2022) which has been warm-started re-using the weights of BERT, and continued pre-trained for MLM following the paradigm described in the original Longformer paper. This pre-trained model is available in Hugging Face at `https://huggingface.co/kiddothe2b/longformer-mini-1024`.

## A.3 Implementation details

All models were implemented using PyTorch and were optimized using *Optuna* (Akiba et al., 2019) with 100 trials for hyperparameter search maximizing the macro averaged F1 score. In each trail, models were trained using AdamW (Loshchilov and Hutter, 2019) optimizer ($\beta_1{=}0.9, \beta_2{=}0.999, \epsilon{=}1e{-}8$) with *learning rate* and number of epochs $n$ searched in $\gamma \in [1e{-}7, 1e{-}3]$ and $n \in [1, 10]$, respectively. In addition, for GCN, the optimization also tried the three feature selection techniques described in the original paper, *auto*, *top-k*, *none* for, respectively, automatic selection based on term weights learned using Logistic Regression, top-$k$ best selection based on *ANOVA F-value* between words and labels with $k \in \{100, 250, 500, 1000, 1500\}$, and no feature selection (full vocabulary). Best obtained hyperparameters for the GCN models are shown in Table A1. Finally, Table A2 presents the parameters obtained for the *longBERT* models,

along with the results of the complementary ablation experiments mentioned at the end of Section 5. Specifically, we divided each interview into two equal parts and performed fine-tuning and evaluation using either the first or the second half. The objective was to reinforce our conclusions regarding the existence of a bias, particularly in the second half of the interviews, as detected by the keywords from the GCN model (Figure 1).

# Parameter-Efficient Fine-Tuning of LLaMA for the Clinical Domain

**Aryo Pradipta Gema**[1]     **Pasquale Minervini**[1]     **Luke Daines**[2]
**Tom Hope**[3,4]     **Beatrice Alex**[5,6]

[1]School of Informatics, University of Edinburgh     [2]Usher Institute, University of Edinburgh
[3]Allen Institute of AI
[4]Hebrew University of Jerusalem
[5]Edinburgh Futures Institute, University of Edinburgh
[6]School of Literatures, Languages and Cultures, University of Edinburgh
{aryo.gema, p.minervini, luke.daines, b.alex}@ed.ac.uk
tomh@allenai.org

## Abstract

Adapting pretrained language models to novel domains, such as clinical applications, traditionally involves retraining their entire set of parameters. Parameter-Efficient Fine-Tuning (PEFT) techniques for fine-tuning language models significantly reduce computational requirements by selectively fine-tuning small subsets of parameters. In this study, we propose a two-step PEFT framework and evaluate it in the clinical domain. Our approach combines a specialised PEFT adapter layer designed for clinical domain adaptation with another adapter specialised for downstream tasks. We evaluate the framework on multiple clinical outcome prediction datasets, comparing it to clinically trained language models. Our framework achieves a better AUROC score averaged across all clinical downstream tasks compared to clinical language models. In particular, we observe large improvements of 4-5% AUROC in large-scale multilabel classification tasks, such as diagnoses and procedures classification. To our knowledge, this study is the first to provide an extensive empirical analysis of the interplay between PEFT techniques and domain adaptation in an important real-world domain of clinical applications.[1]

## 1  Introduction

Large Language Models (LLMs) have consistently achieved state-of-the-art performance across various NLP tasks. However, while these models exhibit impressive generalisation abilities, they often struggle to perform in specialised domains such as clinical applications, primarily due to the absence of domain-specific knowledge. The complexity of medical terminology and the presence of incomplete sentences in clinical notes contribute to this challenge (Lehman and Johnson, 2023). Unfortunately, studies have indicated that even LLMs



Figure 1: An illustration of the proposed two-step PEFT framework. Clinical LLaMA-LoRA fine-tunes the pretrained LLaMA to the clinical domain. Downstream LLaMA-LoRA further fine-tunes the domain-adapted model to downstream clinical tasks.

pretrained with datasets comprising biomedical publications still exhibit suboptimal performance when applied to downstream clinical applications, particularly when compared to LLMs pretrained with clinical notes (Alsentzer et al., 2019; Li et al., 2022; Yang et al., 2022). This observation suggests that there are intrinsic nuances specific to the clinical context that can only be effectively captured if LLMs undergo pretraining using clinical datasets.

The current approach of adapting pretrained LLMs to the clinical domain typically involves fine-tuning the entire model parameters (Alsentzer et al., 2019; Peng et al., 2019; van Aken et al., 2021; Michalopoulos et al., 2021; Lehman and Johnson, 2023). However, due to the rapid increase in the size of LLMs, such a practice demands extensive computational resources, which may not be readily accessible to all researchers. Consequently, this challenge will further exacerbate the disparity between the resource-rich and resource-constrained research institutions (Ruder et al., 2022).

To address the substantial computational demands, studies have proposed various Parameter-

---

[1]The code is accessible via https://github.com/aryopg/clinical_peft.

Efficient Fine-Tuning (PEFT) techniques. These techniques present a practical solution by fine-tuning a small subset of additional parameters while keeping the remaining pretrained parameters fixed. As a result, this strategy significantly alleviates the computational burden while achieving comparable performance to that of full fine-tuning.

In this study, we propose a two-step PEFT framework (see Figure 1). Firstly, we introduce Clinical LLaMA-LoRA, a Low-Rank Adaptation (LoRA, Hu et al., 2022) PEFT adapter built upon the open-source Large Language Model Meta AI (LLaMA) (Touvron et al., 2023). Then, we introduce Downstream LLaMA-LoRA, which is trained on top of the pretrained Clinical LLaMA-LoRA. Downstream LLaMA-LoRA is specifically designed for clinical downstream tasks. The fusion of the two adapters achieves better performance in clinical NLP downstream tasks compared to clinically trained LLMs while considerably reducing the computational requirements. This study presents the following contributions:

- We introduce Clinical LLaMA-LoRA, a PEFT-adapted version of the LLaMA model tailored specifically for the clinical domain.

- We provide comparisons of multiple PEFT techniques in terms of language modelling performance based on perplexity score, shedding light on the optimal PEFT techniques for the clinical domain-adaptive pretraining.

- We introduce Downstream LLaMA-LoRA, built on top of Clinical LLaMA-LoRA and tailored specifically for the clinical downstream tasks.

- We evaluate the proposed mixture of Clinical LLaMA-LoRA and Downstream LLaMA-LoRA on downstream clinical datasets and tasks. Our proposed framework showcases improvements in AUROC scores over the existing clinical LLMs.

## 2 Background

### 2.1 Biomedical Large Language Models

General-domain LLMs continue to face challenges when confronted with domain-specific tasks. The complexity associated with the requisite domain knowledge is recognised as a significant factor (Ling et al., 2023), particularly within the biomedical domain. Consequently, numerous studies have attempted to adapt LLMs specifically for the biomedical domain.

An early example of such adaptation is BioBERT (Lee et al., 2019), which was pretrained using biomedical research articles from PubMed and PubMed Central. This adaptation has shown improved performance across various biomedical NLP tasks. Recognising the significance of biomedical-specific vocabularies, Gu et al. (2022) proposed PubMedBERT, which is pretrained on biomedical data from scratch and initialised the model vocabulary with the biomedical corpus. The growing interest in biomedical NLP research has led to the adaptation of even larger models to the biomedical domain (Luo et al., 2022; Singhal et al., 2022; Wu et al., 2023; Singhal et al., 2023)

While these biomedical LLMs have demonstrated advancements in various biomedical NLP benchmarking tasks, studies have revealed that clinical LLMs still outperform their biomedical counterparts in numerous clinical downstream tasks (Alsentzer et al., 2019; Yang et al., 2022; Li et al., 2022; Lehman and Johnson, 2023). This suggests that domain-adaptive pretraining using clinical data is still the *de facto* protocol in adapting LLMs to the clinical domain.

### 2.2 Clinical Large Language Models

Clinical LLMs are often fine-tuned with clinical data from an LLM that is already pretrained with datasets that encompass broader topics. For instance, Bio+ClinicalBERT (Alsentzer et al., 2019) is domain-adaptively pretrained using clinical notes from the Medical Information Mart for Intensive Care (MIMIC)-III database (Johnson et al., 2016), starting from a pretrained BioBERT (Lee et al., 2019), which itself is pretrained on biomedical articles. BlueBERT (Peng et al., 2019) is domain-adaptively pretrained using PubMed abstracts and MIMIC-III clinical notes from a BERT model (Devlin et al., 2019), that is pretrained with general-domain texts. Similarly, Clinical-T5 (Lehman and Johnson, 2023) is domain-adaptively pretrained using the union of MIMIC-III and MIMIC-IV (Johnson et al., 2023) clinical notes from T5-base (Raffel et al., 2020), another general-domain LLM.

All these studies share a common approach, which is to fine-tune the entire model parameters. With massive LLMs, this method has become cost-prohibitive and inaccessible for many researchers.

Figure 2: Frameworks of domain-adaptive and downstream fine-tuning to adapt a pretrained LLM from the general domain to the clinical domain. As opposed to a full fine-tuning process which can be prohibitively expensive (left), our approach leverages PEFT techniques to introduce a clinically-specialised adapter that is attached to a pretrained general LLM (right). Our proposed framework also introduces another clinical PEFT adapter trained on the downstream clinical tasks, such as clinical note classification.

## 2.3 Parameter-Efficient Fine-Tuning for Large Language Models

Suppose that we have a pretrained LLM $P_\Phi(y|x)$; fine-tuning it can be effectively defined as finding the most appropriate parameter changes $\Delta\Phi$ by optimising the fine-tuning objective. A conventional, full fine-tuning process means that the model needs to learn a $\Delta\Phi$ whose dimension is equal to the entire parameters of the pretrained LLM $|\Delta\Phi| = |\Phi_0|$, which is computationally expensive. PEFT techniques address this by tuning the *delta* $\Delta\Phi$, which corresponds to a very small fraction of additional trainable parameters during the fine-tuning process.

Adapter tuning (Houlsby et al., 2019) is an early PEFT method that involves adding small additional parameters called *adapters* to each layer of the pretrained model and strictly fine-tuning this small set of new parameters. LoRA (Hu et al., 2022) is another PEFT approach that trains low-rank matrices to represent the attention weights update of transformer-based models.

Another group of PEFT approaches leverages the concept of prompting. Prefix Tuning (Li and Liang, 2021) optimises a sequence of continuous task-specific vectors, called a *prefix*, which are trainable parameters that do not correspond to real tokens. P-Tuning (Liu et al., 2021b) uses a similar strategy as Prefix tuning with a focus on text understanding tasks, as opposed to generative tasks. Prompt tuning (Lester et al., 2021) simplifies Prefix tuning by introducing trainable tokens, called *soft prompts*, for each downstream task. Liu et al.

(2021a) introduced P-tuning v2 which uses deep prompt tuning to address the lack of performance gain in the previous prompt tuning techniques.

By fine-tuning a small fraction of additional parameters, all PEFT approaches alleviate the issue of extensive computational resource requirements.

## 2.4 Multi-step Adaptation

Prior studies have explored the two-step adaptation framework, although they have fundamental differences from our proposed setup. For instance, Zhang et al. (2021) introduced a multi-domain unsupervised domain adaptation (UDA) with a two-step strategy, involving domain-fusion training with Masked Language Model loss on a mixed corpus, followed by task fine-tuning with a task-specific loss on the domain corpus. More recently, Malik et al. (2023) introduced UDApter which utilises PEFT adapters to do efficient UDA. However, unsupervised domain matching techniques such as UDApter rely on restrictive assumptions about the underlying data distributions that are often unsatisfied in real-world scenarios (Li et al., 2020). In our study, we experiment with the clinical domain as the target domain that is not available in the LLM's initial pretraining. Consequently, significant discrepancies exist between the distributions of the source and target domains. Leveraging the amount of available clinical notes, we adopt a self-supervised learning paradigm by continually pretraining the LLMs within the target domain rather than relying on the UDA paradigm.

Our approach shares theoretical similarities with the multi-step continual pretraining approach, pro-

posed by Gururangan et al. (2020), which proposes domain- and task-adaptive pretraining. However, the main difference between our proposed approach and Gururangan et al. (2020) is in the discrepancy between the source and the target domains. Gururangan et al. (2020) experimented with adapting general-domain LLMs to domains encountered during their initial pretraining, such as news and biomedical domains. On the other hand, we experiment with the clinical domain which is entirely absent from the LLMs' initial pretraining due to legal constraints which restrict access to sensitive clinical notes. On top of that, adapting to the clinical domain poses a bigger challenge due to the complexity of medical terminology and the presence of incomplete sentences (Lehman et al., 2023).

## 3 Methodology

### 3.1 Problem Statement

Figure 2 shows the comparison between the current and proposed problem definitions. The general problem can be decomposed into two stages:

**Domain-adaptive Pretraining.** Given a pretrained general LLM $P_\Phi(y|x)$ with its parameters $\Phi$ and a training dataset $\mathcal{Z} = \{(x_i, y_i)\}_{i=1,...,N}$. To adapt to the new domain, the model needs to update its weight iteratively from its pretrained state $\Phi_0$ to $\Phi = \Phi_0 + \Delta\Phi$. This process of maximising the objective function can be defined as:

$$\underset{\Phi}{\mathrm{argmax}} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log\left(P_\Phi\left(y_t \mid x, y_{<t}\right)\right)$$

In the current paradigm, a full fine-tuning process means that the model needs to learn a $\Delta\Phi$ whose dimension is equal to the entire pretrained parameters $|\Delta\Phi| = |\Phi_0|$, which is computationally expensive.

In the proposed paradigm, we tune only small additional parameters $\theta$ such that $\Phi = \Phi_0 + \Delta\Phi(\theta)$ whose dimension is very small compared to the original parameters $|\theta| \ll |\Phi_0|$. Thus, the training objective can be redefined as:

$$\underset{\theta}{\mathrm{argmax}} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log\left(P_{\Phi_0 + \Delta\Phi(\theta)}\left(y_t \mid x, y_{<t}\right)\right)$$

In the current paradigm, the outcome of domain-adaptive pretraining would be a clinically-adapted LLM. While in the proposed paradigm, the outcome would be the clinical PEFT component, which can be combined with the untouched pretrained general LLM for downstream applications.

**Downstream Fine-tuning.** In the current paradigm, the pretrained clinical LLM is fine-tuned to the downstream tasks, such as document classification tasks. Suppose that we have a pretrained clinical LLM $P_{\Phi,\Theta}$ with its domain-adapted parameters $\Phi$ and a newly initialised classifier layer $\Theta$, as well as a training dataset $\mathcal{Z} = \{(x_i, y_i)\}_{i=1,...,N}$. We want to maximise a specific loss function, such as a cross-entropy loss:

$$\underset{\Phi,\Theta}{\mathrm{argmax}} \frac{1}{N} \sum_{i=1}^{N} y_i \log\left(P_{\Phi,\Theta}\left(x_i\right)\right)$$

In contrast, in the proposed paradigm, the fine-tuning process only updates the small additional parameters $\Delta\Phi(\theta)$ and the classifier head $\Theta$:

$$\underset{\theta,\Theta}{\mathrm{argmax}} \frac{1}{N} \sum_{i=1}^{N} y_i \log\left(P_{\Phi+\Delta\Phi(\theta),\Theta}\left(x_i\right)\right)$$

In fact, we can also decompose the fine-tuning into an additional "delta-updating" process:

$$\underset{\theta,\phi,\Theta}{\mathrm{argmax}} \frac{1}{N} \sum_{i=1}^{N} y_i \log\left(P_{\Phi+\Delta\Phi(\theta)+\Delta\Phi(\phi),\Theta}\left(x_i\right)\right)$$

Similar to the Domain-adaptive Pretraining stage, the dimensions of the additional parameters $\theta$ and $\phi$ are very small compared to the original parameters. By updating only the additional parameters and the classifier head, the proposed paradigm reduces the computational requirements, making it more efficient and feasible, especially for clinical settings that are often resource-constrained.

### 3.2 Two-step LLaMA-LoRA

In this study, we propose a two-step PEFT framework (as shown on the right-hand side of Figure 2). Firstly, we introduce Clinical LLaMA-LoRA, a LoRA adapter built upon LLaMA (Touvron et al., 2023) that is adapted to the clinical domain. Secondly, we introduce Downstream LLaMA-LoRA, which is trained on top of the pretrained Clinical LLaMA-LoRA and is specifically adapted to the downstream tasks.

**LLaMA models** In this study, we evaluate two LLaMA models; the 7 billion parameters version of LLaMA (Touvron et al., 2023) and the 7 billion parameters version of PMC-LLaMA(Wu et al., 2023). LLaMA was pretrained with an array of texts from multiple sources, such as English CommonCrawl, Wikipedia, ArXiv, and C4 (Raffel et al.,

| Dataset | # Class | Multilabel | # Train | # Valid | # Test |
|---------|---------|------------|---------|---------|--------|
| LOS | 4 | ✗ | 30,421 | 4,391 | 8,797 |
| MOR | 2 | ✗ | 33,954 | 4,908 | 9,822 |
| PMV | 2 | ✗ | 5,666 | 707 | 706 |
| DIAG | 1,266 | ✓ | 33,994 | 4,918 | 9,829 |
| PROC | 711 | ✓ | 30,030 | 4,357 | 8,681 |

Table 1: Statistics and types of downstream clinical document classification tasks: length of stay (LOS), mortality (MOR), prolonged mechanical ventilation (PMV), diagnoses (DIAG), and procedures (PROC).

2020). While, PMC-LLaMA is a domain-adapted LLaMA model that was pretrained on 4.8 million biomedical academic papers from PubMed Central.

**Domain-adaptive Pretraining: Clinical LLaMA-LoRA** Clinical LLaMA-LoRA is trained using a combination of MIMIC-IV de-identified discharge summaries (331,794) and radiology reports (2,321,355), resulting in a collection of 2,653,149 individual clinical notes. We evaluate five PEFT techniques, which include *LoRA* (Hu et al., 2022), *Adaptation Prompt* (Zhang et al., 2023), *Prefix Tuning* (Li and Liang, 2021), *Prompt Tuning* (Lester et al., 2021), and *P-tuning* (Liu et al., 2021b).

Our approach follows the autoregressive language modelling pretraining objective employed in the original LLaMA training. To ensure compatibility with available computational resources, we use fixed model hyperparameters that allow us to fit the LLM into a single NVIDIA A100-80GB GPU (see Appendix A.1). We optimise the hyperparameters specific to each PEFT method using Gaussian Process regression for Bayesian Optimisation (Frazier, 2018) [2] with a maximum of 20 trials. The detailed hyperparameters search space can be found in Appendix A.2. During this stage, we evaluate the perplexity scores of the LLM variants.

**Downstream Fine-tuning: Downstream LLaMA-LoRA** We fine-tune the Clinical LLaMA-LoRA and Downstream LLaMA-LoRA to clinical document classification tasks:

- **Prolonged mechanical ventilation (PMV)**: a binary classification task to predict whether a patient will require mechanical ventilation for more than seven days (Huang et al., 2020; Naik et al., 2022).
- **In-hospital mortality (MOR)**: a binary classification task to predict whether a patient will sur-

vive during their hospital stay (van Aken et al., 2021; Naik et al., 2022).

- **Length of stay (LOS)**: a multiclass classification task to predict the length of a patient's hospital stay, categorised into four time-bins: less than three days, three to seven days, one to two weeks, and more than two weeks (van Aken et al., 2021; Naik et al., 2022).
- **Diagnoses (DIAG)**: a large-scale multilabel classification task to predict the differential diagnoses of a patient, represented by simplified ICD-9 diagnosis codes (van Aken et al., 2021).
- **Procedures (PROC)**: a large-scale multilabel classification task to predict the treatments administered to a patient, represented by simplified ICD-9 procedure codes (van Aken et al., 2021).

The label and split statistics of each dataset can be found in Table 1.

During this downstream fine-tuning process, we use fixed model hyperparameters to ensure compatibility with the available computational resources, a single NVIDIA A100-80GB GPU (see Appendix B.1). We optimise the hyperparameters specific to each PEFT method using Gaussian Process regression for Bayesian Optimisation with a maximum of 20 trials. The detailed hyperparameters search space of the PEFT method can be found in Appendix B.2.

For evaluating the performance of the model on these downstream tasks, we report the Area Under the Receiver Operating Characteristic Curve (AUROC) scores. Additionally, we report the macro-averaged AUROC score across all clinical tasks as commonly done in NLP benchmarking tasks (Wang et al., 2019; Peng et al., 2019; Gu et al., 2022).

### 3.3 Baseline Models

We selected baseline models that have undergone a domain-adaptive pretraining process on clinical notes (MIMIC-III). Thus, these baseline models have been designed to perform specifically on clinical data, providing comparison points for evaluating our proposed approach of two-step adaptation in downstream clinical NLP tasks. The baseline models used in the evaluation are as follows:

- **Bio+ClinicalBERT** (Alsentzer et al., 2019): Bio+ClinicalBERT is pretrained on MIMIC-III clinical notes. It is initialised from a biomedical language model called BioBERT (Lee et al., 2019), which is pretrained on biomedical research articles.

---

[2]Specifically, we use the W&B Sweep APIs: https://docs.wandb.ai/guides/sweeps

- **BlueBERT** (Peng et al., 2019): BlueBERT is pretrained on MIMIC-III clinical notes and PubMed abstracts starting from the pretrained checkpoint of BERT (Devlin et al., 2019), a general-domain language model.
- **CORe** (van Aken et al., 2021): CORe is pretrained on MIMIC-III clinical notes and biomedical articles starting from the pretrained checkpoint of BioBERT (Lee et al., 2019).
- **UmlsBERT** (Michalopoulos et al., 2021): UmlsBERT is pretrained on MIMIC-III clinical notes using the pretrained weights of Bio+ClinicalBERT with modified architecture and pretraining objective that incorporates knowledge from the Unified Medical Language System (UMLS) Metathesaurus (Schuyler et al., 1993).

## 4 Results and Analysis

### 4.1 Domain-adaptive Pretraining

The pretraining results can be found in Table 2. We employ PEFT techniques for domain-adaptive pretraining, requiring a significantly smaller number of parameters ranging from just 0.001% to 0.24% of the original model parameters. This approach substantially reduces the required computational resources and training time. We perform a full-parameter domain-adaptive pretraining of LLaMA, referred to as **Clinical LLaMA**, using four NVIDIA A100-80GB GPUs which took 49.5 hours. Instead, PEFT techniques require less than 24 hours per epoch on average with only a single GPU with a comparable perplexity score.

LoRA emerges as the best-performing PEFT method for both LLaMA and PMC-LLaMA in the clinical domain-adaptive pretraining, achieving the lowest perplexity scores of 2.244 and 2.404, respectively, which are very similar to Clinical LLaMA's perplexity score of 2.210. This pretrained LoRA is referred to as **Clinical LLaMA-LoRA** in the subsequent sections. The following experiments in downstream fine-tuning will utilise this pretrained Clinical LLaMA-LoRA.

### 4.2 Downstream Fine-tuning

From the downstream fine-tuning results shown in Table 3, we can decompose the analysis into multiple research questions:

**Can LoRA help fine-tune LLaMA from other domains (general and biomedical) to achieve higher AUROC scores in clinical tasks?** We compare the results obtained by LLaMA and

LLaMA + LoRA, as well as PMC-LLaMA and PMC-LLaMA + LoRA, as presented in Table 3. The obtained results consistently demonstrate improved AUROC scores when utilising LoRA across all tasks. The macro-averaged AUROC score of LoRA-equipped LLaMA shows a notable 13.01% increase when compared to the LLaMA-only baseline. Similarly, LoRA-equipped PMC-LLaMA exhibits a 12.19% improvement in macro-averaged AUROC compared to the original PMC-LLaMA Both LLaMA and PMC-LLaMA, when equipped with LoRA, show significant AUROC score improvements in all tasks except the PMV prediction task, which is challenging for all model variants.

Furthermore, the marginal difference in AUROC scores between PMC-LLaMA and the general-domain LLaMA may be attributed to two factors. Firstly, the original LLaMA has been exposed to biomedical concepts during its pretraining, reducing the need for domain-adaptive pretraining to the biomedical domain. Secondly, clinical outcome prediction requires an understanding of how to apply biomedical knowledge in an interconnected manner to provide prognostic. We believe that biomedical pretraining may not be sufficient in providing such practical knowledge.

**Can LoRA-equipped LLaMA and PMC-LLaMA perform comparably in comparison to clinically trained LMs?** We compare the AUROC scores obtained by the baseline models, and LoRA-equipped LLaMA and PMC-LLaMA (see Table 3). Among the baseline models, UmlsBERT performs the best with a macro-averaged AUROC score of 72.70%. Compared to UmlsBERT, both LLaMA and PMC-LLaMA underperform with macro-averaged AUROC scores of 58.61% and 60.51%, respectively. This finding highlights the importance of clinical-specific fine-tuning.

Significant improvements can be observed in LoRA-equipped LLaMA and PMC-LLaMA, with macro-averaged AUROC scores of 71.62% and 72.70%, respectively, with noticeable improvements in the diagnoses and procedures prediction tasks. LoRA-equipped LLaMA achieves AUROC scores of 78.37% and 87.49% in the diagnoses and procedures prediction tasks, respectively, compared to 72.08% and 78.32% for UmlsBERT. This represents improvements of 6.29% in diagnoses prediction and 9.17% in procedures prediction. Improvements are also observed in the results obtained by LoRA-equipped PMC-LLaMA, outperforming

| Base Model | PEFT | Trainable Params | Train Ppl | Test Ppl | GPU | Train Time (h:m:s) |
|---|---|---|---|---|---|---|
| Clinical LLaMA | - | 6.7B (100%) | 1.811 | 2.210 | 4x80GB | 49:26:38 |
| LLaMA | **LoRA** | **8.4M (0.12%)** | **1.858** | **2.244** | 1x80GB | **21:37:42** |
| | Adaptation Prompt | 1.2M (0.02%) | 2.561 | 2.865 | 1x80GB | 24:57:17 |
| | Prefix Tuning | 5.2M (0.08%) | 2.815 | 2.748 | 1x80GB | 20:11:07 |
| | Prompt Tuning | 61.4K (0.0009%) | 4.846 | 4.007 | 1x80GB | 23:27:28 |
| | P-tuning | 16.1M (0.24%) | 2.723 | 3.271 | 1x80GB | 23:49:31 |
| PMC-LLaMA | **LoRA** | **2.1M (0.03%)** | **1.938** | **2.404** | 1x80GB | **21:32:59** |
| | Adaptation Prompt | 1.2M (0.018%) | 2.374 | 2.867 | 1x80GB | 23:33:10 |
| | Prefix Tuning | 2.6M (0.04%) | 1.789 | 2.848 | 1x80GB | 20:13:10 |
| | Prompt Tuning | 41K (0.0006%) | 4.821 | 4.385 | 1x80GB | 22:25:32 |
| | P-tuning | 2.2M (0.03%) | 3.491 | 4.572 | 1x80GB | 22:28:15 |

Table 2: Domain-adaptive Pretraining results of LLaMA and PMC-LLaMA trained on MIMIC-IV clinical notes with a language modelling objective. Lower perplexity scores indicate better language modelling performance. The **boldface row** indicates the model with the lowest perplexity score from each base model variant.

UmlsBERT by 6.73% in diagnoses prediction and 8.36% in procedures prediction.

**Can LLaMA and PMC-LLaMA with Clinical LLaMA-LoRA achieve higher AUROC scores than the clinically trained LMs?** The domain-adaptive pretraining step yields the clinically-trained LoRA adapters for LLaMA and PMC-LLaMA, denoted as **Clinical LLaMA-LoRA**. We compare the results of Clinical LLaMA-LoRA-equipped LLaMA and PMC-LLaMA with the baseline models. We evaluate Clinical LLaMA-LoRA with and without fine-tuning, referred to as "Trainable" and "Frozen" respectively.

The results indicate that Clinical LLaMA-LoRA-equipped LLaMA and PMC-LLaMA outperform the baseline models. LLaMA with a trainable Clinical LLaMA-LoRA achieves an AUROC score of 75.13%, surpassing UmlsBERT's score of 72.32%. PMC-LLaMA with a trainable Clinical LLaMA-LoRA achieves a lower AUROC score of 72.23%. LLaMA with a trainable Clinical LLaMA-LoRA also outperforms Clinical LLaMA which achieves an AUROC score of 58.86%.

These findings indicate that the Clinical LLaMA-LoRA contributes to higher AUROC scores for LLaMA and PMC-LLaMA over clinically trained LLMs, while biomedical domain-adaptive pretraining may not be necessary to improve the model's performance in the clinical settings.

**Can LLaMA and PMC-LLaMA with Clinical LLaMA-LoRA achieve higher AUROC scores than the other fine-tuning variants?** We examine the importance of the domain-adapted LoRA by comparing the results obtained by LLaMA and PMC-LLaMA equipped with Clinical LLaMA-

LoRA against the results of LLaMA and PMC-LLaMA fine-tuning, both original and with LoRA.

Firstly, we evaluate the frozen pretrained Clinical LLaMA-LoRA. Both LLaMA and PMC-LLaMA with frozen Clinical LLaMA-LoRA do not exhibit a significant increase in performance compared to the original fine-tuning. This indicates that, despite the domain-adaptive pretraining, the limited number of trainable parameters during the downstream fine-tuning restricts the potential improvement that the model can achieve. A similar finding can also be observed in the Clinical LLaMA fine-tuning whose overall performance does not differ from the original fine-tuning. This finding is further supported by the improvement in the AUROC scores of LLaMA and PMC-LLaMA with trainable Clinical LLaMA-LoRA, which achieve 75.13% and 72.23% macro-averaged AUROC scores, respectively. These represent substantial improvements from the vanilla fine-tuning performance, 58.61% and 60.51% AUROC scores.

**Can a downstream LoRA adapter improve the AUROC scores of LLaMA and PMC-LLaMA equipped with Clinical LLaMA-LoRA?** By considering Clinical LLaMA-LoRA as the "delta-updating" outcome of the domain-adaptive pretraining, we can view the downstream fine-tuning process as an additional "delta-updating" step. To investigate the impact of this approach, we conduct experiments by adding a Downstream LLaMA-LoRA to LLaMA and PMC-LLaMA models that were already equipped with Clinical LLaMA-LoRA. From Table 3, we can observe that Downstream LLaMA-LoRA fails to improve the performance of LLaMA and PMC-LLaMA with frozen Clinical LLaMA-LoRA. On the other

| Model | PMV | MOR | LOS | DIAG | PROC | Macro Average |
|---|---|---|---|---|---|---|
| BlueBERT | 57.31 | 81.34 | 72.92 | 73.39 | 76.62 | 72.32 |
| *UmlsBERT* | *58.29* | *81.83* | *73.02* | *72.08* | *78.32* | *72.70* |
| Bio+ClinicalBERT | 54.00 | 72.67 | 72.21 | 76.65 | 83.21 | 71.75 |
| CORe | 52.11 | 71.52 | 64.17 | 72.40 | 84.51 | 69.40 |
| Clinical LLaMA∗ | 52.28 | 63.22 | 56.06 | 59.31 | 63.42 | 58.86 |
| LLaMA∗ | 51.38 | 66.80 | 57.65 | 60.06 | 63.83 | 58.61 |
| + LoRA | 51.65 | 74.89 | 65.70 | 78.37 | 87.49 | 71.62 |
| + Clinical LLaMA-LoRA (Frozen) | 52.22 | 60.88 | 55.05 | 57.64 | 62.48 | 57.65 |
| + Downstream LLaMA-LoRA | 52.31 | 61.72 | 55.16 | 57.70 | 62.58 | 57.90 |
| + Clinical LLaMA-LoRA (Trainable) | 51.41 | 81.16 | 72.44 | **81.97** | **88.69** | 75.13 |
| + *Downstream LLaMA-LoRA* | *53.81* | ***83.02*** | ***73.26*** | *81.93* | *88.31* | ***76.07*** |
| PMC-LLaMA∗ | 53.06 | 66.77 | 57.94 | 60.17 | 64.63 | 60.51 |
| + *LoRA* | *53.84* | *78.03* | *66.14* | *78.81* | *86.68* | *72.70* |
| + Clinical LLaMA-LoRA (Frozen) | 51.33 | 67.19 | 58.13 | 63.59 | 68.26 | 60.06 |
| + Downstream LLaMA-LoRA | 50.90 | 67.00 | 58.31 | 60.50 | 64.42 | 60.23 |
| + Clinical LLaMA-LoRA (Trainable) | 52.88 | 75.86 | 65.89 | 79.66 | 86.85 | 72.23 |
| + Downstream LLaMA-LoRA | 52.21 | 76.54 | 68.42 | 78.67 | 87.08 | 72.58 |

Table 3: AUROC scores in clinical downstream document classification tasks. The macro-averaged AUROC score is calculated by taking the average of AUROC scores across all tasks. The **boldface cell** indicates the highest AUROC score in a column, the *row in italic* indicates the variant with the highest macro-averaged AUROC in its category. + *LoRA* denotes applying LoRA on top of the pretrained LLM without domain-adaptive pretraining. + *Clinical LLaMA-LoRA* denotes applying Clinical LLaMA-LoRA that is domain-adaptively pretrained on top of the pretrained LLM. + *Downstream LLaMA-LoRA* denotes applying Downstream LLaMA-LoRA on top of the LLM + Clinical LLaMA-LoRA. *Frozen* means that the parameters are not trainable, while *Trainable* means that the parameters are trainable. ∗ Due to restricted computing resources, the fine-tunings of Clinical LLaMA, LLaMA, and PMC-LLaMA were constrained to only training the final classification layer.

hand, improvement can be observed when adding Downstream LLaMA-LoRA to LLaMA with trainable Clinical LLaMA-LoRA. This combination of LLaMA with trainable Clinical LLaMA-LoRA and Downstream LLaMA-LoRA achieves the highest macro-averaged AUROC score of 76.07%. The macro-averaged AUROC score of Clinical LLaMA-LoRA was almost similar to that of PMC-LLaMA with LoRA, suggesting similar efficacy between Clinical LLaMA-LoRA and the full fine-tuning process that PMC-LLaMA has undergone. Moreover, Clinical LLaMA-LoRA offers the advantage of reduced computational resources and training time, which is aligned with the requirements of practical implementation in clinical settings.

Overall, our proposed method manages to achieve better performance in comparison to clinically trained models. We also provide a comparison with the state-of-the-art method of PMV, mortality, and length of stay predictions, called BEEP (Naik et al., 2022), which leverages retrieval augmentation method to provide more contextual information to the model during inference. The comparison is only partial as BEEP models were not evaluated on the diagnosis and procedure prediction tasks. As shown in Appendix C, our best-

performing model achieves a 70.03% averaged AUROC score, which is slightly worse compared to the best-performing BEEP model with 72.26% averaged AUROC score. However, it is worth noting that our proposed method and the state-of-the-art method are complementary to each other. Hence, future work may explore the possibility of combining the two approaches.

## 5   Conclusions

In this study, we propose a two-step PEFT framework. We introduce Clinical LLaMA-LoRA, a LoRA (Hu et al., 2022) adapter built upon LLaMA (Touvron et al., 2023). Then, we introduce Downstream LLaMA-LoRA, a task-specific adapter that is trained on top of the pretrained Clinical LLaMA-LoRA. The fusion of the two adapters achieves an AUROC score of 76.07% macro-averaged across all clinical NLP downstream tasks, which represents a 3.37% improvement over the best-performing clinical LLM. Our proposed framework achieves improvement in performance while reducing the computational requirements, which is suited for clinical settings that are often constrained by their computational power.

## Limitations

This study presents a two-step PEFT framework aimed at effectively adapting LLMs to diverse clinical downstream applications. However, the evaluation of our model was restricted to MIMIC-based datasets, which are constrained to English and obtained exclusively within the Commonwealth of Massachusetts, United States of America. Consequently, despite the promising efficacy demonstrated by our proposed method, it would have been advantageous to directly assess its performance across diverse hospital systems spanning other geographical locations and languages. This would enable a more comprehensive understanding of its applicability and generalizability. However, it is essential to acknowledge that conducting such an analysis would require working within a trusted research environment and obtaining the necessary permissions to access the relevant datasets.

It is crucial to recognise the restrictions imposed on accessing internal clinical datasets, as they limit our ability to evaluate the effectiveness of our approach across different care provider systems. Therefore, we encourage care providers to conduct internal experiments within their trusted research environment to ensure the efficacy of our proposed method within their specific use cases should they adopt this approach.

Despite the demonstrated performance improvements, the proposed model may still be susceptible to spurious correlations. Predicting patient outcomes solely based on clinical notes presents significant challenges due to the other factors that may not be captured within those notes. For instance, the length of a patient's in-hospital stay is not solely correlated with their diagnoses and disease progression. Factors such as the patient's insurance status, which is not typically mentioned in clinical notes, can severely impact the duration of a patient's stay. Therefore, we encourage end users of such clinical LLMs to consider additional measures to ensure predictions that reflect a holistic view of the patient's situation, instead of relying solely on the predictions of LLMs.

## Ethics Statement

In this study, we use MIMIC-based datasets obtained after completing the necessary training. These datasets comply with de-identification standards set by the Health Insurance Portability and Accountability Act (HIPAA) through data cleansing. Due to privacy concerns, we refrain from including direct excerpts of the data in the paper. We also refrain from publicly sharing the pretrained checkpoints.

While our model demonstrates effectiveness, it is important to acknowledge the risks associated with relying solely on clinical outcome prediction models. There are crucial pieces of information that can be found beyond the scope of clinical notes. Considering the potential impact on patient health outcomes, it is crucial to exercise caution when utilising these clinical LLMs. Therefore, we propose that the PEFT adapter generated by our framework, in conjunction with the pretrained LLM, should be used as an aid rather than a replacement for trained clinical professionals.

## References

Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Peter I. Frazier. 2018. A tutorial on bayesian optimization. *CoRR*, abs/1807.02811.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2022. Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *Proceedings of the 36th International Conference on Machine Learning*, page 2790–2799. PMLR.

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*.

Kexin Huang, Abhishek Singh, Sitong Chen, Edward Moseley, Chih-Ying Deng, Naomi George, and Charolotta Lindvall. 2020. Clinical XLNet: Modeling Sequential Clinical Notes and Predicting Prolonged Mechanical Ventilation. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 94–100, Online. Association for Computational Linguistics.

Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J. Pollard, Benjamin Moody, Brian Gow, Li-wei H. Lehman, Leo A. Celi, and Roger G. Mark. 2023. MIMIC-IV, a freely accessible electronic health record dataset. *Scientific Data*, 10(1):1.

Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific Data*, 3(1):160035.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Eric Lehman, Evan Hernandez, Diwakar Mahajan, Jonas Wulff, Micah J. Smith, Zachary Ziegler, Daniel Nadler, Peter Szolovits, Alistair Johnson, and Emily Alsentzer. 2023. Do We Still Need Clinical Language Models?

Eric Lehman and Alistair Johnson. 2023. Clinical-T5: Large Language Models Built Using MIMIC Clinical Text.

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.

Bo Li, Yezhen Wang, Tong Che, Shanghang Zhang, Sicheng Zhao, Pengfei Xu, Wei Zhou, Yoshua Bengio, and Kurt Keutzer. 2020. Rethinking distributional matching based domain adaptation. *CoRR*, abs/2006.13352.

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.

Yikuan Li, Ramsey M. Wehbe, Faraz S. Ahmad, Hanyin Wang, and Yuan Luo. 2022. Clinical-longformer and clinical-bigbird: Transformers for long clinical sequences. *CoRR*, abs/2201.11838.

Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy Chowdhury, Yun Li, Hejie Cui, Xuchao Zhang, Tianjiao Zhao, Amit Panalkar, Wei Cheng, Haoyu Wang, Yanchi Liu, Zhengzhang Chen, Haifeng Chen, Chris White, Quanquan Gu, Carl Yang, and Liang Zhao. 2023. Beyond one-model-fits-all: A survey of domain specialization for large language models. *CoRR*, abs/2305.18703.

Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021a. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *CoRR*, abs/2110.07602.

Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021b. GPT understands, too. *CoRR*, abs/2103.10385.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. BioGPT: generative pre-trained transformer for

biomedical text generation and mining. *Briefings in Bioinformatics*, 23(6). Bbac409.

Bhavitvya Malik, Abhinav Ramesh Kashyap, Min-Yen Kan, and Soujanya Poria. 2023. UDAPTER - efficient domain adaptation using adapters. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2249–2263, Dubrovnik, Croatia. Association for Computational Linguistics.

Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods. https://github.com/huggingface/peft.

George Michalopoulos, Yuanxin Wang, Hussam Kaka, Helen Chen, and Alexander Wong. 2021. Umls-BERT: Clinical domain knowledge augmentation of contextual embeddings using the Unified Medical Language System Metathesaurus. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1744–1753, Online. Association for Computational Linguistics.

Aakanksha Naik, Sravanthi Parasa, Sergey Feldman, Lucy Lu Wang, and Tom Hope. 2022. Literature-augmented clinical outcome prediction. In *Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 438–453. Association for Computational Linguistics.

Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 58–65, Florence, Italy. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67. Citation Key: JMLR:v21:20-074.

Sebastian Ruder, Jonas Pfeiffer, and Ivan Vulić. 2022. Modular and Parameter-Efficient Fine-Tuning for NLP Models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 23–29, Abu Dubai, UAE. Association for Computational Linguistics.

P L Schuyler, W T Hole, M S Tuttle, and D D Sherertz. 1993. The UMLS Metathesaurus: Representing different views of biomedical concepts. *Bulletin of the Medical Library Association*, 81(2):217–222.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Schärli, Aakanksha Chowdhery, Philip Andrew Mansfield, Blaise Agüera y Arcas,

Dale R. Webster, Gregory S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle K. Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2022. Large language models encode clinical knowledge. *CoRR*, abs/2212.13138.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Andrew Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agüera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle K. Barral, Dale R. Webster, Gregory S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. 2023. Towards expert-level medical question answering with large language models. *CoRR*, abs/2305.09617.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.

Betty van Aken, Jens-Michalis Papaioannou, Manuel Mayrdorfer, Klemens Budde, Felix Gers, and Alexander Loeser. 2021. Clinical Outcome Prediction from Admission Notes using Self-Supervised Knowledge Integration. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 881–893, Online. Association for Computational Linguistics.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Further finetuning llama on medical papers. *CoRR*, abs/2304.14454.

Xi Yang, Aokun Chen, Nima M. Pournejatian, Hoo Chang Shin, Kaleb E. Smith, Christopher

Parisien, Colin Compas, Cheryl Martin, Anthony B. Costa, Mona G. Flores, Ying Zhang, Tanja Magoc, Christopher A. Harle, Gloria P. Lipori, Duane A. Mitchell, William R. Hogan, Elizabeth A. Shenkman, Jiang Bian, and Yonghui Wu. 2022. A large language model for electronic health records. *npj Digit. Medicine*, 5.

Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. 2023. LLaMA-Adapter: Efficient Fine-tuning of Language Models with Zero-init Attention.

Rongsheng Zhang, Yinhe Zheng, Xiaoxi Mao, and Minlie Huang. 2021. Unsupervised domain adaptation with adapter. *CoRR*, abs/2111.00667.

## A  Hyperparameters for the Domain-adaptive Pretraining

### A.1  Fixed Model Hyperparameters

| Hyperparameter | Value |
|---|---|
| Learning rate | 3e-4 |
| Warmup steps ratio | 0.06 |
| Maximum sequence length | 512 |
| Gradient accumulation step | 4 |
| Batch size | 10 |

Table 4: Fixed model hyperparameters for language modelling pretraining. These hyperparameters remain unchanged to fit LLaMA into a single GPU.

### A.2  PEFT Hyperparameters Optimisation Search Space

| PEFT | Hyperparameter | Search space |
|---|---|---|
| LoRA | r | [2, 4, 8, 16] |
| | alpha | [4, 8, 16, 32] |
| | dropout | [0.0, 0.1, 0.2] |
| Prefix Tuning | num virtual tokens | [1, 5, 10, 15, 20] |
| | prefix projection | [true, false] |
| Prompt Tuning | num virtual tokens | [1, 5, 10, 15, 20] |
| | prompt init | [text, random] |
| P-Tuning | num virtual tokens | [1, 5, 10, 15, 20] |
| | reparameterisation | ["MLP", "LSTM"] |
| | hidden size | [64, 128, 256, 768] |
| | num layers | [1, 2, 4, 8, 12] |
| | dropout | [0.0, 0.1, 0.2] |
| Adaptation Prompt | adapter length | [5, 10] |
| | adapter layers | [10, 20, 30] |

Table 5: The search space for PEFT Hyperparameters optimisation runs during the domain adaptation fine-tuning with language modelling objective. Each PEFT technique has a specific set of hyperparameters to tune, we selected the combination of hyperparameters which has the lowest perplexity score.

Specifically for Prompt Tuning, we use a common prompt initialisation text "Finish this clinical note:".

## B  Hyperparameters for the Downstream Fine-tuning

### B.1  Fixed Model Hyperparameters

| Hyperparameter | Value |
|---|---|
| Learning rate | 5e-5 |
| Warmup steps ratio | 0.06 |
| Maximum sequence length | 512 |
| Gradient accumulation step | 10 |
| Batch size | 10 |

Table 6: Fixed model hyperparameters for the clinical downstream fine-tuning. These hyperparameters remain unchanged to fit LLaMA into a single GPU.

### B.2  PEFT Hyperparameters Optimisation Search Space

| PEFT | Hyperparameter | Search space |
|---|---|---|
| LoRA | r | [2, 4, 8, 16] |
| | alpha | [4, 8, 16, 32] |
| | dropout | [0.0, 0.1, 0.2] |

Table 7: The search space for PEFT Hyperparameters optimisation runs during the downstream fine-tuning. Each PEFT technique has a specific set of hyperparameters to tune, we selected the combination of hyperparameters which has the highest AUROC score.

## C  Comparison with BEEP (Naik et al., 2022)

| Model | PMV | MOR | LOS | Avg |
|---|---|---|---|---|
| *BEEP* | *59.43* | *84.65* | *72.71* | *72.26* |
| Our method | 53.81 | 83.02 | 73.26 | 70.03 |

Table 8: AUROC scores in a subset of the clinical downstream document classification tasks. The macro-averaged AUROC score is calculated by taking the average of AUROC scores across this subset of tasks. The *row in italic* indicates the model variant with the highest macro-averaged AUROC.

We compared our method with the state-of-the-art clinical outcome prediction model, BEEP (Naik et al., 2022), which leverages a retrieval augmentation technique to enhance the predictive capabilities of clinical language models. A small caveat is that BEEP focused on three downstream tasks: prolonged mechanical ventilation, mortality, and length of stay predictions. We selected the best-performing solution from BEEP, UmlsBERT with weighted voting retrieval augmentation, based on the averaged AUROC score to compare with our solution. While BEEP outperforms our approach, particularly in the prediction of PMV, it is crucial to emphasise that our method achieves its predictions without relying on retrieval augmentation. Future work may explore using retrieval augmentation on top of our proposed method.

## D  Training Configurations

We use HuggingFace's Transformers (Wolf et al., 2020) and PEFT (Mangrulkar et al., 2022) libraries for the experiments. All LLaMA-based models are trained on one NVIDIA A100-80GB GPU, while the baseline models are trained on a single NVIDIA GeForce GTX 1080 Ti-16GB GPU.

# E Artefacts

The pretrained baseline models including BioClinicalBERT (Alsentzer et al., 2019), BlueBERT (Peng et al., 2019), and CORe (van Aken et al., 2021) were released under the Creative Commons designation CC0 1.0 Universal license, whereas UmlsBERT (Michalopoulos et al., 2021) was released under the MIT license. LLaMA (Touvron et al., 2023) was released under a noncommercial license.

MIMIC-III and MIMIC-IV dataset was released under the PhysioNet Credentialed Health Data License 1.5.0 and can only be accessed after one finishes the CITI Data or Specimens Only Research training[3].

---

[3]https://physionet.org/about/citi-course/

# A Multilevel Analysis of PubMed-only BERT-based Biomedical Models

**Vicente Ivan Sanchez Carmona** and **Shanshan Jiang** and **Bin Dong**
Ricoh Software Research Center (Beijing) Co., Ltd
{Vicente.Carmona, Shanshan.Jiang, Bin.Dong}@cn.ricoh.com

## Abstract

Biomedical NLP models play a big role in the automatic extraction of information from biomedical documents, such as COVID research papers. Three landmark models have led the way in this area: BioBERT, MSR Biomed-BERT, and BioLinkBERT. However, their shallow evaluation –a single mean score– forbid us to better understand how the contributions proposed in each model advance the Biomedical NLP field. We show through a Multilevel Analysis how we can assess these contributions. Our analyses across 5000 fine-tuned models show that, actually, BiomedBERT's true effect is bigger than BioLinkBERT's effect, and the success of BioLinkBERT does not seem to be due to its contribution –the Link function– but due to an unknown factor.

## 1 Introduction

Machine reading of biomedical texts has greatly advanced due to pretrained NLP models such as BERT. Biomedical NLP applications are of great value due to their utility in real-world scenarios such as answering questions which require background knowledge or the extraction of complex biomedical entities from astonishing volumes of academic papers related to COVID, for example.

Of special acknowledgement, three BERT-based biomedical models, trained on PubMed abstracts (their only source of biomedical knowledge), led the way to concise research contributions on biomedical NLP, namely, BioBERT (Lee et al., 2019) proposing Domain Adaptive Pretraining (DAPT), MSR BiomedBERT[1] (Gu et al., 2021, which we refer to as BiomedBERT) which challenged DAPT by pretraining BERT from scratch with PubMed abstracts, and BioLinkBERT (Yasunaga et al., 2022) which implemented a way to link hyperlinked documents at pretraining time – the Link function. These 3 contributions resulted

in significant improvements on downstream scores on the BLURB benchmark (Gu et al., 2021).

However, we claim, current evaluation methods are oversimplistic. They reduce to a simple mean score across datasets in the BLURB suite –a single estimate. This forbid us to better understand the contributions proposed by each work such as their effect on scores and interaction with downstream datasets. Moreover, from this single estimate, how can we disentangle the contributions' effects from the effects of other variables such as random seeds, learning rates, or number of epochs? We cannot. And while some works show ablation studies to see the particular effects of the proposed contribution, doing so to isolate it from all possible variables (including those mentioned above) leads to an exponential number of ablations which results in a non-environmentally friendly, unfeasible approach if pretraining is necessary for each ablation.

In this paper, we propose a regression analysis widely used in the fields of Psychology and the Social Sciences –Multilevel Analysis– to account for the effects of all measurable variables, without the need for ablations or further pretraining experiments, in order to disentangle their effects from the true effect of the proposed contributions from BioBERT, BiomedBERT, and BioLinkBERT. Our analyses show that, actually, random seeds have a big effect on downstream scores. Also, while BioBERT's and BiomedBERT's contributions have a big and significant effect by improving on vanilla BERT's score by 2.25 and 4.36 points, respectively, on average across BLURB datasets, BioLinkBERT's Link function shows only a big effect for QA datasets but not for any other dataset.

## 2 Background and Related Work

### 2.1 Multilevel Regression Analysis

Multilevel models (MLMs) are a type of regression analysis where the outcome to be modeled

---

[1]Previously known as PubMedBERT.

(downstream scores in our case) is dependent on a set of independent variables that can pertain to different levels in a hierarchy. In our case, we define our problem as a 2-level hierarchy where the lowest level –Level 1– contains fine-tuned models, which is nested inside the upper level –Level 2– which corresponds to groups of fine-tuned models grouped according to the choice of pre-trained model and downstream dataset; for example, BioBERT-BIOSSES is a group of BioBERT models fine-tuned on the BIOSSES dataset.[2]

Thus, variables at level 1 describe fine-tuning attributes such as learning rates, batch size, and number of epochs. On the other hand, level-2 variables describe attributes of the pretrained models, such as the contribution proposed by a work (for example, the Link function proposed by BioLinkBERT), and the choice of downstream dataset. In this way, a 2-level MLM (de Leeuw and Meijer, 2008) can be expressed as:

$$y = \beta_0 + \sum_{fixed} \beta_i x_i + \sum_{random} \gamma_{ij} x_{ij} + u_{0j} + e \quad (1)$$

where $\beta_0$ is the grand-mean intercept; the first summation corresponds to level-1 and level-2 *fixed-effects* coefficients ($\beta_i$) which represent the average individual effect of each variable ($x_i$) on downstream scores ($y$); the second summation is a key term that distinguishes MLMs from other regression models: level-1 *random-effects*, i.e. an *adjusted* effect ($\gamma_{ij}$) on the level-1 fixed-effects coefficients according to each group (indexed by $j$);[3] and similarly for the random intercepts $u_{0j}$ which are adjusted effects for each group on the grand-mean intercept; finally, $e$ is the residual. This model can be fitted using Maximum Likelihood Estimation or variants.

## 2.2 MLMs for Experimental Analyses

MLMs[4] have been widely used for analysing experimental and observational data by fields such as Psychology (Muradoglu et al., 2023; Judd et al.,

2017), Linguistics (Baayen et al., 2008), and the Social Sciences (Rasbash et al., 2010; de Leeuw and Meijer, 2008). For example, in the field of Education, MLMs analyze the impact of both student (level-1) variables (age, socioeconomic status, gender) and school (level-2) variables (mean socioeconomic status, ethnicity proportions) on students' academic performance (Goldstein et al., 2007).

Works in Psychology have used MLMs to disentangle the effects of different variables at different levels while measuring their impact on participants' reaction time on cognitive tasks (Kliegl et al., 2011, 2010). Moreover, work in Linguistics has leveraged MLMs to model the effect of between-speaker features (age, country, etc.) and within-speaker features (length of sentence, sequential position of phrase, etc.) on articulation rate of spoken sentences (Quené, 2008).

To our knowledge, our work is the first approach towards leveraging MLMs for analysis of biomedical NLP models.

## 3 Dataset and Multilevel Model

### 3.1 Dataset for Multilevel Analysis

To generate a dataset to fit an MLM that explains the impact of variables on downstream scores, we fine-tune[5] BioBERT, BiomedBERT, BioLinkBERT and vanilla BERT (which we use as baseline) on all datasets in the BLURB suite. We use test set scores as the dependent variable. And we use fine-tuning and pretraining features as level-1 and level-2 variables, respectively.

To obtain robust estimates of effects (regression coefficients) we not only include test scores from the best-validation-score models,[6] we also include the scores from a vicinity around the best-validation-score models. This vicinity is defined around the values of variables that lead to the best validation score, (namely learning rate, batch size, and number of epochs), in a way that scores in the vicinity are consistent with the highest validation score but allowing for variation in order to estimate standard errors. We follow this process for 3 different random seeds for each dataset. We obtained 5154 fine-tuned models across datasets and pretrained models.

Table 1 shows a summary of all the variables

---

[2]Therefore, at level 2 we have 52 groups: 4 choices of pretrained models (including BERT) by 13 downstream datasets.

[3]For instance, we may expect random seeds to have a different effect, due to chance, on test scores depending on the choice of group, i.e. depending on the choice of pretrained model and dataset; thus, for each group, we can estimate the number of points, represented by a $\gamma_{ij}$ coefficient, that a random seed deviates from the average effect of that random seed across all groups, represented by a $\beta_i$ coefficient.

[4]Also known as Mixed Models and Hierarchical Linear Models in other fields.

[5]We follow fine-tuning guidelines from BiomedBERT and BioLinkBERT, and we use BioLinkBERT's fine-tuning code.

[6]Models which scored the highest on the validation set of each dataset.

used for the analysis.[7] Most of the variables are indicator (binary) variables which take the value of 1 whenever that variable is used by a particular instance and zero otherwise. On the other hand, the variable num_epochs takes integer values representing the number of epochs used for fine-tuning a specific model.

## 3.2 Multilevel Model

We instantiate Equation 1 with the variables in Table 1. As a common goal in the literature (Frank E. Harrell, 2015), we aim to find which variables have a statistically-significant effect on downstream scores across BLURB datasets.[8] We follow model-building, hypothesis-testing, and evaluation strategies from Robson and Pevalin (2016), Sommet and Morselli (2021), and Brown (2021). To fit MLMs we use the R-package *lmerTest* (Kuznetsova et al., 2017). We use the statistical tests from *lmerTest* to compute significance values ($\alpha = 0.05$ level), AIC, and BIC scores.[9] Furthermore, to estimate the proportion of explained variability in test scores by our variables we compute R-squared effects using the framework of Rights and Sterba (2019) via the R-package *r2mlm* (Shaw et al., 2022).

We added an additional term to our MLM not shown in Equation 1: interaction terms between level-2 variables; these terms are of the form $\beta_m(x_i \times x_k)$, which will help us see if a variable behaves differently for particular datasets in Section 4.

## 4 Multilevel Analysis and Results

We show the results of fitting our MLM. For Tables 2, 3, and 4, the statistical significance code is: p=0 '***', p<0.001 '**', p<0.01 '*'.

**MLM results for level-1 variables:** We first test for the statistical significance of fixed- and random-effects of level-1 variables. We observe in Table 2 that the fixed-effect of only one variable is significant, namely lr_1; this means that the learning rate of 1e-5 has a significant effect across models and datasets: models fine-tuned with this learn-

ing rate, on average, will lose 1 downstream point as shown by the coefficient of lr_1. We also see that the random seeds seed_20 and seed_47 have a small, positive impact on test scores, on average, across models and datasets; nevertheless, these fixed-effects seem to be due to chance since they are not statistically significant. However, likelihood ratio tests show that all random coefficients are statistically significant (Table 4). This means that level-1 variables behave in different ways for each group (combination of pretrained model and dataset) as we explain below.

**Does chance play a role?** All level-1 variables behave differently for each pretrained model; but, we note in particular that seed_20 and seed_47 contribute the biggest variability in test scores as seen in Table 4: on average, scores vary up to ($\pm$) 4.79 and ($\pm$) 6.21 points due to the choice of random seed.[10] If we average all the random coefficients[11] of seed_20 and seed_47 for each pretrained model across datasets, we find that BioBERT loses 2.33 and 3 points when using such random seeds. However, BiomedBERT and BioLinkBERT gain 0.52, 0.22 and 0.09, 0.64 points, respectively, due to such randomness.

**MLM results for level-2 variables:** We observe that most level-2 variables are statistically significant (Table 2), such as the effects of all datasets, meaning that different datasets lead to different results. Also significant are the contributions from BioBERT and BiomedBERT, namely, DAPT and Pretrain_PubMed, respectively, meaning that their effects –an average gain of 2.25 and 4.36 points with respect to vanilla BERT– are consistent across datasets. Surprisingly, the Link function is not significant: probably, its effect is not systematic across datasets. To better understand its effect, we estimated its interaction with all datasets; as we see in Table 3, when the Link function is used with QA datasets, its effect is remarkable: models fine-tuned with BioASQ and PubMedQA datasets gain, on avg., 8.62 and 3.68 points, respectively. However, this figure does not happen with any other dataset. Moreover, the effect of the Link function, besides non-significant, is rather small, which means that whenever the Link function is used, on average, we

---

[7]We include variables for the downstream datasets to take into account the fact that some datasets may be more difficult than others which may impact on the scores.

[8]We chose an MLM over simple linear regression since 1) it allows for multiple levels of analysis, and 2) the fine-tuned models inside a group are not independent from each other and only MLMs can account for such non-independence.

[9]We prefer models that decrease AIC or BIC scores.

[10]These figures represent a comparison of how much variability seed_20 and seed_47 introduce in the test scores with respect to the variability introduced by seed_59.

[11]We do not display the random coefficients since we believe it is more informative to provide an aggregated estimate.

| Variable name | Level | Description |
|---|---|---|
| seed_20, seed_47, seed_59 | 1 | Random seeds used for fine-tuning the Biomedical models |
| lr_1, lr_2, lr_3, lr_4, lr_5 | 1 | Learning rates used for fine-tuning the Biomedical models |
| batch_16, batch_32 | 1 | Batch sizes used for fine-tuning the Biomedical models |
| num_epochs | 1 | Number of epochs for fine-tuning the Biomedical models |
| BioBERT | 2 | Indicator variable for BioBERT |
| BiomedBERT | 2 | Indicator variable for BiomedBERT |
| BioLinkBERT | 2 | Indicator variable for BioLinkBERT |
| DAPT | 2 | Indicator of Domain Adaptive Pretraining on BERT |
| Pretrain_PubMed | 2 | Indicator of pretraining BERT with PubMed data from scratch |
| Link | 2 | Indicator variable of BioLinkBERT's Link function |
| all datasets names | 2 | Indicator variables of the datasets in the BLURB suite |

Table 1: Variables used to model the variability in downstream scores for target Biomedical NLP models across datasets in the BLURB suite. Level 1 corresponds to fine-tuning; level 2 to pretraining; all datasets names: BC2GM, BC5_chem, BC5_disease, NCBI, JNLPBA, PICO, ChemProt, DDI, GAD, BIOSSES, HoC, BioASQ, PubMedQA.



Figure 1: R-squared: Decomposition of variance across fixed and random effects.

would only see an improvement of 0.07 points on any downstream dataset.

**Effects from pretrained models:** If we fit our MLM with indicator variables for each pretrained model, instead of their contributions, we obtain the following effects: BioBERT (1.79**), Biomed-BERT (4.52***), BioLinkBERT (3.47***). The result for BioLinkBERT seems to contradict the non-significant effect of the Link function. It does not. The effect of the BioLinkBERT variable takes into account all possible functions inside BioLinkBERT (without disentangling them) including the variable of Pretrain_PubMed since BioLinkBERT was pretrained from scratch with PubMed data. This means that, overall, BioLinkBERT is highly useful: it surpasses vanilla BERT, on average, by 3.47 points across datasets, though the Link function does not seem to be the main reason for this result due to its small effect size and lack of statistical significance. Surprisingly, we see that BiomedBERT

has the biggest mean effect of all models: 4.52 points improvement over BERT.

**R-squared effects:** As shown in Figure 1, fixed-effects of level-1 and level-2 variables account for around 70% of all the variability in the test scores; however, given that most of the level-1 coefficients are non-significant and moderately small, we would expect them to contribute little to this explanation of variability. Surprisingly, though, level-1 random coefficients (slope variation) account for around 10% of the variance in test scores, a considerable portion of the variability. Finally, we note that around 20% of the variance remains unexplained (the residual part) which may mean two things. First, there is still room for adding variables at either level to better explain the test scores; we hypothesized that other pretraining features, such as batch size, could impact on the scores, however, it was not possible to add them to the analysis since they perfectly correlate with variables already added, leading to the problem of collinearity. And second, fully understanding NLP models is a complex task which requires detailed analyses of several variables.

**Robust estimates:** As shown in Table 2, most of the standard errors (SEs) are small which means that our coefficients estimates are robust, i.e. their estimation is precise due to the low variability represented by the corresponding SE, something that could be more difficult to achieve when only averaging scores from a handful of models across random seeds as is usual in the NLP literature.

| Variable | Coeff. ($\beta$) | SE | t |
|---|---|---|---|
| Intercept | 53.96*** | 0.88 | 61.09 |
| seed_20 | 0.49 | 0.70 | 0.70 |
| seed_47 | 0.21 | 0.89 | 0.24 |
| lr_1 | -1.00** | 0.35 | -2.82 |
| lr_2 | 0.59 | 0.51 | 1.14 |
| lr_3 | 0.65 | 0.48 | 1.35 |
| lr_4 | 0.15 | 0.24 | 0.61 |
| batch_16 | 0.28 | 0.17 | 1.59 |
| num_epochs | -0.02 | 0.02 | -1.00 |
| DAPT | 2.25*** | 0.38 | 5.90 |
| Pretrain_PubMed | 4.36*** | 0.41 | 10.41 |
| Link | 0.07 | 0.32 | 0.21 |
| BC2GM | 27.40*** | 2.02 | 13.54 |
| BC5_chem | 35.42*** | 0.60 | 58.49 |
| BC5_disease | 25.96*** | 0.64 | 40.21 |
| NCBI | 31.79*** | 0.73 | 43.29 |
| JNLPBA | 21.26*** | 0.82 | 25.92 |
| PICO | 16.14*** | 0.89 | 18.00 |
| ChemProt | 17.56*** | 0.73 | 23.97 |
| DDI | 23.07*** | 0.81 | 28.23 |
| GAD | 23.81*** | 0.57 | 41.45 |
| BIOSSES | 20.85*** | 0.67 | 30.74 |
| HoC | 25.48*** | 0.75 | 33.59 |
| BioASQ | 17.58*** | 0.87 | 20.16 |

Table 2: Results of MLM: fixed-effects of variables at levels 1 and 2. Coeff: coefficient. SE: Standard Error. t: t-value. We use seed_59, lr_5, batch_32, and PubMedQA as baselines to avoid collinearity.

| Interaction | Coeff. ($\beta$) | SE | t |
|---|---|---|---|
| DAPT×BIOSSES | -7.49*** | 0.71 | -10.5 |
| PubMed×NCBI | -2.31* | 0.86 | -2.6 |
| PubMed×PICO | -2.22* | 0.95 | -2.3 |
| PubMed×BIOSSES | 13.1*** | 0.63 | 20.5 |
| PubMed×BioASQ | 7.37*** | 1.49 | 4.9 |
| Link×HoC | -3.58*** | 0.85 | -4.2 |
| Link×BioASQ | 8.62*** | 1.83 | 4.6 |
| Link×PubMedQA | 3.68* | 1.38 | 2.65 |

Table 3: Results of MLM: interaction terms. PubMed stands for Pretrain_PubMed. Only statistically significant interactions are displayed.

| Variable | Variance | Std. Dev. |
|---|---|---|
| Intercepts | 9.29*** | 3.04 |
| seed_20 | 22.96*** | 4.79 |
| seed_47 | 38.59*** | 6.21 |
| lr_1 | 3.03*** | 1.74 |
| lr_2 | 9.72*** | 3.11 |
| lr_3 | 8.78*** | 2.96 |
| lr_4 | 0.42* | 0.64 |
| batch_16 | 0.37** | 0.60 |
| num_epochs | 0.01** | 0.10 |

Table 4: Results of MLM: random effects (random intercepts and random coefficients). Variables seed_59, lr_5, batch_32 are used as baselines to avoid collinearity.

single mean score. We hope the community will adopt MLMs as a deeper evaluation method.

## 5 Conclusions

Our multilevel analysis of Biomedical models can disentangle the effects from fine-tuning and pre-training by providing particular effects of each variable with respective statistical significance. As we saw, contrary to expectation, BiomedBERT has the biggest mean effect across datasets from all models. Moreover, even though BioLinkBERT holds as a useful model, its main contribution –the Link function– does not seem to be the main reason for its success, except for QA datasets where the Link function excels. Furthermore, we showed that all fine-tuning variables behave differently for each pretrained model, giving some advantage to some models purely by chance. And this figure, according to R-squared tests, accounts for 10% of all the test scores; thus, we suggest using several random seeds to counterbalance their effects. Finally, we note that it would be nearly impossible to see all these figures with current evaluation methods –a

## Limitations

We note that there may be more independent variables having an effect on downstream scores that we did not take into account due to their difficulty to be measured, or to be known, such as detailed pretraining hyperparameters or data pre-processing methods. Also, we note that our design of the problem as a 2-level hierarchy may not be the most optimal design; there are more design types that can be operationalized via MLMs; however, hierarchical models are the most common and studied in the literature. Furthermore, in this paper, we fine-tuned the base sizes of the pretrained models (e.g. BioLinkBERT-base), we did not analyze the large-size models (e.g. BioLinkBERT-large) which we leave as future work. Also, due to GPU memory limitations, we did not explore more levels of the variables studied such as using a batch of size 64 for fine-tuning which may benefit some of the models.

# References

R.H. Baayen, D.J. Davidson, and D.M. Bates. 2008. Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language*, 59(4):390–412. Special Issue: Emerging Data Analysis.

Violet A. Brown. 2021. An introduction to linear mixed-effects modeling in R. *Advances in Methods and Practices in Psychological Science*, 4(1):1–19.

Jan de Leeuw and Erik Meijer. 2008. *Handbook of Multilevel Analysis*, first edition. Springer New York, NY.

Jr. Frank E. Harrell. 2015. *Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis*, second edition. Springer Cham.

Harvey Goldstein, Simon Burgess, and Brendon McConnell. 2007. Modelling the effect of pupil mobility on school differences in educational achievement. *Journal of the Royal Statistical Society. Series A (Statistics in Society)*, 170(4):941–954.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans. Comput. Healthcare*, 3(1).

Charles M. Judd, Jacob Westfall, and David A. Kenny. 2017. Experiments with more than one random factor: Designs, analytic models, and statistical power. *Annual Review of Psychology*, 68(1):601–625. PMID: 27687116.

Reinhold Kliegl, Michael E. J. Masson, and Eike M. Richter. 2010. A linear mixed model analysis of masked repetition priming. *Visual Cognition*, 18(5):655–681.

Reinhold Kliegl, Ping Wei, Michael Dambacher, Ming Yan, and Xiaolin Zhou. 2011. Experimental effects and individual differences in linear mixed models: Estimating the relationship between spatial, object, and attraction effects in visual attention. *Frontiers in Psychology*, 1.

Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. lmerTest package: Tests in linear mixed effects models. *Journal of Statistical Software*, 82(13):1–26.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Melis Muradoglu, Joseph R. Cimpian, and Andrei Cimpian. 2023. Mixed-effects models for cognitive development researchers. *Journal of Cognition and Development*, 24(3):307–340.

Hugo Quené. 2008. Multilevel modeling of between-speaker and within-speaker variation in spontaneous speech tempo. *The Journal of the Acoustical Society of America*, 123(2):1104–1113.

Jon Rasbash, George Leckie, Rebecca Pillinger, and Jennifer Jenkins. 2010. Children's Educational Progress: Partitioning Family, School and Area Effects. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 173(3):657–682.

J. D. Rights and S. K. Sterba. 2019. Quantifying explained variance in multilevel models: An integrative framework for defining r-squared measures. *Psychological Methods*, 24(3):309–338.

Karen Robson and David Pevalin. 2016. *Multilevel Modeling in Plain Language*, first edition. SAGE Publications Ltd.

Mairead Shaw, Jason D. Rights, Sonya S. Sterba, and Jessica Kay Flake. 2022. r2mlm: An r package calculating r-squared measures for multilevel models. *Behavior Research Methods*, 55:1942–1964.

Nicolas Sommet and Davide Morselli. 2021. Keep calm and learn multilevel linear modeling: A three-step procedure using spss, stata, r, and mplus. *International Review of Social Psychology*, 34(1).

Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. LinkBERT: Pretraining language models with document links. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8003–8016, Dublin, Ireland. Association for Computational Linguistics.

# A Privacy-Preserving Corpus for Occupational Health in Spanish: Evaluation for NER and Classification Tasks

**Claudio Aracena**[1,5]**, Luis Miranda**[2,5]**, Thomas Vakili**[3]**, Fabián Villena**[4,5]**,**
**Tamara Quiroga**[2,5]**, Fredy Núñez-Torres**[6]**, Victor Rocco**[7]**, and Jocelyn Dunstan**[2,5]

[1]Faculty of Physical and Mathematical Sciences, University of Chile
[2]Department of Computer Science, Pontifical Catholic University of Chile
[3]Department of Computer and Systems Sciences, Stockholm University
[4]Department of Computer Science, University of Chile
[5]Millennium Institute Foundational Research on Data (IMFD), Chile
[6]Department of Language Science, Pontifical Catholic University of Chile
[7]Chilean Safety Association (ACHS), Chile
`claudio.aracena@uchile.cl, lmirandn@uc.cl, thomas.vakili@dsv.su.se,`
`fvillena@imfd.cl, t.quiroga@uc.cl, frnunez@uc.cl, varoccoc@achs.cl,`
`jdunstan@uc.cl`

## Abstract

Annotated corpora are essential to reliable natural language processing. While they are expensive to create, they are essential for building and evaluating systems. This study introduces a new corpus of 2,869 medical and admission reports collected by an occupational insurance and health provider. The corpus has been carefully annotated for personally identifiable information (PII) and is shared, masking this information. Two annotators adhered to annotation guidelines during the annotation process, and a referee later resolved annotation conflicts in a consolidation process to build a gold standard subcorpus. The inter-annotator agreement values, measured in $F_1$, range between 0.86 and 0.93 depending on the selected subcorpus. The value of the corpus is demonstrated by evaluating its use for NER of PII and a classification task. The evaluations find that fine-tuned models and GPT-3.5 reach $F_1$ of 0.911 and 0.720 in NER of PII, respectively. In the case of the insurance coverage classification task, using the original or de-identified corpus results in similar performance. The annotated data are released in de-identified form.

## 1 Introduction

Text plays a relevant role in healthcare since it is one of the richest forms of information inside electronic health records (Dalianis, 2018). Therefore, developing tools for processing and analyzing clinical text is an important goal of clinical natural language processing (NLP). However, one of the challenges when processing clinical text is the appearance of PII, such as names, locations, and identification numbers. To develop tools that can help the clinical community process text, researchers and developers need to access clinical text in a privacy-preserving manner for the patients involved. Otherwise, patients' rights are being violated.

A common way to share clinical text without violating patients' rights is to publish a de-identified version of a clinical corpus. Some of the most known clinical datasets are the MIMIC (Multi-parameter Intelligent Monitoring for Intensive Care) databases (Moody and Mark; Saeed et al., 2011; Johnson et al., 2016, 2023). These databases contain not just clinical text from critical care units but also the whole structure and data from their databases.

The previously described datasets are uncommon in languages other than English (Névéol et al., 2018). In particular, for Spanish, few clinical annotated corpus have been released. Some examples are: CANTEMIST (Miranda-Escalada et al., 2020), an annotated corpus of oncology reports; CT-EBM-SP (Campillos-Llanos et al., 2021), an annotated corpus of clinical trials; NUBes (Lima Lopez et al., 2020) an annotated corpus with negation and uncertainty entities in anonymized health records; and the Chilean waiting list corpus (Báez et al., 2020; Báez et al., 2022), an annotated corpus of referrals for the Chilean waiting list.

This work presents a corpus for occupational health in Spanish. Occupational health is an area of work in public health to promote and maintain the highest degree of physical, mental, and social well-being of workers in all occupations (World Health Organization, 2023). Occupational insurance and health providers collect patient data

111

whenever patients face a work-related accident or disease. This data is used to deliver better treatment and to decide if an occupational insurer will cover a patient.

The corpus presented in this work is similar to MEDDOPROF (Lima-López et al., 2021), an annotated corpus of occupations in clinical texts in Spanish. However, MEDDOPROF focuses on annotating only occupations, while our corpus also annotates PII. In that sense, our corpus is comparable to MEDDOCAN (Marimon et al., 2019), one of the few freely available clinical datasets for PII identification in Spanish. But there are two main differences, MEDDOCAN is a synthetic corpus, while our corpus uses actual data, and PII are masked.

Our team holds an agreement with one of the biggest occupational insurance and health providers in Chile, giving us access to their data for research purposes. Hence, the corpus introduced in this paper is a clinical corpus containing information that must be protected. The annotation procedures outlined in Section 3 describe the precautions taken to minimize the privacy risks described in Section 2. Later, Section 4 reports the experiments run with the original and de-identified corpus, including NER and classification tasks, and Section 5 shows the results and discussion of the experiments. Finally, Section 6 states the main conclusion and future work that can be done.

The main contributions of this work are:

- Publicly available pseudonymous corpus of 2,869 medical and admission reports.

- Performance comparison between fine-tuning in existing synthetic clinical corpus and our corpus for NER of PII.

- Performance comparison of a downstream task between fine-tuning in our corpus with and without PII.

## 2 Related Research

### 2.1 Privacy in NLP

With the dominance of data-driven approaches to NLP, state-of-the-art results are attained by relying on large corpora. This tendency has been further compounded with the introduction of transformer models. It is not uncommon to read about models trained using many gigabytes of textual data. However, datasets of that scale are too large to be manually audited. This means they typically contain large amounts of PII, which is a privacy risk. The parameter sizes of modern transformer models compound this risk by providing ample opportunity for training data to be memorized.

The risks of memorization in transformer models have been demonstrated through mounting attacks on pre-trained language models. Carlini et al. (2021) demonstrated that it was possible to extract memorized sequences of PII from the model GPT-2 (Brown et al., 2020). These kinds of training data extraction attacks have been repeated for other models as well, with varying success (Huang et al., 2022). Other researchers have focused on determining whether models are susceptible to membership inference attacks. These attacks are less ambitious, aiming to determine if a given datapoint was used to train a model. Such attacks have been successful even when targeting models for which training data extraction has failed (Lehman et al., 2021; Vakili and Dalianis, 2021), as demonstrated by Mireshghallah et al. (2022).

Although privacy in the context of language models is difficult to measure (Vakili and Dalianis, 2023) or even define (Brown et al., 2022), any risk of training data leakage threatens privacy. While privacy is a right that should always be protected, it is an especially pertinent value when dealing with data from sensitive sources, as is often the case in the clinical domain.

### 2.2 De-Identification

One way of reducing the privacy risks of using sensitive corpora for training is by de-identifying the data. This entails finding sensitive spans of texts and sanitizing them. When corpora are large, this can be done through automated means. Automatic de-identification is a process that typically relies on NER models to detect sensitive entities and then handle them in various ways. Automatic de-identification has been shown to decrease privacy risks while preserving the utility of the data both for fine-tuning and pre-training purposes (Verkijk and Vossen, 2022; Vakili et al., 2023). However, the impact on utility may vary depending on the task, the sanitization strategy and the quality of the underlying NER model (Berg et al., 2020; Lothritz et al., 2023). Crucially, a well-performing automatic de-identifier needs a high-quality PII dataset to train a sufficiently powerful NER model.

MEDDOCAN (Marimon et al., 2019) is one of

few freely available clinical dataset for PII identification in Spanish. The corpus comprises 1,000 synthetic documents describing fictional patients and is annotated for a wide range of PII. It was created for a shared task in which several systems were able to attain impressive $F_1$ scores reaching over 0.96. However, the documents were synthetically created for the shared task. A consequence of this is that the documents have certain artifacts that may make classification easier, but that may be absent in data encountered elsewhere. For example, MEDDOCAN documents always begin with a structured list of PII describing the patient. These include the patient's name, address, and the date of their imagined visit.

As with MEDDOCAN, the corpora annotated for PII typically originate from one or a few sources. This means that there is a risk that the models trained using the data overfit to peculiarities found in the specific datasets. Thus, it is not always clear that the NER models will generalize and be as effective at detecting sensitive information in data from other institutions unseen during training. Previous studies (Yang et al., 2019; Bridal et al., 2022) have found that performance may decrease when using data from new sources and that mismatches in annotation guidelines may make results difficult to interpret. Cross-institutional evaluations are challenging because of legal and ethical barriers to data sharing. In this paper, we not only perform such an evaluation but make our data available to other researchers interested in evaluating the cross-institutional validity of their systems. Furthermore, the data are carefully de-identified and audited by humans, meaning there is a high degree of confidence that the data are safe to share.

## 3 Corpus

In Chile, occupational insurance and health providers actively address work-related health problems during commuting or within the workplace. The core of this procedure involves creating a document known as an admission report, in which an administrative employee compiles a narrative summary of the events surrounding the incident. After this process, a medical report called anamnesis is generated. This new document is a clinical report where healthcare professionals register the clinical details of the affected patient and the specifics of the problem from a medical perspective.

In this work, we compiled a dataset of 3,000 work-related accidents. Typically, each case includes both a medical report and an admission report; however, there are instances where only one of the reports is available. As a result, we constructed an annotated corpus consisting of 2,869 documents, divided into 1,383 medical reports (anamnesis) and 1,486 admission reports. These documents are presented in a free-text format, enabling a rich and diverse range of textual content (it is noteworthy that many contain PII).

Table 1 provides a detailed analysis of corpus statistics, differentiating between medical and admission reports, and drawing a comparison between our comprehensive annotated corpus and the MEDDOCAN annotated corpus. While the MEDDOCAN corpus comprises of a smaller number of documents, it contains over twice the number of tokens and more than three times the quantity of entities compared to our dataset. This disparity can be attributed to the synthetic nature of the MEDDOCAN corpus, intentionally designed to incorporate a substantial volume of PII. Nevertheless, as outlined in Section 2.2, it is important to note that MEDDOCAN is a semi-structured corpus, which is reflected in its comparatively lower lexical diversity in contrast to our corpus.

Conversely, within our dataset, we noted that the admission report typically demonstrates a more pronounced structural organization than to the medical report. This results in a reduced lexical variety, as illustrated in Table 1.

### 3.1 Annotation Procedure

The annotation process consisted of three distinct stages. We developed a preliminary version of the annotation guidelines in the initial stage by thoroughly reviewing existing guidelines and studies about NER in Spanish or NER of PII (Dalianis and Velupillai, 2010; Báez et al., 2020; Marimon et al., 2019). We also integrated insights from the Health Insurance Portability and Accountability Act (HIPAA) (Office of the Federal Register, National Archives and Records Administration, 1996), a U.S. law defining 18 personal identifiers in Clinical Health Records.

In the second stage, one annotator annotated the entire corpus. This task included continually refining the annotation guidelines by examining encountered scenarios and ongoing discussions.

In the third and final stage, armed with well-

| Metric | Total | Med. | Adm. | MEDDOCAN |
|---|---|---|---|---|
| Documents | 2,869 | 1,383 | 1,486 | 1,000 |
| Tokens | 243,537 | 125,404 | 118,147 | 508,340 |
| Vocabulary | 18,261 | 14,483 | 6,018 | 19,699 |
| Lexical diversity | 7.5% | 11.5% | 5.1% | 3.8 % |
| Tok. per doc. | 85± 37 | 91± 53 | 79±18 | 508 ± 47 |
| Ent. per doc. | 2.1 ± 1.7 | 2.3±1.8 | 1.8±1.7 | 32.5±1.9 |
| Annotated tokens | 8,447 | 5,194 | 3,253 | 42,254 |
| Entities | 5,895 | 3,152 | 2,743 | 22,795 |

Table 1: Corpus statistics divided by medical and admission reports and comparison with MEDDOCAN.

consolidated annotation guidelines, a second annotator successfully annotated 956 documents within the corpus. This comprised 496 admission reports and 460 medical reports. This phase marked a significant milestone in our annotation process, allowing us to refine further and enhance the quality of our annotated data.

After the annotation process, we implemented a consolidation process to resolve disagreements between the first and second annotators. Each annotation underwent a comprehensive review by a team of three researchers: the two annotators and a referee responsible for making the final decision. This team examined and discussed each annotation, engaging in detailed deliberations to reach a consensus. This review process resulted in the creation of a gold standard dataset comprising 956 documents.

## 3.2 Annotation Scheme

The annotation scheme for this research exclusively encompasses non-overlapping entities. In other words, each token can have at most one associated entity. After careful consideration, we have utilized 11 entities shown in Table 2.

We incorporated all the entities proposed by Dalianis and Velupillai (2010) plus extra ones described in the next paragraphs. The annotators in this study drew inspiration from HIPAA guidelines to shape these entities, making specific modifications through their discussions. However, concerning the *Location* label, the authors unified the tags for *Country*, *Municipality*, *Street address*, and *Town* into a single category called *Location*. In contrast, we decided to preserve the *Organization* as a distinct entity, which we named *Institution*. The fact that our dataset frequently included institution names influenced this choice, usually related to the institution where the person works but

does not necessarily correspond to a location.

Given the frequent occurrence of patient occupation data within our dataset, we introduced the *Occupation* label, as outlined in the MEDDOCAN guidelines (Marimon et al., 2019). The primary motive behind its inclusion is the presence of particular occupations in the procedure annotations, suggesting that individuals could be identified based on their occupation.

Furthermore, due to the [country redacted for anonymity] civil registration origin of the data, we introduced the *Personal ID* label, which denotes a unique identification number allocated to individuals and legal entities for tax and spread use for administrative purposes.

Analyzing the annotated entities in Table 2, we detail the number of entities categorized by their respective entity classes across each subcorpus. Notably, the quantity of entities significantly fluctuates depending on their class. For instance, the most frequently occurring entity is *Full Date*, predominantly present in admission reports, and the second most prevalent entity is *Occupation*, primarily sourced from medical reports. In contrast, the *Phone Number* tag is exceptionally rare, appearing only three times throughout the entire corpus.

Furthermore, Figure 1 illustrates the token frequency distribution for each entity within the subcorpus and the distribution of annotated entities per document. Concerning token frequency, it is noteworthy that distinct subcorpora exhibit varying distributions. Generally, entities consist of a single token, but there are multi-token entities. In the admission report corpus, entities like *Location* and *Health Care Unit* are mostly composed of more than one token. Additionally, in the medical report subcorpus, entities such as *Occupation*, *Institution*, and *Location* are multi-token.

Figure 1: Frequency distribution of (left) annotated entities per document by subcorpus, and (right) tokens per entity across the subcorpus.

| Entity | Total | Med. | Adm. |
|--------|-------|------|------|
| Age | 195 | 194 | 1 |
| Institution | 242 | 217 | 25 |
| Health Care Unit | 394 | 348 | 46 |
| Date Part | 485 | 473 | 12 |
| Full Date | 1981 | 563 | 1418 |
| First Name | 402 | 21 | 381 |
| Last Name | 358 | 53 | 305 |
| Location | 197 | 63 | 134 |
| Occupation | 1634 | 1214 | 420 |
| Phone Number | 3 | 2 | 1 |
| Personal ID | 4 | 4 | 0 |

Table 2: Number of entities by entity class and if they are in the medical or administrative subcorpus.

Conversely, when considering the distribution of annotated entities per document, Figure 1 reveals that, on the whole, documents tend to contain a relatively small number of entities. However, it's worth noting that admission reports, on average, contain one *Full Date* entity per document, while medical reports, on average, feature one *Occupation* entity per document.

### 3.3 Annotation Guidelines

Three researchers collaboratively drafted a comprehensive document outlining the annotation guidelines: the annotator responsible for the entire corpus, a linguist, and a computer science pro-

fessor. It resulted from a thorough review of literature (Báez et al., 2020; Marimon et al., 2019; Dalianis and Velupillai, 2010; Office of the Federal Register, National Archives and Records Administration, 1996) and discussions on annotation casuistry, where regular meetings were held to ensure that the guidelines maintained linguistic and syntactic coherence while enhancing privacy protection without undermining the texts' narrative. The current version of the annotation guidelines is freely available[1].

Building upon the framework established by Báez et al. (2020) for annotation guidelines, we categorize the rules into two sections: general rules, which have universal application to all entities, and specific rules customized for each entity. Within the general and specific rule sections, we further distinguish between positive rules (guiding what should be annotated) and negative rules (clearly outlining what should not be tagged or what constitutes an incorrect annotation). Finally, the guidelines provide informative explanations regarding typical scenarios encountered within the dataset.

We elucidated two general rules, refraining from incorporating trailing punctuation marks or white spaces after entities. Furthermore, we emphasized the importance of tagging each entity

---

[1] https://totoiii.github.io/clinical_deidentification_guideline/

with the utmost specificity to ensure the most comprehensive coverage of the entity.



Figure 2: Example of an annotated document where PII has been modified. Translation: *Entry - the 62 years old, PMH: asthma, DOA: 05/20/2032, Allergies: None. Principal at South High School. The patient reports that on Wednesday, 11/02, while working in a classroom, he began experiencing shortness of breath, prompting him to seek care at Saint John Hospital. He has been experiencing recurrent asthma attacks recently and is not using an inhaler.*

Finally, Figure 2 presents an example document with annotations drawn from the existing corpus and modified for explanatory purposes, with all PII appropriately modified.

### 3.4 Inter-Annotator Agreement

We evaluated the challenge of achieving consistent annotations by assessing inter-annotator agreement (IAA). Specifically, the macro $F_1$ was employed to assess and compare the annotations. Table 3 depicts the agreements for each comparison within the different subcorpora. These comparisons entail assessments between annotator 1 and annotator 2 and between each annotator and the gold standard corpus.

Furthermore, Figure 3 visualizes the IAA for various entity classes, except for (*Age*, *Phone Number*, and *Personal ID*) that have too few instances. The figure highlights that, in most cases, more favorable agreement results are evident in the admission report as compared to the medical report. This can be attributed to the slightly more structural organization found in the admission report subcorpus when compared to the medical report subcorpus.

### 3.5 Masking Procedure

A masking process was carried out to share our corpus without private or sensitive information. The masking process adds a mask using the tag `"__entity_name__"` for every entity, where the entity name corresponds to the name of an entity, e.g., First Name.

| | Global | Medical | Admission |
|---|---|---|---|
| A1 - A2 | 0.90 | 0.86 | 0.93 |
| GS - A1 | 0.97 | 0.94 | 0.98 |
| GS - A2 | 0.92 | 0.90 | 0.93 |

Table 3: Macro $F_1$ agreements for each comparison and subcorpus. Where A1 is annotator 1, A2 is annotator 2, and GS is gold standard corpus.



Figure 3: Macro $F_1$ score for IAA for each entity and subcorpus.

## 4 Experiments

The three experiments described in this section demonstrated the value of the corpus:

1. GPT-3.5 is used to detect PII and compare it with human annotators.

2. The corpus is used to train and evaluate NER models for privacy-preserving purposes.

3. The corpus is used to perform an insurance coverage classification task.

### 4.1 NER via Few-Shot In-Context Learning

The utilization and experimentation with Large Language Models (LLMs) in NER tasks hold profound significance in the realm of NLP. These models, equipped with their vast contextual understanding, have the potential to greatly enhance the performance and efficiency of identifying named entities within text.

In this experiment, we employed the gpt-3.5-turbo model (OpenAI, 2023) through Microsoft

Azure to conduct few-shot NER on the entire corpus. The prompt given to the model involved explaining the 11 entities outlined in Section 3.2, along with providing brief descriptions of the content associated with each tag. Additionally, the prompt included the desired output format, which involves annotating the input text using a markup language format, as follows: `...<Entity class>Named entity body</Entity class>....`

Furthermore, we generated and instructed the model with five distinct examples for each subcorpus. These examples were crafted in alignment with the unique characteristics of their corresponding subcorpus.

## 4.2 Training and Evaluating NER

As explained in Section 2.2, an important use case for NER is for automatic de-identification of sensitive data. However, assessing the cross-institutional validity of such models is difficult due to data scarcity, which is especially dire in languages other than English. In this experiment, the transferability of performance gained through training models using the MEDDOCAN corpus was evaluated using our new corpus. Models were trained using either MEDDOCAN data or our corpus and then evaluated on the curated gold standard part of the corpus.

A wide range of models for Spanish language modeling were selected. The best base models for Spanish NER suggested by Agerri and Agirre (2023) were fine-tuned for PII detection. The models chosen were the multilingual model mDeBERTaV3 (He et al., 2023) and the monolingual Spanish model IXABERTes-v2[2] that is based on RoBERTa (Liu et al., 2019). Fine-tuned models were created using MEDDOCAN as well as the corpus introduced in this study.

The PII tags defined for MEDDOCAN and our corpus differed in a few ways. Before training the models, the tagsets were harmonized. This involved translating MEDDOCAN tags into their counterparts in our corpus. If a MEDDOCAN tag lacked a counterpart, it was ignored. The procedure also involved collapsing labels found in our corpus that were not distinguished in MEDDO-CAN. The distinction between first and last names, and between partial and full dates, was present in our corpus but not in MEDDOCAN. Conse-

quently, they were collapsed into the tags *Name* and *Date*. Both datasets were then converted into the IOB format[3].

After training each model configuration for five epochs, the best checkpoint was selected based on the $F_1$ score on the validation set. The selected models were evaluated on the gold standard described in Section 3 and a held-out test set from MEDDOCAN. The performance difference when testing models on the unseen dataset indicates how much they generalize to novel data. Our corpus's classification results were compared with those obtained by prompting GPT-3.5.

## 4.3 Insurance Coverage Classification

The insurance coverage classification task is selected to assess the impact of de-identifying PII on downstream tasks' performance. This task aims to classify the insurance coverage decision of the occupational insurance provider. Following Aracena et al. (2023), the process of building a classifier consists of using the pre-trained model bsc-bio-ehr-es[4] (Carrino et al., 2022) as a base model. Then, a fine-tuning step is carried out, in which the documents from the corpus with their corresponding label for the insurance coverage decision are used for this purpose [5]. Lastly, the fine-tuned model is evaluated in other cases not part of the corpus.

The previous process is implemented for the original and de-identified corpus, and also for the admission subcorpus, the medical subcorpus, and both combined. The de-identification was performed by replacing each sensitive entity with its class name.

## 5 Results

This section shows the results of the experiments and discusses the implications for NER and downstream tasks.

## 5.1 NER for De-Identification

Four fine-tuned models were trained based on the pre-trained mDeBERTaV3 and IXABERTes-v2 models. Each model was trained on either the MEDDOCAN data, or our corpus. Table 4 shows

---

[2]http://www.deeptext.eus/es/node/3

[3]Specifically, we use the version of IOB that reserves *B* for entities spanning multiple tokens.

[4]https://huggingface.co/PlanTL-GOB-ES/roberta-base-biomedical-clinical-es

[5]The labels for the classification task are not part of the released corpus.

| Base model | Training data | Test $F_1$ score | |
| --- | --- | --- | --- |
| | | MEDDOCAN | Our corpus |
| mDeBERTaV3 | Our corpus | 0.400 | 0.853 |
| mDeBERTaV3 | MEDDOCAN | 0.990 | 0.498 |
| IXABERTes-v2 | Our corpus | 0.368 | 0.834 |
| IXABERTes-v2 | MEDDOCAN | 0.990 | 0.381 |
| gpt-3.5-turbo | N/A | N/A | 0.720 |

Table 4: $F_1$ scores for each combination of model, training dataset and testing dataset. gpt-3.5-turbo was not fine-tuned but was accessed through an API and prompted using a few-shot approach targeting the new corpus.

| Base model | Fine-tuning data | Test $F_1$ score | |
| --- | --- | --- | --- |
| | | Original | De-identified |
| bsc-bio-ehr-es | Admission | 0.726 ± 0.015 | 0.708 ± 0.004 |
| | Medical | 0.738 ± 0.006 | 0.743 ± 0.008 |
| | Admission+Medical | 0.750 ± 0.002 | 0.763 ± 0.006 |

Table 5: $F_1$ scores for classification task in the test set.

the results of evaluating the models on the test version of their training data and the test set of the other dataset. The models perform substantially worse in all four cases when evaluated on novel data. This indicates a clear mismatch between the two datasets, even though the task they represent is ostensibly equivalent.

It is not obvious if the mismatch between our corpus and MEDDOCAN is due to the synthetic nature of MEDDOCAN or stems from an inherent diversity in how PII are represented in the clinical domain. A truly cross-institutional PII tagger for de-identification purposes should perform well on a diverse range of datasets. The new corpus thus functions as a source of training data, and as a benchmark to evaluate the generalizability of NER models trained on other data sources.

Additionally, we show the performance of gpt-3.5-turbo when performing few-shot NER through in-context learning on the new corpus. Even though it does not show the best results, it performs better than the cross-institutional taggers, which is still remarkable considering that just a few examples were given to understand the task. However, similar to a previous experience (Wang et al., 2023), the performance of gpt-3.5-turbo for NER tasks is not state-of-the-art.

The gpt-3.5-turbo outputs sometimes deviated from expectations by altering the original text in various ways. These alterations included fixing

misspelled words or introducing punctuation not in the original text. This posed a significant challenge, resulting in misaligning the original annotations with the model-generated ones. To address this issue, we analyzed in detail the disparities between the original text and the model-modified text. We then adjusted the positions of tokens for each annotated entity in the model output, enabling us to make precise comparisons between the annotations.

### 5.2 Insurance Coverage Classification

Six fine-tuning configurations were used to train models, three with the original corpus and three with the de-identified corpus. Admission subcorpus, medical subcorpus, and both combined were used to fine-tune models in each type of corpus. For every fine-tuning configuration, three random seeds were used to check the variability of the results, and for each run, three epochs were used. Table 5 shows the insurance coverage classification results. The positive class is the decision not to cover a patient, as this is the less frequent class.

Under the described conditions, none to little differences were found between using the original or the de-identified corpus. These results suggest that using a de-identified corpus for downstream tasks should not impact the performance. However, depending on the task under study, this situation may vary. Aracena et al. (2023) reported bet-

ter performance for the same task, reaching 0.963 of AUC. This is due to the amount of data used for fine-tuning, which is more than 200 times bigger than this study.

# 6 Conclusions

This work introduces a novel corpus of admission and medical reports retrieved from an insurance and health provider specialized in occupational health. The annotation of PII and the subsequent de-identification process highlight the importance of releasing data considering ethical and privacy matters. This corpus is released[6] in de-identified form, where all sensitive entities are replaced with their class names.

Our exploration of the corpus has revealed its inherent value in named entity recognition and classification tasks. The insights gained through these analyses not only contribute to the existing body of knowledge but also hold practical implications for improving information extraction within the specified domain.

As future work, one promising avenue involves exploring synthetic data to replace masked PII entities. This approach has the potential to not only safeguard privacy but also the possibility of building robust models that can be applied in diverse real-world scenarios with synthetic data.

# Acknowledgements

# References

Rodrigo Agerri and Eneko Agirre. 2023. Lessons learned from the evaluation of Spanish Language Models. *Procesamiento del Lenguaje Natural*, 70(0):157–170. Number: 0.

Claudio Aracena, Nicolás Rodríguez, Victor Rocco, and Jocelyn Dunstan. 2023. Pre-trained language models in Spanish for health insurance coverage. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 433–438, Toronto, Canada. Association for Computational Linguistics.

Pablo Báez, Felipe Bravo-Marquez, Jocelyn Dunstan, Matías Rojas, and Fabián Villena. 2022. Automatic extraction of nested entities in clinical referrals in spanish. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(3):1–22.

Hanna Berg, Aron Henriksson, and Hercules Dalianis. 2020. The Impact of De-identification on Downstream Named Entity Recognition in Clinical Text. In *Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis*, pages 1–11, Online. Association for Computational Linguistics.

Olle Bridal, Thomas Vakili, and Marina Santini. 2022. Cross-Clinic De-Identification of Swedish Electronic Health Records: Nuances and Caveats. In *Proceedings of the Workshop on Ethical and Legal Issues in Human Language Technologies and Multilingual De-Identification of Sensitive Data In Language Resources within the 13th Language Resources and Evaluation Conference*, pages 49–52, Marseille, France. European Language Resources Association.

Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr. 2022. What Does it Mean for a Language Model to Preserve Privacy? In *2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, pages 2280–2292, New York, NY, USA. Association for Computing Machinery.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Pablo Báez, Fabián Villena, Matías Rojas, Manuel Durán, and Jocelyn Dunstan. 2020. The Chilean Waiting List Corpus: a new resource for clinical Named Entity Recognition in Spanish. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 291–300, Online. Association for Computational Linguistics.

Leonardo Campillos-Llanos, Ana Valverde-Mateos, Adrián Capllonch-Carrión, and Antonio Moreno-Sandoval. 2021. A clinical trials corpus annotated with UMLS entities to enhance the access to evidence-based medicine. *BMC Medical Informatics and Decision Making*, 21(1):69.

---

[6]The data are available upon request at:https://zenodo.org/records/11035754

Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. Extracting Training Data from Large Language Models. In *Proceedings of the 30th USENIX Security Symposium*, pages 2633–2650.

Casimiro Pio Carrino, Joan Llop, Marc Pàmies, Asier Gutiérrez-Fandiño, Jordi Armengol-Estapé, Joaquín Silveira-Ocampo, Alfonso Valencia, Aitor Gonzalez-Agirre, and Marta Villegas. 2022. Pretrained Biomedical Language Models for Clinical NLP in Spanish. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 193–199, Dublin, Ireland. Association for Computational Linguistics.

Hercules Dalianis. 2018. *Clinical Text Mining: Secondary Use of Electronic Patient Records*. Springer International Publishing.

Hercules Dalianis and Sumithra Velupillai. 2010. De-identifying swedish clinical text - refinement of a gold standard and experiments with conditional random fields. *Journal of Biomedical Semantics*, 1(1):6.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. DeBERTaV3: Improving DeBERTa using ELECTRA-Style Pre-Training with Gradient-Disentangled Embedding Sharing.

Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022. Are Large Pre-Trained Language Models Leaking Your Personal Information? In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2038–2047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J. Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, Li wei H. Lehman, Leo A. Celi, and Roger G. Mark. 2023. MIMIC-IV, a freely accessible electronic health record dataset. *Scientific Data*, 10(1).

Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific Data*, 3(1).

Eric Lehman, Sarthak Jain, Karl Pichotta, Yoav Goldberg, and Byron Wallace. 2021. Does BERT Pretrained on Clinical Notes Reveal Sensitive Data? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 946–959, Online. Association for Computational Linguistics.

Salvador Lima Lopez, Naiara Perez, Montse Cuadros, and German Rigau. 2020. NUBes: A Corpus of Negation and Uncertainty in Spanish Clinical Texts. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5772–5781, Marseille, France. European Language Resources Association.

Salvador Lima-López, Eulàlia Farré-Maduell, Antonio Miranda-Escalada, Vicent Brivá-Iglesias, and Martin Krallinger. 2021. NLP applied to occupational health: MEDDOPROF shared task at IberLEF 2021 on automatic recognition, classification and normalization of professions and occupations from medical texts. *Procesamiento del Lenguaje Natural*, 67(0):243–256.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv:1907.11692 [cs]*. ArXiv: 1907.11692.

Cedric Lothritz, Bertrand Lebichot, Kevin Allix, Saad Ezzini, Tegawendé Bissyandé, Jacques Klein, Andrey Boytsov, Clément Lefebvre, and Anne Goujon. 2023. Evaluating the Impact of Text De-Identification on Downstream NLP Tasks. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 10–16, Tórshavn, Faroe Islands. University of Tartu Library.

Montserrat Marimon, Aitor Gonzalez-Agirre, Ander Intxaurrondo, Heidy Rodriguez, Jose Lopez Martin, Marta Villegas, and Martin Krallinger. 2019. Automatic De-identification of Medical Texts in Spanish: the MEDDOCAN Track, Corpus, Guidelines, Methods and Evaluation of Results. In *IberLEF@ SEPLN*, pages 618–638.

A Miranda-Escalada, E Farré, and M Krallinger. 2020. Named entity recognition, concept normalization and clinical coding: Overview of the Cantemist track for cancer text mining in spanish, corpus, guidelines, methods and results. In *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2020), CEUR Workshop Proceedings*.

Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri. 2022. Quantifying privacy risks of masked language models using membership inference attacks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8332–8347, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

G.B. Moody and R.G. Mark. A database to support development and evaluation of intelligent intensive care monitoring. In *Computers in Cardiology 1996*. IEEE.

Aurélie Névéol, Hercules Dalianis, Sumithra Velupillai, Guergana Savova, and Pierre Zweigenbaum.

2018. Clinical Natural Language Processing in languages other than English: opportunities and challenges. *Journal of Biomedical Semantics*, 9(1):12.

Office of the Federal Register, National Archives and Records Administration. 1996. Public law 104 - 191 - health insurance portability and accountability act of 1996.

OpenAI. 2023. Models. https://platform.openai.com/docs/models/gpt-3-5. [Accessed 20-10-2023].

Mohammed Saeed, Mauricio Villarroel, Andrew T. Reisner, Gari Clifford, Li-Wei Lehman, George Moody, Thomas Heldt, Tin H. Kyaw, Benjamin Moody, and Roger G. Mark. 2011. Multiparameter intelligent monitoring in intensive care II (MIMIC-II): A public-access intensive care unit database. *Critical Care Medicine*, 39(5):952–960.

Thomas Vakili and Hercules Dalianis. 2021. Are Clinical BERT Models Privacy Preserving? The Difficulty of Extracting Patient-Condition Associations. In *Proceedings of the AAAI 2021 Fall Symposium on Human Partnership with Medical AI: Design, Operationalization, and Ethics (AAAI-HUMAN 2021)*.

Thomas Vakili and Hercules Dalianis. 2023. Using Membership Inference Attacks to Evaluate Privacy-Preserving Language Modeling Fails for Pseudonymizing Data. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, NEALT Proceedings Series, pages 318–323, Tórshavn, Faroe Islands. University of Tartu Library.

Thomas Vakili, Aron Henriksson, and Hercules Dalianis. 2023. End-to-End Pseudonymization of Fine-Tuned Clinical BERT Models.

Stella Verkijk and Piek Vossen. 2022. Efficiently and Thoroughly Anonymizing a Transformer Language Model for Dutch Electronic Health Records: a Two-Step Method. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1098–1103, Marseille, France. European Language Resources Association.

Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*.

World Health Organization. 2023. Occupational health. https://www.who.int/health-topics/occupational-health. [Online; accessed 16-October-2023].

Xi Yang, Tianchen Lyu, Qian Li, Chih-Yin Lee, Jiang Bian, William R. Hogan, and Yonghui Wu. 2019. A study of deep learning methods for de-identification of clinical notes in cross-institute settings. *BMC Medical Informatics and Decision Making*, 19(5):232.

# DERA: Enhancing Large Language Model Completions
# with Dialog-Enabled Resolving Agents

**Varun Nair**[*]    **Elliot Schumacher**[*]    **Geoffrey Tso**    **Anitha Kannan**
Curai Health

## Abstract

Large language models (LLMs) have emerged as valuable tools for many natural language understanding tasks. In safety-critical applications such as healthcare, the utility of these models is governed by their ability to generate factually accurate and complete outputs. In this work, we present dialog-enabled resolving agents (DERA). DERA is a paradigm made possible by the increased conversational abilities of LLMs. It provides a simple, interpretable forum for models to communicate feedback and iteratively improve output. We frame our dialog as a discussion between two agent types – a *Researcher*, who processes information and identifies crucial problem components, and a *Decider*, who has the autonomy to integrate the *Researcher*'s information and makes judgments on the final output.

We test DERA against three clinically-focused tasks, with GPT-4 serving as our LLM. DERA shows significant improvement over the base GPT-4 performance in both human expert preference evaluations and quantitative metrics for medical conversation summarization and care plan generation. In a new finding, we also show that GPT-4's performance (70%) on an *open-ended* version of the MedQA question-answering (QA) dataset (Jin et al. (2021), USMLE) is well above the passing level (60%), with DERA showing similar performance. We will release the open-ended MedQA dataset.

## 1 Introduction

Large language models (LLMs; Brown et al. (2020); Lewis et al. (2020)) are deep-learning models trained to predict natural language text conditioned on an input. These models have led to advances in natural language performance far beyond traditional language modeling tasks, including on few-shot learning (Brown et al., 2020) and multimodal tasks (Driess et al., 2023). Within the realm of medicine, LLM-powered methods have shown improvements in medical tasks such as question answering (Singhal et al., 2022; Liévin et al., 2022), information extraction (Agrawal et al., 2022), and summarization (Chintagunta et al., 2021).

LLM-powered methods use natural language instructions called *prompts*. These instruction sets often include a task definition, rules the predictions must follow, and few-shot examples of the task input and output (Reynolds and McDonell, 2021; Brown et al., 2020). The ability of generative language models to create output based on natural language instructions (or prompts) removes the need for task-specific training (Min et al., 2022) and allows non-experts to build upon this technology.

While many tasks can be formulated as a single prompt, later work has shown that breaking down single tasks into sub-tasks (called *chaining*) has benefits in terms of task performance and interpretability (Wu et al., 2022). Chain-of-thought (CoT) (Wei et al., 2022) is one example of a chaining strategy in which the model is prompted to think through a problem as an expert might approach it, leading to improvements in some tasks (Liévin et al., 2022; Wang et al., 2022; Tafjord et al., 2022; Huang et al., 2022). Other chaining strategies specific to particular domains have also been developed, such as in Agrawal et al. (2022) for basic clinical tasks and in Zhu et al. (2023) for image captioning.

All of these chaining approaches attempt to coerce the correct generation from a base language model. However, one fundamental limitation of this strategy is that they are usually sequential and manually engineered for every task. Even with this complexity, chained approaches struggle with generating factually accurate text and often can include hallucinations and omissions (Maynez et al., 2020; Dziri et al., 2022; Berezin and Batura, 2022).

---

[*] The first two authors contributed equally to this work. For correspondence, please contact elliot@curai.com.

**1)** The *Decider* agent (◉) first computes some initial output for a given task.

**2)** The *Decider* and *Researcher* agent (◉) then discuss changes for alignment to task goals.

**3)** Finally, the *Decider* uses the discussed changes to compute the final resolved output.

Figure 1: Overview of DERA. The method consists of two agents–a *Researcher* and a *Decider*. The *Decider* generates an initial output for the task (step 1). Then, the *Decider* and *Researcher* work through the problem via conversation (step 2), with the *Researcher* tasked to help identify crucial problem components. The *Decider* has the autonomy to integrate the *Researcher*'s inputs and makes judgments on the final output (step 3). Neither agent has knowledge of the ideal final output.

This poses a significant hurdle when applying them to real-world scenarios, especially in the clinical domain.

The increasingly robust and realistic conversational capabilities of LLMs (OpenAI, 2023; Pal et al., 2022) leads us to ask – *can reformulating language tasks as conversations between LLM agents improve generative output?* We present a framework, DERA (Dialog-Enabled Resolving Agents), for improving performance on natural language tasks using agents tasked with refining task output through dialog. We pair an agent that generates the initial task output with one that can guide the other by suggesting areas of focus in each round of the conversation.

DERA is a task-agnostic framework that refines text generation issues such as hallucinations and omissions. The dialogue medium adds interpretability to the process and allows the generation to be refined holistically. We propose that scoping each agent in the dialog to a specific role will better enable them to focus on discrete portions of the task and ensure their partner agent stays aligned with the overall goal.

Our paper makes the following contributions:

- We introduce DERA (§ 2) - a framework for agent-agent dialog to improve performance on natural language tasks.

- We evaluate DERA on three different types of clinical tasks. Specifically, these include a medical doctor-patient conversation summarization task (§ 3), a provider-facing, careplan generation task (§4), and medical open-ended question answering tasks (§5). Each of these requires different types of textual inputs and types of knowledge to solve.

- In both human-annotated evaluations, we find that DERA outperforms base GPT-4 performance in the careplan generation and medical conversation summarization tasks on a variety of metrics. In quantitative evaluations, we find that DERA successfully corrects medical conversation summaries with large amounts of errors. Conversely, we find small to no improvement between GPT-4 performance and DERA on question-answering.

- We theorize this approach is well suited for longer-form generation tasks in which there are a lot of fine-grained details.

- To further research, we release the open-ended version of medical question-answering dataset (MedQA; Jin et al. (2021)).

## 2 DERA: Overview

DERA is a general chat framework that leverages dialog-capable agents to collaboratively work through a task (Figure 1). We focus on agent setups that work to probe knowledge sources, whether internal or external (from text, documents, etc.). We propose that pairing an information-focused agent with a decision-maker agent will lead to a higher-quality output. Furthermore, this approach allows for DERA to alternate between processing knowledge and acting upon information, as opposed to doing them concurrently.

First, we propose the use of a *Researcher* agent, shown in Orange in Figure 1. The goal of a researcher agent is to review pieces of information – which can be internal to an LLM or external – and make suggestions on what is likely to be crucial in solving the problem. As we do not have a definitive source of what is and is not relevant, we rely on an LLM's ability to identify relevancy in light of the current task. We do not treat this agent as the definitive source of truth. Rather, we task it with being helpful and constructive during the dialog.

Second, we propose the use of a *Decider* agent, shown in Green in Figure 1. In addition to starting the conversation as shown in the left part of the figure, this agent is tasked with responding to the information provided by the *Researcher* agent, and deciding whether to integrate that information into the task output. This allows an LLM to make discrete decisions in reaction to the information highlighted by the *Researcher*. At no point, however, does the *Decider* defer to the *Researcher*. This agent is ultimately responsible for the final decision. While it is tasked with reviewing all information highlighted by *Researcher*, it does not have to use any of that information.

The specific directives of each agent can vary for different tasks. For Question Answering, the *Researcher* is tasked with pulling information from the question, using the internal knowledge of an LLM alone. For summarization, the *Researcher* has access to external texts which contain the full patient encounter. Conversely, the edits to the text generation task are made incrementally by the *Decider* in the summarization task, while they are made more discretely in the question-answering task. In some settings, agents take a hybrid role, each having access to different information and jointly making decisions. Overall, the goal remains the same – that this approach allows for informa-

tion to be processed in a role-defined and iterative manner, producing better quality output. We use GPT-4 (OpenAI, 2023) as the LLM for this paper, but we propose that this approach can generalize to other LLMs[1].

We apply DERA to three natural language generation tasks. The first, medical conversation summarization (§3), probes the ability of DERA to create a summary of a doctor-patient chat. This requires the ability to identify and rewrite medically-relevant information in a concise format. The second, care plan generation (§4), tests whether DERA can generate doctor-facing suggestions for potential actions to address patient concerns. This requires similar abilities, with the added challenge of knowing the appropriate next steps for a variety of medical conditions. Finally, medical question-answering (§5) tests the ability of DERA to generate a wide variety of medical knowledge in a short format.

## 3 Medical Conversation Summarization

**Overview**  The task of medical conversation summarization is to encapsulate a patient-doctor conversation (Enarvi et al., 2020; Joshi et al., 2020; Zhang et al., 2021; Chintagunta et al., 2021). Doctors use these summaries for downstream tasks such as clinical decision-making, and hence it is important that the generated summaries are both factually accurate (no hallucinations) and complete (no omissions). We focus on summarizing patient-doctor chats into six independent sections: *Demographics and Social Determinants of Health*, *Medical Intent*, *Pertinent Positives*, *Pertinent Negatives*, *Pertinent Unknowns*, and *Medical History*. This structured format requires the model to summarize the chat while placing each piece of information in the appropriate section.

**DERA Setup**  We formulate the DERA setup for medical conversation summarization as follows. Both *Decider* and *Researcher* have access to the full medical conversation between the patient and the physician. Both agents are prompted to converse with one another. The *Decider* agent generates an initial summary of the medical conversation (Prompt 1) and shares it with the *Researcher* agent. The *Researcher* agent's role (Prompt 4) is to "read" the summary and point out any discrepancies to *Decider*. *Decider*, using Prompt 3, either

---

[1]At the time of writing, we did not have access LLMs of comparable performance.

accepts or rejects those discrepancies by agreeing with the suggestion or disagreeing and responding with some reasoning. Instead of regenerating the summary at each step of the conversation, *Decider* writes the accepted suggestions to a shared *scratchpad*, which acts like a memory that it uses at the end of the conversation to generate the final summary. The conversation terminates once *Researcher* is satisfied with the suggestions made to the scratchpad or a maximum conversation length is reached (set to 15 turns total). As the final step, the *Decider* generates (Prompt 5) the final summary using the contents of the scratchpad and the original summary. GPT-4 prompts are run with the settings mentioned in Table 5.

**Dataset** We randomly sampled 500 medical encounters from a chat-based telehealth platform. Each encounter contains the patient's age, sex, and chat conversation with a licensed medical provider. Encounters in this dataset cover a wide variety of common presentations in telehealth, including urinary tract infections, back/abdominal pains, toothaches, and others. All data is de-identified prior to experimentation. Conversations contain 27 dialog turns on average (min of 9 turns, max of 82 turns) and average 646 unigram tokens per encounter (min 42 tokens, max 2031 tokens).

**Human Expert Evaluation** To evaluate the effectiveness of DERA to generate better summaries, we conducted human evaluation studies with four licensed physicians on a random subset of 50 out of the 500 encounters described above. We sampled a smaller, random subset due to the high labeling cost induced by using expert physicians.

The licensed physicians were provided with the encounter and the two summaries. These included the initial GPT-4 generated summary and the final generated summary produced using DERA. Each physician was asked to answer three main questions in the light of the summary's clinical utility for themselves or another physician: **(1)** *Which summary do you prefer to use for the given patient and encounter?* **(2)** *What percentage of the overall clinical information in the dialog is captured by the summary?* **(3)** *What percentage of the suggestions added to the DERA scratchpad do you agree with?*

Figure 2 shows the results of our human expert evaluation. Physicians notably choose the summary produced after DERA over the initially generated summary 90% - 10%. Their preference for the DERA-produced summary is further corroborated

by the fraction of medical information captured in the final DERA summary vs. initial, as final summaries were rated as capturing "All" medical information from the patient-physician dialog in 86% of encounters vs. the initial summaries capturing "All" medical information in just 56% of encounters. In general, we also find broad agreement for the suggestions in each encounter's scratchpad: they agreed with "All" corrections suggested for a given encounter's summary 63% of the time, "Most" 14% of the time, "Some" 5% of the time, and "None" 18% of the time. On average, each scratchpad contains 2-3 suggestions.

In addition to these questions, we also asked the physician-experts the following: *If this summary were acted upon by another clinical provider, does this summary contain information that could potentially be harmful to the patient given their presentation?* (Options: Yes, No). The number of summaries containing "harmful" information drops from 2% in the initial summary to 0% in the final DERA summary. We caution against drawing generalizations from these harmfulness numbers. Our evaluations are both limited in number and drawn from a patient population specific to the telehealth platform; thus cannot predict the generalizability of these findings in other settings.

**Quantitative Evaluation** We also perform a more large-scale study without the need for human annotation. We generate GPT-4 summaries for all the 500 encounters and assume them to be ground truth. Then, we synthetically induce "corruptions" into the generated summary and use that as the initial input. These mistakes artificially lower the summary's quality and produce significant hallucinations and omissions. The goal is to quantitatively evaluate DERA's ability to write medical summaries by measuring the degree to which the *Researcher* and *Decider* agents can identify and fix "corruptions" introduced to the medical summary.

Prompt 2 contains specific instructions for generating the corruptions. We can control the level of corruption desired by passing one of three levels of corruption as a variable to our corruption prompt: low ($\frac{3}{10}$), medium ($\frac{5}{10}$), or high ($\frac{7}{10}$). The higher the corruption, the more symptoms could be rearranged. Similarly, hallucinated symptoms could be introduced, among other corruptions. See Fig. 5 for a qualitative example of this process of generating an initial summary, corrupting it, resolving with DERA, and generating a final summary.

Figure 2: Results from physician-expert evaluations on the medical conversation summarization task. (Left) Physicians choose the final summary produced by DERA over the initial GPT-4 generated summary 90% to 10%. (Center) Final DERA summaries capture far more clinical information than initial GPT-4 generated summaries, with physicians rating "All" relevant clinical information from the patient-physician chat captured in 86% of DERA summaries vs. 56% of initial GPT-4 summaries. (Right) For summary correction suggestions in the scratchpad, physicians rate agreement with All suggestions in 63% of encounters, Most in 14%, Some in 5%, and None in 18%.

| Corruption Level | Summ. Version | Pertinent Positives | Pertinent Negatives | Pertinent Unknowns | Medical History | Average |
|---|---|---|---|---|---|---|
| low ($\frac{3}{10}$) | Initial | 89.38 | 83.05 | 87.42 | 80.88 | 85.18 |
| | Baseline | 93.90 | 89.33 | 90.11 | 89.91 | 90.81 |
| | DERA | 95.65 | 96.77 | 97.10 | 97.35 | **96.71** |
| medium ($\frac{5}{10}$) | Initial | 83.12 | 81.60 | 71.14 | 73.82 | 77.42 |
| | Baseline | 92.79 | 86.57 | 89.44 | 88.38 | 89.30 |
| | DERA | 94.29 | 95.31 | 96.17 | 98.12 | **95.97** |
| high ($\frac{7}{10}$) | Initial | 68.35 | 70.07 | 68.79 | 57.27 | 66.12 |
| | Baseline | 88.34 | 83.98 | 86.52 | 86.72 | 86.39 |
| | DERA | 92.96 | 90.86 | 94.81 | 95.16 | **93.45** |

Table 1: Medical conversation summarization task: Quantitative evaluation (GPT-F1 scores) of the initial summary with errors and the DERA corrected version. We show that by introducing synthetic corruption (hallucinations, omissions, etc.) into medical summaries, DERA can resolve these corruptions at low, medium, and high levels of corruption. GPT-F1 scores for the DERA-produced summary are consistently higher than the initial summaries.



Figure 3: Care plan generation task: Results from physician-expert evaluations. (Left) Physicians choose the final care plan produced by DERA over the initial GPT-4 generated care plan 84% to 16%. (Center) Final DERA care plans capture far more of the necessary care management steps than initial GPT-4 generated care plans, with physicians rating "All" relevant steps inferred from the patient-physician chat generated in 92% of DERA care plans vs. 64% of initial GPT-4 care plans. (Right) For care plan correction suggestions in the scratchpad, physicians rate agreement with "All" suggestions in 72% of encounters, Most" in 14%, "Some" in 0%, and "None" in 14%.

126

Recent research has shown that traditional summarization metrics are not sufficient to capture nuanced changes in performance (Goyal et al., 2022). Therefore, we measure the degree to which corruptions are present by using a GPT-based metric that tracks the medical concept coverage of the medical summary, **GPT-F1** (Nair et al., 2023). GPT-F1 is computed as the harmonic mean of two sub-metrics: GPT-Recall and GPT-Precision. A GPT-F1 score of 100 implies a perfect match in medical concepts present in the query and reference text. We further describe these metrics in Appendix Section A.1. The results of our quantitative evaluation using the GPT-F1 metric are shown in Table 1.

We compare GPT-F1 on the initial summary with errors to a baseline method and the DERA corrected summary. The baseline method is a simplified version DERA in which allows just a single pass at corrections, effectively ablates the importance of back-and-forth dialogue between the DERA agents.

Note first how the higher levels of corruption manifest in the initial summary GPT-F1. As the corruption level of the initial summary increases, the initial GPT-F1 score drops. We find that DERA can produce significantly improved summaries in low, medium, and high levels of corruption, as evidenced by increases in GPT-F1 over both the initial and baseline method summaries. This suggests that the collaborative interaction between the *Researcher* and *Decider* agents identifies hallucinations and omissions and resolves them through dialog, even when many such corruptions are present.

## 4 Care Plan Generation

We also analyze the performance of DERA on the task of generating a care management plan. This care plan contains suggestions that are meant to be *physician-facing* - that is, we generate suggestions that a physician would be required to approve of and then communicate to a patient. Our care plans contain five sections: Medications, Referrals, Tests, Lifestyle, and Supportive Care.

**DERA setup** As in the medical conversation summarization task, the goal of DERA is to improve the quality of the generated care plan by suggesting more appropriate home care for the patient, recommending additional lab tests, or otherwise better aligning the generated summary. The DERA setup is the same as the medical conversation summarization task with care plan-specific prompts.

The *Decider* starts with an initial care plan. The *Researcher* is prompted (Prompt 10) to converse with the *Decider* (Prompt 9). Finally, the *Decider* generates the final care plan (Prompt 11). by combining the initial care plan with the content of the 'scratchpad' accumulated during the conversation.

We run DERA on the care plan generation task using GPT-4 with the settings mentioned in Table 5. We used the same set of 50 medical encounters we used for the human expert evaluation of the medical conversation summarization task.

**Human Experts Evaluation** We evaluated the effectiveness of DERA to generate care plans through human evaluation with four licensed physicians. We explicitly instructed the physician evaluators that the generated plan is defined as "meant to be provider-facing, meaning that not all suggested interventions will necessarily be recommended to the patient or followed by the patient." The physicians who evaluated the quality of these care plans were not those who provided care to the patients in the original encounter.

The experts were provided with the encounter and the two careplans – the baseline GPT-4 generated summary and the DERA generated summary starting from GPT-4 generated summary. They were asked to answer three questions similar to those described in section 3. For brevity, these are included in Appendix A.2.

Figure 3 shows the results. In a head-to-head comparison, the physicians prefer the final care plan produced by DERA 84% of the time. Furthermore, when asked to give what fraction of care plan corrections were useful, they fully agreed with 72% of suggestions. They agree with none of the suggestions only 14% of the time. Finally, they rated 92% of care plans as complete, compared to 64% of initial care plans. In summation, the application of DERA to care plan generation increased the resulting quality substantially.

In addition to these questions, we also asked the physician-experts the following: *If this care plan were acted upon by the patient, does this care plan contain information that could potentially be harmful to the patient given their presentation?* (Options: Yes, No). The amount of careplan containing "harmful" information drops from 2% in the initial careplan to 0% in the final DERA summary. As stated in section 3, we caution against drawing generalizations from these harmfulness numbers, especially in sub-topics beyond tele-medicine.

**Qualitative Examples** We show a qualitative example of the care plan generation task with DERA in Appendix Figure 4. The initial care plan generated by the *Decider* was originally rated as containing "Most" necessary care management steps by our physician-expert evaluator, suggesting there were still some improvements possible. In the DERA dialog, the *Researcher* highlights potential drug interactions with the patient's current medications and the recommendation to educate the patient on safe sexual practices. These corrections were accepted by the *Decider*, as evidenced by the notes written to the scratchpad. In turn, the corrections were manifested in the final care plan, with the three changes **bolded**. This final care plan was rated as containing "All" necessary care management steps by our physician-expert evaluator.

## 5 Open-Ended Medical Question Answering

We also investigate the use of DERA for short-form medical reasoning. A commonly used dataset for this task is MedQA (Jin et al., 2021) which consists of USMLE-style practice multiple-choice questions. Previous approaches for this dataset have included using RoBERTa (Liu et al., 2019), refining chain-of-thought using GPT-3 (Liévin et al., 2022), and fine-tuning PaLM (Chowdhery et al., 2022; Singhal et al., 2022). While most previously-reported results achieved passing results, recent GPT-4 is shown to work at a near-expert level (Nori et al., 2023).

In all previous work for this dataset, the primary focus was on the multiple-choice question format which has limited applicability in the real world. If these models are to support doctors in decision-making, they need to operate without any options provided. To mimic this setting, we extend the MedQA dataset to be open-ended to evaluate the model in a more realistic and harder setting. In an open-ended form, the model must generate the correct answer free-form and not choose from a given bank of options. We also evaluate a set of continuing education questions from the New England Journal of Medicine (NEJM), again in an open-ended setting.

A method that can perform at a high level on this task requires several attributes. First, it must be able to recall a large set of knowledge across multiple domains of medicine. Second, it must be able to reason over long questions, which will likely

| Model | Accuracy |
|---|---|
| PaLM (Singhal et al., 2022) | 0.676 |
| Nori et al. (2023) | 0.814 |
| GPT-4 0-shot | 0.834 |
| DERA | 0.840 |

Table 2: MedQA multiple-choice (4-option)

include both irrelevant and crucial facts needed to arrive at the solution.

**Experimental Setup** For our DERA setup, we include multiple prompts for the *Decider* agent, including one that generates a distribution of answers based on a self-consistency approach, one that discusses the question with the *Researcher*, and one that answers the question given the question and chat. We also formulate the *Researcher* agent with a single prompt. We report results on two question-answering datasets that were rewritten as open-ended questions using GPT-4. We include further details in Appendix Section A.4 about both the DERA setup and datasets.

To measure the relatedness between generated answers and the gold standard answer, we use a GPT-4 prompt (Prompt 19). Similarly, we use a separate prompt to make a binary exact match determination (Prompt 20). Finally, we evaluate the generated and gold answer similarity using BERTScore (Zhang et al. (2019), model scibert-basevocab-uncased), although this approach has limitations (Hanna and Bojar, 2021; Sun et al., 2022).

**Results** We compare DERA to single-shot performance using GPT-4, where $n = 5$ answers are detected, and the one with the most votes is selected as the answer. Due to the costs involved with running the experiments, we only report single runs. We include quantitative results for open-ended question answering in Table 3, and for multiple-choice question answering in Table 2.

For the multiple-choice results, we find that GPT-4 outperforms the best previously published approaches out of the box on MedQA. This is in line with that reported by Nori et al. (2023), which uses a very similar approach. We suspect that our results are slightly higher due to our use of a self-consistency approach. We do not see significant improvements when applying DERA compared to the multiple choice setting. We include further analysis in Appendix A.4.2.

| | MEDQA | | | NEJM | | |
|---|---|---|---|---|---|---|
| | BERTScore | GPT-4 Exact | GPT-4 Sim | BERTScore | GPT-4 Exact | GPT-4 Sim |
| GPT-4 1-shot | 0.746 | 0.698 | 0.65 | 0.676 | 0.703 | 0.711 |
| DERA | 0.744 | 0.703 | 0.67 | 0.670 | 0.711 | 0.724 |

Table 3: MedQA and NEJM Open-Ended. We evaluate the quality of the generated answers by using GPT-4 prompts that identify exact and similar matches (using a 0-1 scale) and average BERTScore $F_1$.

In the open-ended setting, we see strong performance in both one-shot GPT-4 and DERA for both NEJM and MedQA. Liévin et al. (2022) notes that the passing grade for the MedQA test set is 60%. For both GPT-4 one-shot and DERA, we see that GPT-4 Exact Matching is above 60% and BERTScore and Similarity measures are above 0.6. This marks an impressive ability to generate open-ended answers to questions. Yet there still exists a gap between open-ended and multiple-choice performance, suggesting opportunities for future work.

Similarly to the multiple choice setting, DERA shows small to no improvement over GPT-4, depending on the metric. The largest gain for DERA is in the similarity metric for both MedQA and NEJM, which suggests that DERA can lead to answers that are closer to the ground truth. Examples of the open-ended question-answering chats are included in Appendix Section A.5.

We include a qualitative evaluation in Appendix Section A.4.1. We note that DERA changes the answer in a majority of cases, although sometimes it maintains the same answer. More enlightening is the fact that DERA often adds additional details to the answer (e.g. responding with two tests instead of one) that further removes it from the more general gold answer. This illustrates the difficulty of determining the correct level of specificity for open-ended question-answering scoring.

## 6 Discussion and Conclusion

We introduce DERA, a framework improving LLM generations. This approach reduces the need for an LLM to produce a high-fidelity generation in one or two passes. We find that using LLM-powered agent dialog is an effective forum to improve output. We use two types of agents – *Researcher*, tasked with reviewing and selecting information, and *Decider*, tasked with integrating that information into the final output. However, we propose that this approach can generalize to other agent setups as dictated by the task.

As found in Sections 3 and 4, we find DERA im-

proves the quality of the generated text in a variety of metrics. Importantly, this reduces the number of hallucinations and omissions in the resulting text. This finding is important given the ability of LLMs to generate text that is fluent but potentially prone to errors, especially with GPT-4. The ability of DERA to identify and correct these hallucinations and omissions is critical when applying these models to real-world scenarios. A key feature is that the same LLM can be harnessed in both roles.

We did not find similar improvements in the question-answering task. As discussed in Section 5, DERA produced little to no improvement over a GPT-4 baseline. We suggest this is due to several factors, including the requirement to generate a single, granular answer. DERA often adds information to an answer, which is not helpful for short text generation. These findings, paired with those discussed above, suggest this method is well-suited for longer-generation tasks.

The chat-based format of DERA allows for increased interpretability when auditing the results. Even though LLMs such as GPT-4 may achieve high performance in zero-shot or one-shot settings, generating long-form explanations does not provide a granular forum for understanding resulting generations. Conversely, the chat-based format allows for discussions that are granular and could be verified by an end user for mistakes. We believe these insights are applicable to other domains and tasks given the plug-and-play nature of DERA.

The DERA setup could be altered to include human input in the discussion. Alternatively, different problems may dictate the inclusion of different types of agents. Overall, we believe that while LLM-based tools are critical in increasing the quality of natural language performance, research is required to ensure they are consistent and auditable. Finally, we reiterate the need for further research in automated metrics for evaluating LLM output. Human-led qualitative evaluations can provide important insights, but it remains a challenge to measure improvement given the limited tools currently

available.

# 7 Limitations

The experiments in this paper were performed using OpenAI's API, mostly using GPT-4 models. While these models generate text at a higher quality than other previous models, there are still limitations. First, we do not have access to what the model has and has not been trained on. Specifically, we do not know if openly-released datasets, such as MedQA, were included in the training data. Second, we report results using the latest version of GPT-4 available at the time. As OpenAI does not persist models, this may make reproducing results challenging.

While we include a variety of quantitative evaluations, the task of automatically evaluating generated text needs further research. We highlight the need for the broader community to build robust, generalizable metrics, and not limited to a single LLM. Similarly, while we find that the ability of DERA to reduce the presence of harmful text is promising, we encourage future users to conduct their own harmfulness study.

Broadly, dataset construction and usage are challenging in the clinical space. Several evaluation datasets cannot be openly released for data privacy or licensing reasons. Additionally, some other openly available datasets cannot be directly used with API-based models (Agrawal et al., 2022), further limiting options. We also acknowledge that while MedQA does probe medical knowledge, it likely does so in a different form than is likely to be applied in a regular clinical setting.

# 8 Ethical Considerations

The datasets used for the Summarization and Care Plan tasks contain Patient Health Information (PHI). Research on this dataset was conducted as a quality improvement activity as defined in the United States of America 45CRF §46.104(d)(4)(iii). This data will not be shared publicly due to patient privacy and HIPAA compliance. All data is de-identified and scrubbed for protected health information prior to experimentation.

# References

Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. Large language models are zero-shot clinical information extractors. *arXiv preprint arXiv:2205.12689*.

Sergey Berezin and Tatiana Batura. 2022. Named entity inclusion in abstractive text summarization. In *Proceedings of the Third Workshop on Scholarly Document Processing*, pages 158–162, Gyeongju, Republic of Korea. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Bharath Chintagunta, Namit Katariya, Xavier Amatriain, and Anitha Kannan. 2021. Medically aware GPT-3 as a data generator for medical dialogue summarization. In *Proceedings of the Second Workshop on Natural Language Processing for Medical Conversations*, pages 66–76, Online. Association for Computational Linguistics.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.

Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. 2023. Palm-e: An embodied multimodal language model. *Preprint*, arXiv:2303.03378.

Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and Siva Reddy. 2022. On the origin of hallucinations in conversational models: Is it the datasets or the models? In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5271–5285, Seattle, United States. Association for Computational Linguistics.

Seppo Enarvi, Marilisa Amoia, Miguel Del-Agua Teba, Brian Delaney, Frank Diehl, Stefan Hahn, Kristina Harris, Liam McGrath, Yue Pan, Joel Pinto, Luca Rubini, Miguel Ruiz, Gagandeep Singh, Fabian Stemmer, Weiyi Sun, Paul Vozila, Thomas Lin, and Ranjani Ramamurthy. 2020. Generating medical reports from patient-doctor conversations using sequence-to-sequence models. In *Proceedings of the First Workshop on Natural Language Processing for Medical*

*Conversations*, pages 22–30, Online. Association for Computational Linguistics.

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *arXiv preprint*.

Michael Hanna and Ondřej Bojar. 2021. A fine-grained analysis of BERTScore. In *Proceedings of the Sixth Conference on Machine Translation*, pages 507–517, Online. Association for Computational Linguistics.

Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *Preprint*, arXiv:2210.11610.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

Anirudh Joshi, Namit Kataria, Xavier Amatriain, and Anitha Kannan. 2020. Dr. summarize: Global summarization of medical dialogue by exploiting local structures. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3755–3763, Online. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Valentin Liévin, Christoffer Egeberg Hother, and Ole Winther. 2022. Can large language models reason about medical questions? *arXiv preprint arXiv:2207.08143*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.

Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint*.

Varun Nair, Elliot Schumacher, and Anitha Kannan. 2023. Generating medically-accurate summaries of patient-provider dialogue: A multi-stage approach using large language models. *Preprint*, arXiv:2305.05982.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of GPT-4 on Medical Challenge Problems.

OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 248–260. PMLR.

Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–7.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2022. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*.

Tianxiang Sun, Junliang He, Xipeng Qiu, and Xuanjing Huang. 2022. BERTScore is unfair: On social bias in language model-based metrics for text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3726–3739, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Oyvind Tafjord, Bhavana Dalvi Mishra, and Peter Clark. 2022. Entailer: Answering questions with faithful and truthful chains of reasoning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2078–2093, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Boshi Wang, Xiang Deng, and Huan Sun. 2022. Iteratively prompt pre-trained language models for chain of thought. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2714–2730, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.

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.

Longxiang Zhang, Renato Negrinho, Arindam Ghosh, Vasudevan Jagannathan, Hamid Reza Hassanzadeh, Thomas Schaaf, and Matthew R Gormley. 2021. Leveraging pretrained models for automatic summarization of doctor-patient conversations. *arXiv preprint arXiv:2109.12174*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Deyao Zhu, Jun Chen, Kilichbek Haydarov, Xiaoqian Shen, Wenxuan Zhang, and Mohamed Elhoseiny. 2023. Chatgpt asks, blip-2 answers: Automatic questioning towards enriched visual descriptions. *Preprint*, arXiv:2303.06594.

## A Appendix

### A.1 Long-Form Text Generation Metrics

We measure the degree to which corruptions are present by using a GPT-based metric that tracks the medical concept coverage of the medical summary, **GPT-F1**. To compute GPT-F1, we compute the harmonic mean of two sub-metrics: GPT-Recall and GPT-Precision. We describe each sub-metric below.

**GPT-Recall**: To compute, we first extract medical entities from both the predicted text and ground-truth text[2] of the same summary section (using Prompt 6) and use a verification prompt (Prompt 7) to infer if the entities extracted from the ground-truth section are also present in the predicted text, This produces $tp_{gt}$ and $f_n$ values, which is used to calculate GPT-Recall $= \frac{tp_{gt}}{tp_{gt}+f_n}$.

**GPT-Precision**: To compute, we also first extract medical entities from the corresponding predicted and ground-truth summary sections and verify concepts extracted from the predicted section are also present in the ground-truth text, either as exact matches or re-phrasings. This produces $tp_{pred}$ and $f_p$, which is used to calculate GPT-Precision $= \frac{tp_{pred}}{tp_{pred}+f_p}$.

We present the results of our quantitative evaluation using the GPT-F1 metric in Table 1. Specifically, we compare GPT-F1 on the initial summary with errors to the DERA corrected summary. Note first how the higher levels of corruption manifest in the initial summary GPT-F1. As the corruption level of the initial summary increases, the initial GPT-F1 score drops. We find that DERA can produce significantly improved summaries in low, medium, and high levels of corruption, as evidenced by increases in GPT-F1. This suggests that the interaction between the *Researcher* and *Decider* agents is identifying hallucinations and omissions and resolving them through dialog, even when many such corruptions are present.

### A.2 Qualitative Evaluation of Care Plan

Physicians were asked the following three questions in their evaluations of care plans:

1. Which careplan do you prefer to use for the given patient and encounter? (Options: Initial, Final).

2. What fraction of the necessary care plan management steps are captured? (Options: All, Most, Some, None)

3. What percentage of the suggestions added to the DERA scratchpad do you agree with? (Options: All, Most, Some, None)

### A.3 Care Plan Example

---

[2]Note that the terms "predicted" and "ground-truth" are used loosely here, and can also refer to synthetic text that is used as ground-truth (as in the case of the corruption experiments).

**Doctor-Patient Chat**

Patient: UTI

Doctor: Hi NAME, thank you for starting a visit. My name is NAME, and I am your clinical associate for this live chat. . . . I am sorry to hear that.

Patient: Okay thank you . . .

Doctor: When was your last sexual encounter?

Patient: Two days ago. I'm not having any since then because of the pain with my bladder infection.

Doctor: Was the encounter unprotected?

Patient: No

Doctor: How many sexual partners have you had in the past 6 months?

Patient: Several. I did however get tested for HIV and STIs a couple weeks ago and all came back negative. I get tested regularly since I am sexually active.

Doctor: Do you have a new sexual partner?

Patient: No new ones

Doctor: Do you douche or use feminine hygiene products?

Patient: No Sometimes I use baby wipes
. . .

Doctor: For what medical conditions you are taking BuSpar, LaMICtal, Zoloft?

Patient: Buspar for Bipolar, Lamictal for anxiety, Zoloft for depression

---

**Initial Care Plan**

*Tests*: (1) Perform a urinalysis and urine culture to confirm the presence of a urinary tract infection and to identify . . .

*Medications*: (1) Prescribe an appropriate antibiotic, such as nitrofurantoin, fosfomycin, or trimethoprim/sulfamethoxazole, to treat the infection. Note that the choice of antibiotic may be adjusted based on the results of the urine culture. (2) Continue the use of over-the-counter pain relief medication, such as AZO, until the antibiotic treatment relieves symptoms.

*Lifestyle*: (1) Encourage the patient to increase fluid intake, particularly water, to help flush out bacteria from the urinary tract. (2) Suggest urinating frequently and fully emptying the bladder to help clear the infection. (3) Recommend proper hygiene practices, such as . . . (4) Advise the use of a urinary health supplement . . .

*Referrals*: (1) If symptoms worsen or do not improve after a few days of antibiotic treatment, refer the patient to . . .

*Supportive Care*: (1) Provide education on recognizing the signs and symptoms of recurrent urinary tract infections and the importance of seeking timely medical care. (2) Offer reassurance and support regarding the patient's mental health and . . .

---

**DERA Dialog**

*Researcher*: Starting with the Tests section, everything looks good with the urinalysis and urine culture being recommended. However, in the Medications section, it would be helpful to include that antibiotic choice should consider the patient's previous antibiotic exposure and allergies. It will also be useful to mention any potential drug interactions with BuSpar, LaMICtal, and Zoloft. In the Lifestyle section, there is no mention of the patient's sexual practices, which could contribute to recurrent UTIs. Considering the patient's history of multiple sexual partners, you might want to add a recommendation to educate the patient on safe sexual practices and the importance of using a barrier method during sex to reduce the risk of UTIs. . . .

*Decider*: Thank you for your suggestions. I agree with the recommendation to consider the patient's previous antibiotic exposure and allergies, as well as potential drug interactions with their other medications. I also see the importance of addressing safe sexual practices and barrier method usage to reduce the risk of UTIs. In terms of supportive care, providing guidance on managing UTIs and potential medication interactions concerning their mental health medications is crucial.

*Researcher*: Thank you for considering my suggestions. Please make sure to discuss these points with the patient, and . . .

---

**Scratchpad**

1) Add a note in the Medications section to consider the patient's previous antibiotic exposure, allergies, and potential drug interactions with BuSpar, LaMICtal, and Zoloft.

2) In the Lifestyle section, include a recommendation to educate the patient on safe sexual practices and the importance of using barrier methods during sex to reduce the risk of UTIs.

3) In the Supportive Care section, offer guidance on managing UTIs and potential medication interactions related to their mental health medications.

> **Final Care Plan**
>
> *Tests*: (1) Perform a urinalysis and urine culture to confirm the presence of a urinary tract infection and to identify …
> *Medications*: (1) Prescribe an appropriate antibiotic, such as nitrofurantoin, fosfomycin, or trimethoprim/sulfamethoxazole, to treat the infection. Note that the choice of antibiotic may be adjusted based on the results of the urine culture. **Consider the patient's previous antibiotic exposure, allergies, and potential drug interactions with BuSpar, LaMICtal, and Zoloft.** (2) Continue the use of over-the-counter pain relief medication, such as AZO, until the antibiotic treatment relieves symptoms.
> *Lifestyle*: (1) Encourage the patient to increase fluid intake, particularly water, to help flush out bacteria from the urinary tract. (2) Suggest urinating frequently and fully emptying the bladder to help clear the infection. (3) Recommend proper hygiene practices, such as … (4) Advise the use of a urinary health supplement … **(5) Educate the patient on safe sexual practices and the importance of using barrier methods during sex to reduce the risk of UTIs.**
> *Referrals*: (1) If symptoms worsen or do not improve after a few days of antibiotic treatment, refer the patient to …
> *Supportive Care*: (1) Provide education on recognizing the signs and symptoms of recurrent urinary tract infections and the importance of seeking timely medical care. (2) Offer reassurance and support regarding the patient's mental health … **(3) Offer guidance on managing UTIs and potential medication interactions related to their mental health medications.**

Figure 4: Qualitative example of care plan generation with DERA abridged for space (continued from previous page). Starting with the initial chat, the *Decider* generates an initial care plan. The *Researcher* and *Decider* agents in DERA then converse with one another, visible in DERA dialog. The *Decider* adds accepted suggestions to a scratchpad, which collects the final changes to make to the care plan. The final care plan is generated by the *Decider* using this scratchpad. Note the points in **bold** that were added to the final care plan.

## A.4 Question Answering Experimental Details

**DERA setup**  To generate an initial answer for DERA to discuss, we use a single-shot prompt which outputs a short answer (Prompt 14). We use a single-shot prompt to ensure a consistent output, which we were unable to achieve with a zero-shot prompt. Earlier work (Singhal et al., 2022) has shown that using a self-consistency strategy provides stronger results. We adopt this approach by running 5 completions of our single-shot prompt and selecting the answer with the most votes as the *single-shot* answer, and consider this as our baseline[3].

Instead of initializing our *Decider* with a single answer, we provide it with the distribution of votes. This approach provides DERA with the distribution better captures the underlying uncertainty of the model[4]. A prompt (Prompt 15) is tasked with writing a reasoning behind the votes, which is used as the initial *Decider* message.

Starting with the initial *Decider* message, both *Decider* (Prompt 17) and *Researcher* have access only to the question and their own conversation as they iteratively discuss the problem and attempt to achieve the right answer. The *Researcher* can stop the dialogue when they have exhausted all relevant information, otherwise, it is set to end after $n = 3$ turns. At each turn, the *Decider* must state what

their current answer is and explain their reasoning, and they may choose to either confirm or change their answer.

We instruct both prompts to act as medical doctors who work at an expert level. To arrive at a final answer, a prompt is given the same information as the original one-shot prompt, with the exception that it is also given the full chat history to use as additional context. We generate $n = 5$ answers and use the most frequently generated answer as our final answer (see Prompt 18). If there are ties, the first completion of the highest-ranking answers is selected.

We run DERA on open-ended question answering with the parameters noted in Table 6. For the multiple-choice setting, we use a very similar configuration. The primary prompt changes are limited to the fact that *Decider* is given a set of options and asked to generate the letter (A-D) instead of a short phrase.

**Datasets**  We evaluate our approach using two Medical Question answering datasets - MedQA US dataset (Jin et al., 2021) and New England Journal of Medicine Test Questions (NEJM). Both datasets consist of questions taken from practice or real medical exams (United States Medical Licensing for MedQA, and continuing education questions for NEJM). For both datasets, the questions are originally written in multiple-choice format (*e.g, Which of the following is the best diagnosis?*). Our goal is to test DERA 's performance on open-ended question answering, where the task will be to generate the answer free-form.

---

[3]We do not account for variations in the text, each lexical form is counted separately.

[4]This also handles cases where closely related lexical forms receive separate votes, as the *Decider* output will conclude that the options are similar.

Therefore, we use GPT-4 to alter the questions to be open-ended. In most cases, this requires a simple rephrasing of the final sentence. For example, the previous question could be re-written as *What is the best diagnosis?*. In these cases, we restrict GPT-4 to rewrite only the final sentence of the question, so as to guard against hallucinations. When a more complex rewrite is required, we prompt GPT-4 to rewrite the entire question and find that it only changes the relevant sentence. Some questions could already be answered open-ended and required no rewriting. Although we performed quality checks, as the entire process is automated, there may be some errors. The prompts for rewriting the final sentence 13 and the full question 12 are included in the Appendix. We also release the full MedQA open-ended dataset at [REDACTED] We cannot release the NEJM dataset due to licensing issues.

For MedQA, we sample a portion of the training set (1178 questions) as a development set and maintain the integrity of the test set (1273 questions) as formulated by the authors. For NEJM, we split the datasets by area, reserving 7 areas [5] as a development set (consisting of 639 questions), with the remainder serving as a test set (1112 questions). We do not exclude questions containing images. The GPT-4 and DERA results multiple-choice results in Table 2 used the model available in Feb. 2023.

### A.4.1 Open-Ended Analysis

We include the first 10 examples from the MedQA development set (which we randomly drew from their training set) in Appendix Table 4[6]. In our analysis of these development examples, we see several patterns.

First, sometimes the agent successfully changes an incorrect answer to the correct answer. For example, in Question 4 shown in Appendix Section A.5, the original answer is *Inherited bleeding disorder*, and DERA changes it to the more specific *Von Willebrand Disease*. In other cases, DERA leaves the answer as the same in the original 1-shot generation (*e.g,* Questions 5, 9, 55, 94, 98). We also note that this does not occur in a majority of cases, as only 542 of the 1273 MedQA training examples

---

[5]Reproductive, Gastrointestinal, Neurologic/Psychogenic, Special Sensory, Endocrine, Musculoskeletal, and Maternity Care

[6]These results were generated with an earlier version of GPT-4 available in February 2023.

have the exact same answer between DERA and one-shot.

In other cases, such as in Question 54, DERA adds additional details to the 1-shot answer (1-shot *Smoking cessation counseling and support* to the *Decider*'s final answer *Assessing for occupational lung disease and providing smoking cessation*. There are some clear challenges with open-ended question answering that show in both the DERA and 1-shot generations. Specifically, often both give a more general answer than is given in the gold standard answer. For example, in Question 74, the gold standard answer text is a specific medication (*Deantrolene*), while both DERA and 1-shot produce more general answers (*e.g, Immediate hospitalization and supportive care*).

Overall, without the inclusion of a specific set of options, it is difficult for GPT-4 to generate an answer at a correct level of specificity (*e.g,* a specific treatment instead of a general approach) and a correct length (*e.g,* answering in short answer format instead of long sentences). In some settings, these attributes may be useful, but it results in a challenging approach to evaluate. We predict the need for additional work in methods that automatically evaluate the output of large language model-powered tools, given the inherent complexity present.

### A.4.2 Multiple Choice Analysis

The comparison to the open-ended question answering task, the multiple choice task is even more discrete. Therefore, is usually no room for a nuanced answer change which is DERA's main performance improvement. For example, DERA's final answer is only different 8.9% (113) of the time compared to the voting-based initial answer. Of those questions that do differ between DERA and the voting-based initial answer, DERA changes 47% (53) of questions from an incorrect initial answer to the correct answer, while it changes 35% (39) of questions from the correct answer to the incorrect answer. The remainder of the time, both DERA and the voting-based answer are incorrect. In addition, the baseline performance of the voting method is high, and therefore there is less room for improvement. We believe the remaining errors would benefit from approaches such as retrieval augmented generation to provide correct context to the model.

## A.5 Open-Ended Question Answering Examples

We provide several examples of chats between the *Decider* and *Researcher* agents taken from our MedQA development set. We include the Gold Answer, 1-shot answer, and DERA answer. In the following example, the single-shot answer is not specific enough and the DERA dialog changes the answer to the correct one.

---

***Question*** [ID 4]:A 20-year-old woman presents with menorrhagia for the past several years. She says that her menses "have always been heavy", and she has experienced easy bruising for as long as she can remember. Family history is significant for her mother, who had similar problems with bruising easily. The patient's vital signs include: heart rate 98/min, respiratory rate 14/min, temperature 36.1°C (96.9°F), and blood pressure 110/87 mm Hg. Physical examination is unremarkable. Laboratory tests show the following: platelet count 200,000/mm3, PT 12 seconds, and PTT 43 seconds. What is the most likely cause of this patient's symptoms? *Gold Answer: Von Willebrand Disease, 1-shot : Inherited Bleeding Disorder*

- *Decider*: The patient's symptoms of menorrhagia and easy bruising, along with a family history of similar problems, suggest an inherited bleeding disorder. However, Von Willebrand disease, which is the most common inherited bleeding disorder, cannot be ruled out. Both options are possible causes of the patient's symptoms, but the presence of a family history of bleeding problems slightly increases the likelihood of an inherited bleeding disorder.
- *Researcher*: Remember to consider the fact that the patient has experienced menorrhagia for several years and easy bruising for as long as she can remember, along with a family history of similar problems
- *Decider*: After considering the patient's symptoms of menorrhagia and easy bruising, along with a family history of similar problems, I believe that Von Willebrand disease is the most likely cause of this patient's symptoms.
- *Researcher*: Take note of the laboratory test results, specifically the platelet count, PT, and PTT values, as these can help determine the cause of the patient's symptoms.
- *Decider*: After taking into consideration the laboratory test results, specifically the platelet count, PT, and PTT values, I still believe that Von Willebrand disease is the most likely cause of this patient's symptoms. The patient's symptoms of menorrhagia and easy bruising, along with a family history of similar problems, further support this diagnosis.
- *Researcher*: Pay attention to the patient's vital signs, specifically the heart rate, respiratory rate, temperature, and blood pressure, as these can also provide clues to the cause of her symptoms.
- *Decider*: After taking into consideration the patient's vital signs, specifically the heart rate, respiratory rate, temperature, and blood pressure, I still believe that Von Willebrand disease is the most likely cause of this patient's symptoms. The patient's symptoms of menorrhagia and easy bruising, along with a family history of similar problems and laboratory test results further support this diagnosis. *DERA answer: Von Willebrand Disease*

---

In the following example, the 1-shot prompt arrives at the correct answer, but DERA expands the answer to include additional information.

---

***Question*** [ID 574]:Your colleague has been reading the literature on beta-carotene supplementation and the risk of heart disease. She thinks they may share a clinically relevant association and would like to submit an editorial to a top journal. Upon final literature review, she discovers a newly published study that refutes any association between beta-carotene and heart disease. Your colleague is upset; you suggest that she, instead, mathematically pool the results from all of the studies on this topic and publish the findings. What type of study design are you recommending to your colleague? *Gold Answer: Meta-analysis, 1-shot : Meta-analysis*

- *Decider*: Both options are viable for pooling results from multiple studies, but a meta-analysis (80%) is more likely to be recommended as it allows for a quantitative synthesis of the data. A systematic review and meta-analysis (20%) would also be a good

option, as it includes a comprehensive literature search and qualitative analysis, but it may be more time-consuming. Ultimately, the choice depends on the specific goals and resources of the researcher.

*Researcher*: Consider the fact that your colleague wants to 'mathematically pool the results from all of the studies on this topic' and think about which study design would best allow for this type of analysis.

*Decider*: Based on the fact that my colleague wants to mathematically pool the results from all of the studies, I believe a meta-analysis would be the best study design as it allows for a quantitative synthesis of the data.

*Researcher*: Think about the goal of your colleague, which is to 'submit an editorial to a top journal' and consider which study design would be most appropriate for achieving this goal.

*Decider*: After considering the goal of my colleague to submit an editorial to a top journal, I believe a systematic review and meta-analysis would be the most appropriate study design. This design includes both a comprehensive literature search and qualitative analysis, in addition to the quantitative synthesis of data, which would make it more suitable for publication in a top journal.

*Researcher*: Consider the fact that your colleague discovered a 'newly published study that refutes any association between beta-carotene and heart disease' and think about how this might impact the choice of study design.

*Decider*: After considering the fact that my colleague discovered a newly published study that refutes any association between beta-carotene and heart disease, I still believe a systematic review and meta-analysis would be the most appropriate study design. This design would allow for the inclusion of this new study in the comprehensive literature search and analysis, providing a more complete picture of the current state of research on this topic. *DERA answer: Systematic review and meta-analysis*

In the following example, the 1-shot prompt arrives at the wrong answer, and while DERA changes the answer, it is also not correct.

---

*Question* [ID 2599]:A 15-year-old girl is brought to her pediatrician's office complaining of frequent diarrhea, fatigue, and inability to gain weight. Her vital signs are within normal limits, and her BMI is 17. She describes her stools as pale, malodorous, and bulky. She often has abdominal bloating. Her symptoms are most prominent after breakfast when she typically consumes cereal. After several weeks of careful evaluation and symptomatic treatment, the pediatrician recommends an esophagogastroduodenoscopy. A diagnostic biopsy shows blunting of intestinal villi and flat mucosa with multiple intraepithelial lymphocytes. What is the patient likely deficient in? *Gold Answer: IgA, 1-shot : Gluten tolerance*

*Decider*: Based on the patient's symptoms and the results of the esophagogastroduodenoscopy, it is likely that the patient has a gluten intolerance or Celiac disease. However, without further testing and confirmation, it is difficult to determine the exact diagnosis with certainty. Therefore, the options of Gluten intolerance/Celiac Disease, Gluten intolerance or celiac disease, and Gluten intolerance or Celiac disease are all rated at 20%. The option of Gluten tolerance is rated at 40% as it is possible that the patient's symptoms are caused by another condition.

*Researcher*: Consider the patient's symptoms of frequent diarrhea, fatigue, and inability to gain weight, along with her description of pale, malodorous, and bulky stools.

*Decider*: After considering the patient's symptoms of frequent diarrhea, fatigue, and inability to gain weight, along with her description of pale, malodorous, and bulky stools, I believe that Celiac Disease is the most likely diagnosis. The results of the esophagogastroduodenoscopy showing blunting of intestinal villi and flat mucosa with multiple intraepithelial lymphocytes also support this conclusion.

*Researcher*: Think about the patient's symptoms being most prominent after breakfast when she typically consumes cereal, and how this could be related to her condition.

*Decider*: After considering the patient's symptoms being most prominent after breakfast

when she typically consumes cereal, I believe this further supports my conclusion that Celiac Disease is the most likely diagnosis. The consumption of cereal, which often contains gluten, could be triggering her symptoms.

*Researcher*: Consider the diagnostic biopsy results which show blunting of intestinal villi and flat mucosa with multiple intraepithelial lymphocytes, and how this could be related to the patient's condition.

*Decider*: After considering the patient's symptoms, the timing of her symptoms being most prominent after consuming cereal, and the diagnostic biopsy results showing blunting of intestinal villi and flat mucosa with multiple intraepithelial lymphocytes, I believe Celiac Disease is the most likely diagnosis.

*DERA answer: Celiac Disease*

Table 4: Examples from the MedQA Training set (used as a development set in our experiments). We include the id (or line number), the rewritten open-ended question, the correct original gold answer, the student predicted answer, and the 1-shot answer.

| id | question | gold text | DERA answer | 1-shot answer |
|---|---|---|---|---|
| 4 | A 20-year-old woman presents with menorrhagia for the past several years. She says that her menses "have always been heavy", and she has experienced easy bruising for as long as she can remember. Family history is significant for her mother, who had similar problems with bruising easily. The patient's vital signs include: heart rate 98/min, respiratory rate 14/min, temperature 36.1°C (96.9°F), and blood pressure 110/87 mm Hg. Physical examination is unremarkable. Laboratory tests show the following: platelet count 200,000/mm3, PT 12 seconds, and PTT 43 seconds. What is the most likely cause of this patient's symptoms? | Von Willebrand disease | Von Willebrand disease | Inherited bleeding disorder |
| 5 | A 40-year-old zookeeper presents to the emergency department complaining of severe abdominal pain that radiates to her back, and nausea. The pain started 2 days ago and slowly increased until she could not tolerate it any longer. Past medical history is significant for hypertension and hypothyroidism. Additionally, she reports that she was recently stung by one of the zoo's smaller scorpions, but did not seek medical treatment. She takes aspirin, levothyroxine, oral contraceptive pills, and a multivitamin daily. Family history is non-contributory. Today, her blood pressure is 108/58 mm Hg, heart rate is 99/min, respiratory rate is 21/min, and temperature is 37.0°C (98.6°F). On physical exam, she is a well-developed, obese female that looks unwell. Her heart has a regular rate and rhythm. Radial pulses are weak but symmetric. Her lungs are clear to auscultation bilaterally. Her lateral left ankle is swollen, erythematous, and painful to palpate. An abdominal CT is consistent with acut... | Scorpion sting | Scorpion venom-induced acute pancreatitis | Scorpion venom-induced acute pancreatitis |

Continued on next page

| id | question | answer text | DERA answer | 1-shot answer |
|---|---|---|---|---|
| 9 | A 35-year-old male presents to his primary care physician with complaints of seasonal allergies. He has been using intranasal vasoconstrictors several times per day for several weeks. What is a likely consequence of the chronic use of topical nasal decongestants? | Persistent congestion | Rhinitis Medicamentosa (rebound nasal congestion) | Rhinitis medicamentosa (rebound nasal congestion) |
| 54 | A 60-year-old man comes to the physician for an examination prior to a scheduled cholecystectomy. He has hypertension treated with hydrochlorothiazide. His mother had chronic granulomatous disease of the lung. He works in a glass manufacturing plant. He has smoked two packs of cigarettes daily for 38 years. His vital signs are within normal limits. Examination shows no abnormalities. Laboratory studies are within the reference range. An x-ray of the chest is shown. What is the most appropriate next step in management? | Request previous chest x-ray | Assessing for occupational lung disease and providing smoking cessation counseling. | Smoking cessation counseling and support. |
| 55 | You are examining a 3-day-old newborn who was delivered vaginally without any complications. The newborn presents with vomiting, hyperventilation, lethargy, and seizures. Blood work demonstrates hyperammonemia, elevated glutamine levels, and decreased blood urea nitrogen. A CT scan demonstrates cerebral edema. Which enzyme defects would result in a clinical presentation similar to this infant? | Carbamoyl phosphate synthetase I | "Urea cycle enzyme deficiencies" | Urea cycle enzyme deficiencies |
| 64 | An 18-year-old man comes to the clinic with his mom for "pins and needles" of both of his arms. He denies any past medical history besides a recent anterior cruciate ligament (ACL) tear that was repaired 1 week ago. The patient reports that the paresthesias are mostly located along the posterior forearms, left more than the right. What would you expect to find on physical examination of this patient? | Loss of wrist extension | Decreased sensation and possible weakness in both posterior forearms, with the left side being more affected than the right. | Decreased sensation and possible weakness in the posterior forearms. |

Continued on next page

| id | question | answer text | DERA answer | 1-shot answer |
|---|---|---|---|---|
| 74 | A 16-year-old girl is brought to the emergency department by her friends who say that she took a whole bottle of her mom's medication. They do not know which medication it was she ingested. The patient is slipping in and out of consciousness and is unable to offer any history. Her temperature is 39.6°C (103.2°F), the heart rate is 135/min, the blood pressure is 178/98 mm Hg, and the respiratory rate is 16/min. On physical examination, there is significant muscle rigidity without tremor or clonus. What is the best course of treatment for this patient? | Dantrolene | Immediate stabilization and supportive care with emergency toxicology consultation. | Immediate hospitalization and supportive care. |
| 77 | A 3-week-old boy is brought to the emergency department by his parents because of a 3-day history of progressive lethargy and difficulty feeding. He was born at term and did not have difficulty feeding previously. His temperature is 39.4°C (103°F), pulse is 220/min, respirations are 45/min, and blood pressure is 50/30 mm Hg. Pulse oximetry on 100% oxygen shows an oxygen saturation of 97%. Examination shows dry mucous membranes, delayed capillary refill time, and cool skin with poor turgor. Despite multiple attempts by the nursing staff, they are unable to establish peripheral intravenous access. What is the most appropriate next step in management for this 3-week-old boy? | Intraosseous cannulation | Establishing intraosseous access for fluid resuscitation and medication administration. | Intraosseous needle placement for fluid resuscitation and antibiotics. |
| 94 | A 70-year-old man comes to the physician because of a 4-month history of epigastric pain, nausea, and weakness. He has smoked one pack of cigarettes daily for 50 years and drinks one alcoholic beverage daily. He appears emaciated. He is 175 cm (5 ft 9 in) tall and weighs 47 kg (103 lb); BMI is 15 kg/m2. He is diagnosed with gastric cancer. What cytokine is the most likely direct cause of this patient's examination findings? | IL-6 | Tumor necrosis factor-alpha (TNF-$\alpha$) | Tumor necrosis factor-alpha (TNF-$\alpha$) |

Continued on next page

| id | question | answer text | DERA answer | 1-shot answer |
|----|----------|-------------|-------------|---------------|
| 98 | Three days after starting a new drug for malaria prophylaxis, a 19-year-old college student comes to the physician because of dark-colored urine and fatigue. He has not had any fever, dysuria, or abdominal pain. He has no history of serious illness. Physical examination shows scleral icterus. Laboratory studies show a hemoglobin of 9.7 g/dL and serum lactate dehydrogenase of 234 U/L. Peripheral blood smear shows poikilocytes with bite-shaped irregularities. What drug has the patient most likely been taking? | Primaquine | Primaquine | Primaquine |

| Prompt | temp. | max_tokens | top_p | freq. penalty | num. turns |
|---|---|---|---|---|---|
| Summarization - Initial (1) | 1 | 512 | 1 | 0 | - |
| Summarization - Decider (3) | 1 | 512 | 1 | 0 | 15 |
| Summarization - Researcher (4) | 1 | 512 | 1 | 0 | 15 |
| Summarization - Corruption (2) | 1 | 512 | 1 | 0 | - |
| Summarization - Final (5) | 1 | 512 | 1 | 0 | - |
| GPT-F1 Metric - Concept Extractor (6) | 0 | 200 | 1 | 0 | - |
| GPT-F1 Metric - Concept Verifier (7) | 0 | 200 | 1 | 0 | - |
| Care Plan - Initial (8) | 1 | 512 | 1 | 0 | - |
| Care Plan - Decider (9) | 1 | 512 | 1 | 0 | 15 |
| Care Plan - Researcher (10) | 1 | 512 | 1 | 0 | 15 |
| Care Plan - Final (11) | 1 | 512 | 1 | 0 | - |

Table 5: Experimental details for Summarization and Care Plan Experiments. For each prompt, we include the Temperature, Maximum Number of Tokens for Generation, Top-P, the Frequency Penalty, and the Number of Turns for the prompt(s) (where applicable). In all cases, the final parameters of the model were decided by qualitative evaluation of the output.

| Prompt | temp. | max_tokens | top_p | num comp. | freq. penalty | num. turns |
|---|---|---|---|---|---|---|
| 1-shot (14) | 0.7 | 50 | 1 | 5 | 0 | 1 |
| Decider Initial Msg. (15) | 0 | 400 | 1 | 1 | 0 | 1 |
| Res.-Dec. Msg. (17 & 16) | 0.3 | 250 | 1 | 1 | 0.5 | 3 |
| Final Decider Answer (18) | 0 | 100 | 1 | 5 | 0 | 1 |
| Pred - Gold Ans. Sim. (19) | 0 | 100 | 1 | 1 | 0 | 1 |

Table 6: Experimental details for Question Answering. For each prompt, we include the Temperature, Maximum Number of Tokens for Generation, Top-P, the Number of Completions (*e.g,* the number of generations we prompted from GPT-4), the Frequency Penalty, and the Number of Turns for the prompt(s). In all cases, the final parameters of the model were decided by qualitative evaluation of the output. For GPT-3 CoT, the two turns refers to 1) the generation of the chain of thought, and 2) the generation of the answer.

**Doctor-Patient Chat**
Patient: I have a sinus infection and need something to knock it out.
Doctor: Hi NAME thank you for starting a visit. I am so sorry to hear about your sinus infection. How long have you had the symptoms for?
Patient: Since Sunday
Doctor: Ah I see. Which symptoms do you have at present?
Patient: My face is swollen my cheeks hurt my eyelids are swollen and I am running a slight fever and I can feel something draining down the back of my throat. . .

**Initial Summary, Corruptions in Red**
*Demographics and Social Determinants of Health*: The patient's specific demographics are not mentioned in the encounter. The patient experiences sinus infections multiple times a year. . .
*Patient Intent*: The patient seeks care for a sinus infection, which they have identified themselves. They are . . .
*Pertinent Positives*: The patient has been experiencing sinus infection symptoms since Sunday. Their face is swollen, cheeks hurt, eyelids are swollen, and they have a slight fever at night. The patient can also feel something draining down the back of their throat. They have a history of 3-6 sinus infections a year, with the **last one occurring a couple of months ago → last two occurring simultaneously a couple of months ago**.
*Pertinent Unknowns*: **There are no pertinent unknowns mentioned in this medical encounter. → The patient is unsure whether their current condition is a sinus infection or a viral infection.**.
*Pertinent Negatives*: **The patient has not been vaccinated against COVID-19 → The patient has received the COVID-19 vaccine.**. Their temperature is currently normal but spikes at night.
*Medical History*: **The patient has a history of seasonal allergies and sinus infections → The patient has had chickenpox, seasonal allergies, and sinus infections**. They typically use an Albuterol inhaler for asthma due to their allergies causing wheezing in the spring and fall, however they currently do not have one. They have had COVID-19 before and were previously prescribed **amoxicillin → penicillin** for sinus infections. **The patient currently uses Tylenol severe sinus medication and ibuprofen to alleviate symptoms → The patient currently uses ibuprofen and a nasal spray to alleviate symptoms.**.

**DERA dialog**
*Researcher*: Let's go through the summary section by section. First, the Pertinent Positives section - there is no mention of the patient taking Tylenol severe sinus medication, which was mentioned in the dialogue. Also, the last sentence about the last two sinus infections occurring simultaneously might be misleading as it's not explicitly mentioned in the dialogue. Could you please add the Tylenol information and consider rephrasing the last part?
*Decider*: Thank you for pointing that out. I'll add the information about Tylenol severe sinus medication to the summary. I agree that the last part about the last two sinus infections occurring simultaneously can be misleading; I will rephrase it. Do you have any other suggestions or corrections to make in the summary?
*Researcher*: Great, let's move on to the Pertinent Unknowns section. The patient actually mentioned that they know it's a sinus infection, so we shouldn't include that they are unsure about their condition. Could you please remove this part from the Pertinent Unknowns section?
*Decider*: You're right, the patient did mention that they know it's a sinus infection. I apologize for the confusion. I'll remove that part from the Pertinent Unknowns section. . . .

**Scratchpad**
Add Tylenol severe sinus medication to the summary, and rephrase the sentence about the last two sinus infections.
Remove the patient being unsure about their condition from Pertinent Unknowns section.
Correct the COVID-19 vaccine information in the Pertinent Negatives section.
Update the summary to mention the patient using the Albuterol inhaler for wheezing and change medication name to amoxicillin.

**Final Summary**
*Demographics and Social Determinants of Health*: The patient's specific demographics are not mentioned in the encounter. The patient experiences sinus infections multiple times a year. . .
*Patient Intent*: The patient seeks care for a sinus infection, which they have identified themselves. They are . . .
*Pertinent Positives*: The patient has been experiencing sinus infection symptoms since Sunday. Their face is swollen, cheeks hurt, eyelids are swollen, and they have a slight fever at night. The patient can also feel something draining down the back of their throat. They have a history of 3-6 sinus infections a year, with the last one occurring a couple of months ago and being treated with amoxicillin.
*Pertinent Unknowns*: There are no pertinent unknowns mentioned in this medical encounter.
*Pertinent Negatives*: The patient has not been vaccinated against COVID-19. Their temperature is currently normal but spikes at night.
*Medical History*: The patient has a history of seasonal allergies and sinus infections. They typically use an Albuterol inhaler for wheezing due to their allergies causing wheezing in the spring and fall, however they currently do not have one. They have had COVID-19 before. The patient currently uses Tylenol severe sinus medication and ibuprofen to alleviate symptoms.

Figure 5: Medical conversation summarization - Quantitative evaluation: The above is a condensed example taken from our summarization dataset. Starting with the initial chat, we generate an initial summary, which is then corrupted (see red text) using Prompt 2. The DERA dialog generates a scratchpad of edits, which results in the final corrected summary. The quantitative evaluation compares the initial summary with the final summary .

```
1  ----
2  Below is a medical encounter between an {age}
3   and {sex} patient and a doctor done over chat.
4  Chief Complaint: "{cc}".
5  ----
6  Medical Encounter
7  ----
8  {chat}
9  ----
10 Summary Instructions
11 ----
12 Provide a summary of the medical encounter between the doctor and the {
       age_and_sex} patient in 6 sections (Demographics and Social Determinants of
        Health, Patient Intent, Pertinent Positives, Pertinent Unknowns, Pertinent
        Negatives, Medical History). The definitions of each section are listed
        below. Write a paragraph under each section, not bullet points.
13
14 Demographics and Social Determinants of Health:
15 // Definition of section
16
17 Patient Intent:
18 // Definition of section
19
20 Pertinent Positives:
21 // Definition of section
22
23 Pertinent Unknowns:
24 // Definition of section
25
26 Pertinent Negatives:
27 // Definition of section
28
29 Medical History:
30 // Definition of section
31
32 ----
33 Summary of Medical Encounter
34 ----
```

Prompt 1: Prompt for generating initial summary.

```
1  ---
2  Below is a medical encounter between a {age_and_sex} patient and a doctor done
       over chat.
3  Chief complaint: "{cc}".
4  ----
5  Medical Encounter
6  ----
7  {chat}
8  ----
9  Below is a summary of the conversation that was written using the following
       instructions:
10
11 // Definition of medical summary (same as in initial summarization prompt)
12 ----
13 Summary of Medical Encounter
14 ----
15 {summary}
16 ----
17 Using the above dialogue and provided summary, corrupt the summary slightly.
       This could include moving a positive symptom to be a negative symptom,
       making up medical history mentioned, etc.
18
19 Corruptions should only occur on the Pertinent Positives, Pertinent Unknowns,
       Pertinent Negative, or Medical History section.
20
21 The lower the desired corruption level, the fewer the changes made. Note that a
        0 would be not changing the summary at all, and a 10 would be completely
       corrupting the summary.
22
23 Note that any changes/corruption should make the summary less factual.
24
25 Desired Corruption Level: {corruption_level}/10
26 ----
27 Corrupted Summary of Medical Encounter
28 ----
```

Prompt 2: Prompt for generating corruptions based off of the initial summary.

```
 1  You (Person A) are a very good summary writer for medical dialogues between
       physicians and patients.
 2
 3  This is the medical dialogue you summarized for a {age} and {sex} patient:
 4  -Medical Dialogue-
 5  {chat}
 6  -Medical Dialogue-
 7
 8  You are discussing the summary you wrote for this dialogue with another summary
        writer (Person B) whose job it is to verify your summary for correctness.
 9
10  Person B will give you points for correction and it will be your job to add the
        points of correction to a scratchpad if you agree with them.
11
12  This is your original version of the summary:
13  -Your Original Summary-
14  {summary}
15  -Your Original Summary-
16
17  Here is your current scratchpad of corrections to make to the summary:
18  -Correction Scratchpad-
19  {scratchpad}
20  -Correction Scratchpad-
21
22  You are generally very confident about the summary you wrote, however, when
       presented with compelling arguments by the verifying summary writer, you
       add to the correction scratchpad. You also suggest any edits of your own in
        case you notice a mistake.
23
24  This is the summary discussion so far:
25  -Summary Discussion-
26  {discussion}
27  -Summary Discussion-
28
29  Question: What do you say next? Respond to Person B in the tag [RESPONSE: "<
       your_response_here>"] and output any corrections to add to the scratchpad
       in the tag [SCRATCHPAD: "<things_to_add_to_the_scratchpad_here>"]. Make
       sure to use the "[]" when outputting tags.
30  Answer:
```

Prompt 3: Prompt for decider agent used in DERA summarization experiments.

```
 1  ---
 2  You (Person B) are a very good summary editor for medical dialogues between
        physicians and patients.
 3
 4  This is the medical dialogue you will be referencing for a {age} and {sex}
        patient:
 5  -Medical Dialogue-
 6  {chat}
 7  -Medical Dialogue-
 8
 9  You are discussing the summary that another summary writer (Person A) wrote for
         this dialogue one section at a time.
10
11  You will be giving Person A points for correction based on any mistakes/
        discrepancies you see between the dialogue and summary one section at a
        time. Person A will add the points of correction that they agree on to a
        scratchpad to later make edits.
12
13  However, you will only go through the Pertinent Positives, Pertinent Negatives,
         Pertinent Unknowns, and Medical History sections.
14
15  This is Person A's original version of the summary:
16  -Person A's Original Summary-
17  {summary}
18  -Person A's Original Summary-
19
20  Here is Person A's current scratchpad of corrections to make to the summary:
21  -Correction Scratchpad-
22  {scratchpad}
23  -Correction Scratchpad-
24
25  Go through each section of the summary one at a time and point out any text
        that does not have a grounding in the dialogue. It must be possible to
        directly tie any span of the summary to the dialogue.
26
27  Make sure to make accurate, useful suggestions for corrections.
28
29  Person A may not initially agree with you, but if you are confident there is an
         error do your best to convince Person A of the mistake.
30
31  Once you have gone through each section and have confirmed each section with
        Person A, and you are satisfied with all of the corrections added to the
        scratchpad and/or all of Person A's reasoning to reject additional
        corrections, output the tag "[STOP]".
32
33  This is the summary discussion with Person A so far:
34  -Summary Discussion-
35  {discussion}
36  -Summary Discussion-
37
38  Question: What do you say next? Respond to Person A in the tag [RESPONSE: "<
        your_response_here>"]. If you are done correcting and are satisfied, output
         the "[STOP]" tag.
39  Answer:
```

Prompt 4: Prompt for researcher agent used in DERA summarization experiments.

```
 1  ---
 2  You are a very good summary writer for medical dialogues between physicians and
       patients.
 3
 4  This is the medical dialogue you summarized for a {age} and {sex} patient:
 5  -Medical Dialogue-
 6  {chat}
 7  -Medical Dialogue-
 8
 9  This is your original version of the summary:
10  -Original Summary-
11  {summary}
12  -Original Summary-
13
14  Here is your current scratchpad of corrections to make to the summary:
15  -Correction Scratchpad-
16  {scratchpad}
17  -Correction Scratchpad-
18
19  Make all changes mentioned in the scratchpad to the original summary to output
       the corrected summary.
20
21  Output the tag "[STOP]" when finished writing the corrected summary.
22
23  -Corrected Summary-
```

Prompt 5: Prompt for final summarization step (incorporating scratchpad of corrections into the original summary) used in DERA summarization experiments.

```
 1  Given the following snippet of a medical dialogue summary, extract the medical
       concepts (symptoms, diseases, conditions, allergies, lab tests, etc.)
       present.
 2
 3  The heading of the section from which the summary was extracted will also be
       provided.
 4
 5  ---Example 1---
 6  Pertinent Negatives:  Patient reports no <concept_1>, no <concept_2>, <
       concept_3>, and <concept_4>. Patient also reports having no trouble with <
       concept_5>.
 7
 8  Medical Concepts: [<concept_1>, <concept_2>, <concept_3>, <concept_4>, <
       concept_5>]
 9  ---Example 1---
10
11  ---Example 2---
12  Pertinent Positives:  Patient ongoing <concept_1> for the past 5 days, <
       concept_2>, and some <concept_3>. Patient had <concept_4> done in May 2021.
13
14  Medical Concepts: [<concept_1>, <concept_2>, <concept_3>, <concept_4>]
15  ---Example 2---
16
17  ---Example 3---
18  Pertinent Unknowns:  Patient is unsure about <concept_1> and <concept_2>.
19
20  Medical Concepts: [<concept_1>, <concept_2>]
21  ---Example 3---
22
23  ---Example 4---
24  Medical History: Patient reports some <concept_1> in the past, and had last <
       concept_2> on DATE_1.
25
26  Medical Concepts: [<concept_1>, <concept_2>]
27  ---Example 4---
28
29  Here is the example to extract medical concepts from:
30
31  {section_heading}: {section_value}
32
33  Medical Concepts:
```

Prompt 6: Prompt for extracting medical concepts from the summary used to compute the GPT-F1 metric.

```
 1  Given a snippet (snippet) from a medical dialogue summary and a corresponding
        list (list_a) of medical concepts extracted from that snippet, evaluate
        what medical concepts from a separate list (list_b) can be found in either
        list_a or snippet.
 2
 3  Note that on some occasions a medical concept from list_b may not be found in
        list_a, but can be appropriate to be present given the snippet. This could
        include rephrasings of medical concepts that are clinically equivalent (Ex:
         COVID and COVID-19).
 4
 5  ---Example---
 6  snippet: <snippet>
 7  list_a: [<concept_1>, <concept_2>, <concept_3>, <concept_4>, <concept_5>, <
        concept_7>]
 8  list_b: [<concept_0>, <concept_1>, <concept_3>, <concept_4>, <concept_5>, <
        concept_6>]
 9
10  found_b: [<concept_1>, <concept_3>, <concept_4>, <concept_5>]
11  not_found_b: [<concept_0>, <concept_6>]
12
13  ---Example---
14
15  Here is the snippet, list_a. Evaluate the medical concepts in list_b as above.
16
17  snippet: {snippet}
18  list_a: {list_a}
19  list_b: {list_b}
20
21  found_b:
```

Prompt 7: Prompt for verifying medical concepts from a summary section used to compute the GPT-F1 metric.

```
1  ----
2  Care Plan Instructions
3  ----
4  You are a primary care physician tasked with writing a care plan, which lists
       the next steps in care management that the patient and the physician will
       perform.
5  Categorize the next steps into five sections: Medications, Referrals, Tests,
       Lifestyle and Supportive Care. Definitions and scopes of each section are
       defined below.
6
7  Medications:
8  // Definition of section
9  Referrals:
10 // Definition of section
11 Tests:
12 // Definition of section
13 Lifestyle:
14 // Definition of section
15 Supportive Care:
16 // Definition of section
17
18 {example}
19 ----
20 Care Plan Instructions
21 ----
22 Now that you've seen an example, you will now write a care plan of the same
       format (five sections: Medications, Referrals, Tests, Lifestyle and
       Supportive Care).
23
24 The dialogue you will use to write a care plan about is a medical encounter
       between a {age} and {sex} patient and a doctor done over chat:
25 ----
26 Dialogue
27 ----
28 {chat}
29 ----
30 Care Plan
31 ----
```

Prompt 8: Prompt for generating initial care plan

```
 1  ---
 2  You (Person A) are a very good writer of care plans for patients following
        their discussion with a physician. The full instructions are presented
        below.
 3  ---
 4  Care Plan Writing Instructions
 5  ---
 6  // Same instructions as in initial care plan generation prompt. Removed for
        brevity.
 7  ---
 8  Given the instructions, this is the medical dialogue you see for a  {{age}} {{
        sex}} patient:
 9  ---
10  Medical Dialogue
11  ---
12  {chat}
13  ---
14  You are discussing the care plan you wrote for this dialogue with another care
        plan writer (Person B) whose job it is to verify your care plan for
        soundness.
15
16  Person B will give you points for correction and it will be your job to add the
         points of correction to a scratchpad if you agree with them.
17
18  This is your original version of the care plan:
19  ---
20  Your Original Care Plan
21  ---
22  {careplan}
23  ---
24  Here is your current scratchpad of corrections to make to the care plan:
25  ---
26  Correction Scratchpad
27  ---
28  {scratchpad}
29  ---
30  You are generally very confident about the care plan you wrote, however, when
        presented with compelling arguments by the verifying care plan writer, you
        add to the correction scratchpad. You also suggest any edits of your own in
         case you notice a mistake.
31
32  This is the care plan discussion so far:
33  ---
34  Care Plan Discussion
35  ---
36  {discussion}
37  ---
38  Question: What do you say next? Respond to Person B in the tag [RESPONSE: "<
        your_response_here>"] and output any corrections to add to the scratchpad
        in the tag [SCRATCHPAD: "<things_to_add_to_the_scratchpad_here>"]. Make
        sure to use the "[]" when outputting tags. All text should be within the
        tag brackets.
39  An example answer would be: [RESPONSE: "I think we should remove ... from the
        care plan"] [SCRATCHPAD: "Remove ... from the care plan because ..."]
40  ---
41  Answer:
```

Prompt 9: Prompt for decider agent used in DERA care plan experiments.

```
 1  ---
 2  You are a primary care physician and very good editor of care plans for
        patients following their discussion with a physician. The full instructions
         for writing care plans are presented below.
 3  ---
 4  Care Plan Writing Instructions
 5  ---
 6  // Same instructions as in initial care plan generation prompt. Removed for
        brevity.
 7  ---
 8  Given the instructions, this is the medical dialogue you see for a {age_and_sex
        } patient:
 9  ---
10  Medical Dialogue
11  ---
12  {chat}
13  ---
14
15  You are discussing the care plan that another care plan writer (Person A) wrote
         for this dialogue one section at a time.
16
17  You will be giving Person A points for correction based on any reconsiderations
         you see between the dialogue and care plan one section at a time. Person A
         will add the points of correction that they agree on to a scratchpad to
        later make edits.
18
19  This is Person A's original version of the care plan:
20  ---
21  Person A's Original Care Plan
22  ---
23  {careplan}
24  ---
25  Here is Person A's current scratchpad of corrections to make to the care plan:
26  ---
27  Correction Scratchpad
28  ---
29  {scratchpad}
30  ---
31  Go through each section of the care plan one section at a time and point out
        any suggestions that does not have a grounding in the dialogue. All
        suggestions must be grounded in information from the dialogue.
32
33  Remember to make sure the care plan is congruent with the Care Plan Writing
        Instructions.
34
35  Make sure to make accurate, useful suggestions for corrections.
36
37  Person A may not initially agree with you, but if you are confident there is an
         error do your best to convince Person A of the mistake.
38
39  Once you have gone through each section and have confirmed each section with
        Person A, and you are satisfied with all of the corrections added to the
        scratchpad and/or all of Person A's reasoning to reject additional
        corrections, output the tag "[DONE]".
40
41  This is the care plan discussion with Person A so far:
42  ---
43  Care Plan Discussion
44  ---
45  {discussion}
46  ---
47  Question: What do you say next? Respond to Person A in the tag [RESPONSE: "<
        your_response_here>"]. If you are done correcting, are satisfied, and want
        to end the conversation, output "DONE".
48  ---
49  Answer:
```

Prompt 10: Prompt for researcher agent used in DERA care plan experiments.

```
1  ---
2  You are a very good writer of care plans for patients following their
       discussion with a physician. The full instructions are presented below.
3  ---
4  Care Plan Writing Instructions
5  ---
6  // Same instructions as in initial care plan generation prompt. Removed for
       brevity.
7  ---
8  Given the instructions, this is the medical dialogue you see for a {age} and {
       sex} patient:
9  ---
10 Medical Dialogue
11 ---
12 {{chat}}
13 ---
14 You have been discussing the care plan you wrote for this dialogue with another
        care plan writer (Person B) whose job it is to verify your care plan for
       soundness.
15
16 You added corrections to a scratchpad after discussing them with Person B, and
       you will later be tasked with updating the original care plan based off of
       the correctness suggested in the scratchpad.
17
18 This is your original version of the care plan:
19 ---
20 Your Original Care Plan
21 ---
22 {careplan}
23 ---
24 Here is your current scratchpad of corrections to make to the care plan:
25 ---
26 Correction Scratchpad
27 ---
28 {scratchpad}
29 ---
30 Make all changes mentioned in the scratchpad to the original care plan to
       output the corrected care plan. Make sure all changes are congruent to the
       Care Plan Writing Instructions.
31
32 Output the tag "[STOP]" when finished writing the corrected care plan.
33 ---
34 Corrected Care Plan
35 ---
```

Prompt 11: Prompt for final care plan generation step (incorporating scratchpad of corrections into the original care plan) used in DERA care plan experiments.

```
 1  The following question was written as a multiple choice question.  Rewrite it
       as posing an open-ended question. If it is already an open-ended question
       and the question requires no rewrite, output "[OPEN]" only.  Do not change
       any details or facts in the question, and only change the phrasing of the
       question.
 2  --Example--
 3  Question: A 60-year-old man comes to the physician for an examination prior to
       a scheduled cholecystectomy. He has hypertension treated with
       hydrochlorothiazide. His mother had chronic granulomatous disease of the
       lung. He works in a glass manufacturing plant. He has smoked two packs of
       cigarettes daily for 38 years. His vital signs are within normal limits.
       Examination shows no abnormalities. Laboratory studies are within the
       reference range. An x-ray of the chest is shown. Which of the following is
       the most appropriate next step in management?
 4  Rewrite: A 60-year-old man comes to the physician for an examination prior to a
        scheduled cholecystectomy. He has hypertension treated with
       hydrochlorothiazide. His mother had chronic granulomatous disease of the
       lung. He works in a glass manufacturing plant. He has smoked two packs of
       cigarettes daily for 38 years. His vital signs are within normal limits.
       Examination shows no abnormalities. Laboratory studies are within the
       reference range. An x-ray of the chest is shown. What is the most
       appropriate next step in management?
 5  --Example--
 6  Question: Several patients at a local US hospital present with chronic
       secretory diarrhea. Although there are multiple potential causes of
       diarrhea present in these patients, which of the following is most likely
       the common cause of their chronic secretory diarrhea?
 7  Rewrite: Several patients at a local US hospital present with chronic secretory
        diarrhea.  Although there are multiple potential causes of diarrhea
       present in these patients, what is most likely the common cause of their
       chronic secretory diarrhea?
 8  --Example--
 9  Question: A 39-year-old male presents to your office with nodular skin lesions
       that progress from his right hand to right shoulder. The patient reports
       that the initial lesion, currently necrotic and ulcerative, developed from
       an injury he received while weeding his shrubs a couple weeks earlier. The
       patient denies symptoms of respiratory or meningeal disease. Which of the
       following most likely characterizes the pattern of this patient's skin
       lesions:
10  Rewrite: A 39-year-old male presents to your office with nodular skin lesions
       that progress from his right hand to right shoulder. The patient reports
       that the initial lesion, currently necrotic and ulcerative, developed from
       an injury he received while weeding his shrubs a couple weeks earlier. The
       patient denies symptoms of respiratory or meningeal disease. How would you
       characterize the pattern of this patient's skin lesions?
11  --Example--
12  Question: A 71-year-old man presents to the clinic with complaints of right
       wrist pain for 2 days. On examination, redness and swelling were noted on
       the dorsal aspect of his right wrist. He had  pain with extreme range of
       motion of the wrist. His history includes 2 hip replacements, 2 previous
       episodes of gout in both first metatarsophalangeal joints, and hypertension
       . Two days later, the swelling had increased in the dorsal aspect of his
       right wrist and hand. Wrist flexion was limited to 80% with severe pain.
       The pain was present on palpation of the scaphoid bone. Due to the
       suspicion of fracture, the patient was referred to his general practitioner
        for radiographs. These findings were consistent with gouty arthritis. What
        is the most likely cytokine involved in this process?
13  Rewrite: [OPEN]
14  ---
15  Question: {{question}}
16  Rewrite:
```

Prompt 12: Prompt for rewriting the question in full (temperature at 0 and otherwise uses default parameters)

```
 1  The following question was written as a multiple choice quesiton.  For the
       sentence in the question poses a multiple choice, rewrite it as posing an
       open-ended question. If the relevant is a compound sentence, re-write the
       entire sentence.  If it is already an open-ended question and the question
       requires no rewrite, output "[OPEN]" only.  Do not change any details or
       facts in the question, and only change the phrasing of the question.
 2  --Example--
 3  Question: A 60-year-old man comes to the physician for an examination prior to
       a scheduled cholecystectomy. He has hypertension treated with
       hydrochlorothiazide. His mother had chronic granulomatous disease of the
       lung. He works in a glass manufacturing plant. He has smoked two packs of
       cigarettes daily for 38 years. His vital signs are within normal limits.
       Examination shows no abnormalities. Laboratory studies are within the
       reference range. An x-ray of the chest is shown. Which of the following is
       the most appropriate next step in management?
 4  Original: Which of the following is the most appropriate next step in
       management?
 5  Rewrite: What is the most appropriate next step in management?
 6  --Example--
 7  Question: Several patients at a local US hospital present with chronic
       secretory diarrhea. Although there are multiple potential causes of
       diarrhea present in these patients, which of the following is most likely
       the common cause of their chronic secretory diarrhea?
 8  Original: Although there are multiple potential causes of diarrhea present in
       these patients, which of the following is most likely the common cause of
       their chronic secretory diarrhea?
 9  Rewrite: Although there are multiple potential causes of diarrhea present in
       these patients, what is most likely the common cause of their chronic
       secretory diarrhea?
10  --Example--
11  Question:A 39-year-old male presents to your office with nodular skin lesions
       that progress from his right hand to right shoulder. The patient reports
       that the initial lesion, currently necrotic and ulcerative, developed from
       an injury he received while weeding his shrubs a couple weeks earlier. The
       patient denies symptoms of respiratory or meningeal disease. Which of the
       following most likely characterizes the pattern of this patient's skin
       lesions:
12  Original: Which of the following most likely characterizes the pattern of this
       patient's skin lesions:
13  Rewrite: How would you characterize the pattern of this patient's skin lesions?
14  --Example--
15  Question: A 71-year-old man presents to the clinic with complaints of right
       wrist pain for 2 days. On examination, redness and swelling were noted on
       the dorsal aspect of his right wrist. He had  pain with extreme range of
       motion of the wrist. His history includes 2 hip replacements, 2 previous
       episodes of gout in both first metatarsophalangeal joints, and hypertension
       . Two days later, the swelling had increased in the dorsal aspect of his
       right wrist and hand. Wrist flexion was limited to 80% with severe pain.
       The pain was present on palpation of the scaphoid bone. Due to the
       suspicion of fracture, the patient was referred to his general practitioner
        for radiographs. These findings were consistent with gouty arthritis. What
        is the most likely cytokine involved in this process?
16  Original: What is the most likely cytokine involved in this process?
17  Rewrite: [OPEN]
18  ---
19  Question: {{question}}
20  Original:
```

Prompt 13: Prompt for rewriting the question by changing the last sentence only (temperature at 0 and otherwise uses default parameters).

```
1  Given the following medical question, respond with the phrase that best answers
      the question.
2
3  --Example--
4  Question: A mother brings her 3-week-old infant to the pediatrician's office
      because she is concerned about his feeding habits. He was born without
      complications and has not had any medical problems up until this time.
      However, for the past 4 days, he has been fussy, is regurgitating all of
      his feeds, and his vomit is yellow in color. On physical exam, the child's
      abdomen is minimally distended but no other abnormalities are appreciated.
      What embryologic error could account for this presentation?
5
6  What phrase best answers the question posed?
7
8  Answer: Abnormal migration of ventral pancreatic bud
9  ----
10 Question: {question}
11
12 What phrase best answers the question posed?
13
14 Answer:
```

Prompt 14: Prompt for generating the single-shot answer.

```
1  {question}
2
3  {options_filtered_str}
4
5  You think the relative likelihood of each option is {relative_likelihood}.
      Write a 3-4 sentence message explaining why you rate the options in that
      way, without taking a decisive stand.
6
7  Message:
```

Prompt 15: Prompt for generating the explanation for the single-shot answer distribution.

```
1  You are an expert medical doctor who is guiding a medical student through
       thinking about which of several answers is best for a given question.  You
       cannot give the student the answer.  Your role is to help the student think
       through the question, specifically by pointing out portions of the
       question that are important in understanding the problem.
2  Rules;
3  - All responses should include a quote from the question.
4  - Consider what you, as the teacher, have said in the previous conversation,
       and do not repeat yourself.
5  - Responses should be at most 4 sentences long.
6  - Stop only when you, as the teacher, have pointed out all important aspects of
        the question in the previous discussion.  To stop, respond with 'STOP' at
       the next turn.
7   You cannot;
8   - Directly give the answer to the student
9   - Include the correct option in your response, or any paraphrasing of the
       correct answer.
10  - Do not narrow down the options in your response.
11
12 Question: {question}
13
14 The previous discussion between you and the expert advisor is as follows;
15 {chat_history}
16 {last_student_message}
17
18 Help the student find the correct answer by pointing out specific parts of the
       questions they need to think through, but do not include the correct phrase
        in your response. Your response should be no more than 3-4 sentences.  If
       you have pointed out all challenging aspects of the question in the
       previous conversation, respond with "STOP" after the student's next turn.
19
20 Response:
```

Prompt 16: Prompt for question-answering *Researcher*.

```
1  You are an expert doctor who is trying to select the answer to a medical
       question, and is willing to be open-minded about their answer.   The
       questions are taken from a short-answer medical exam, and your role is to
       arrive at the correct answer.
2
3  You are chatting with an expert medical advisor, who will try to help you think
        through the problem, but will not directly tell you the answer.  They will
        help you by pointing out aspects of the question that are important in
       finding the answer.  Do not assume that the teacher knows the answer; only
       that they know how to think through the question. You can change your
       answer at any point, but do not assume that the expert knows the exact
       answer and is providing leading questions. Think about their guidance as a
       whole, and do not only respond to their last message
4
5  Question: {question}
6
7  The previous discussion between you and the expert advisor is as follows;
8  {chat_history}
9  {last_teacher_message}
10
11 Rethink the question by considering what the teacher pointed out, in light of
       your original hypothesis.  Remember they do not know the answer, but only
       how to think through the question. You can change your mind on the correct
       answer, but remember that unless the question explicitly asks for multiple
       answers, you can only provide a single answer. Respond with the option you
       believe most likely to be the right answer ("Answer:<SHORT ANSWER>") and a
       response to that message ("Response:<MESSAGE>"):
12
13 Answer:
```

Prompt 17: Prompt for question-answering *Decider*.

```
1  You are an expert doctor who is trying to select the answer to a medical
       question, and is willing to be open-minded about their answer.   The
       questions are taken from a short-answer medical exam, and your role is to
       arrive at the correct answer.
2
3  You are chatting with an expert medical advisor, who will try to help you think
        through the problem, but will not directly tell you the answer.  They will
        help you by pointing out aspects of the question that are important in
       finding the answer.  Do not assume that the teacher knows the answer; only
       that they know how to think through the question. You can change your
       answer at any point, but do not assume that the expert knows the exact
       answer and is providing leading questions. Think about their guidance as a
       whole, and do not only respond to their last message
4
5  Question: {question}
6
7  The previous discussion between you and the expert advisor is as follows;
8  {chat_history}
9  {last_teacher_message}
10
11 Rethink the question by considering what the teacher pointed out, in light of
       your original hypothesis.  Remember they do not know the answer, but only
       how to think through the question. You can change your mind on the correct
       answer, but remember that unless the question explicitly asks for multiple
       answers, you can only provide a single answer. Respond with the option you
       believe most likely to be the right answer ("Answer:<SHORT ANSWER>") and a
       response to that message ("Response:<MESSAGE>"):
12
13 Answer:
```

Prompt 18: Prompt for question-answering final answer.

```
1  Assign a dxSimilarityScore to each of the following pairs where the first
       diagnosis is an "expectedDx" and the second diagnosis is the "
       providedDiagnosis".
2
3  Expected Vs Provided Dx Pairs:
4  {answer_text} | {predicted_answer_text}
5  {answer_text} | {zero_shot_option_index}
6
7  Output each pair in one line using this format "dx1" "|" "dx2" "|" "
       dxSimilarityScore"
8  output:
```

Prompt 19: Prompt similar to that used for similarity score between generated and gold answers. Note that occasionally this outputs a number outside of 0-1. Unless these are all 100s we set these to 0s. This commonly occurs with math problems.

```
1  Question:{question}
2
3  Do the following two answers refer to the same medical concept? Respond with an
        answer ("Answer:True" or "Answer:False") followed by an explanation ("
       Explanation:")
4
5  {answer_text}
6  {predicted_answer_text}
7
8  Answer:
```

Prompt 20: Prompt for exact matching between generated and gold answers.

# LlamaMTS: Optimizing Metastasis Detection with Llama Instruction Tuning and BERT-Based Ensemble in Italian Clinical Reports

**Livia Lilli[1,2], Stefano Patarnello[1], Carlotta Masciocchi[1], Valeria Masiello[3],**
**Fabio Marazzi[3], Luca Tagliaferri[3], Nikola Dino Capocchiano[1]**

[1] Real World Data Facility, Gemelli Generator, Gemelli Hospital of Rome
[2] Catholic University of the Sacred Heart of Rome
[3] Department of Diagnostic Imaging, Radiation Oncology and Hematology,
UOC of Radiation Oncology, Gemelli Hospital of Rome
`livia.lilli@policlinicogemelli.it`

## Abstract

Information extraction from Electronic Health Records (EHRs) is a crucial task in healthcare, and the lack of resources and language specificity pose significant challenges. This study addresses the limited availability of Italian Natural Language Processing (NLP) tools for clinical applications and the computational demand of large language models (LLMs) for training. We present LlamaMTS, an instruction-tuned Llama for the Italian language, leveraging the LoRA technique. It is ensembled with a BERT-based model to classify EHRs based on the presence or absence of metastasis in patients affected by Breast cancer. Through our evaluation analysis, we discovered that LlamaMTS exhibits superior performance compared to both zero-shot LLMs and other Italian BERT-based models specifically fine-tuned on the same metastatic task. LlamaMTS demonstrates promising results in resource-constrained environments, offering a practical solution for information extraction from Italian EHRs in oncology, potentially improving patient care and outcomes.

## 1 Introduction

Electronic health records (EHRs) represent the principal data source for hospital centers, housing invaluable information regarding medical histories, treatments, examinations, disease progression and symptoms of a patient. However, efficiently extracting this data with high accuracy and minimal computational resources presents a growing challenge, particularly in the context of the Italian language. While solutions specialized in the clinical domain are readily available for the English language (Lee et al., 2020; Luo et al., 2022; Labrak et al., 2024; Wang et al., 2024), the exploration of similar solutions for the Italian language remains limited, with only a handful of alternatives (Buonocore et al., 2023). Consequently, our objective is to investigate

novel approaches that could be implemented in real-world clinical contexts, to extract specific outcomes from Italian textual data. For this purpose, we searched for methods that enable fine-tuning of large language models for specific tasks while minimizing computational resource consumption. Recent studies have showcased the effectiveness of implementing instruction-tuning on pre-trained large language models (Wei et al., 2021; Chung et al., 2022; Liu et al., 2024; Wang et al., 2022), also leveraging techniques such as LoRA (Hu et al., 2022). Through this approach, the number of trainable parameters is reduced, and the model is trained to respond to specific instructions provided during training.

In this paper, we introduce LlamaMTS (Figure 1), a fine-tuned Llama model, through the LoRA instruction tuning technique. Our model is designed to identify the presence of tumoral metastasis by analyzing EHRs from patients diagnosed with breast cancer. Llama was fine-tuned by using as base model Camoscio (Santilli and Rodolà, 2023), which is a Llama adapter for the Italian language, trained on the Italian translation of the Stanford Alpaca Dataset (Taori et al., 2023). To further enhance model performance, we employed an ensemble approach by incorporating a BERT-based model fine-tuned on the same classification task. Additionally, to allow the model to learn from entire EHRs (which may exceed the maximum token limit allowed by Llama during training), we implemented text summarization on both the training and testing datasets. This enabled information extraction from shorter and more concise texts, reducing the noise that long texts may cause.

To evaluate LlamaMTS performances, we compare it with several benchmarks, exploring zero-shot LLMs configurations and fine-tuning known BERT-based model for text classification. Results show that our approach, which leverages instruction-tuning and model ensembling, outper-

162

Figure 1: LlamaMTS framework

forms all the other baselines on our metastatic classification task.

## 2 Background

### 2.1 Clinical Text Classification

In the Italian language domain, the availability of pre-trained language models for text classification, especially in the clinical field, is currently limited. Notable mentions include BioBit, MedBit and MedBIT-r3-plus, which are different versions of pre-trainings on Italian clinical texts, proposed by Buonocore et al. (2023). In particular, BioBit relies on Italian translations of PubMed abstracts, MedBit is trained on medical textbooks originally written in Italian, while MedBIT-r3-plus is trained on Italian textbooks augmented with web-crawled data. Other works for the Italian language of interest for our study are: AlBERTo (Polignano et al., 2019), an Italian version of BERT (Devlin et al., 2018) trained on Italian tweets, GePpeTto (De Mattei et al., 2020), an Italian fine-tune version of GPT-2 base (117 million parameters), IT5 (Sarti and Nissim, 2022), a T5 model tailored for Italian and BART-IT (La Quatra and Cagliero, 2022), an Italian variant of BART (Lewis et al., 2019). Finally Abdaoui et al. (2020) proposed a set of multilingual models (including the Italian language), pre-trained on a reduced number of parameters.

### 2.2 Instruction Tuning

Recent works demonstrated the efficacy of implementing instruction-tuning on a pre-trained large language model, to increase the downstream performances (Wei et al., 2021; Chung et al., 2022; Liu et al., 2024; Wang et al., 2022). A first step in this direction was made by Taori et al. (2023), who presented Stanford Alpaca, an instruction-tuned version of Llama in the English language. Following this approach, further instruction-tuned Llama models have been trained with LoRA (Hu et al., 2022), as the English Alpaca Lora (Wang, 2023), the Portuguese Cabrita (Larcher et al., 2023) and the Italian Camoscio (Santilli and Rodolà, 2023). In addition to Camoscio, Bacciu et al. (2023) presented Fauno, a language model trained on a corpus of self-chat performed by ChatGPT. Compared to Camoscio, Fauno is a conversational agent for the Italian language. Similarly, Michael (2023) released Stambecco, an instruction-tuned version of LLaMA on a translation to Italian of the GPT-4-LLM dataset (Peng et al., 2023).

This study is inspired by the approach of Hromei et al. (2023), implementing the LoRA instruction-tuning on the Italian Camoscio adapter of Santilli and Rodolà (2023). In this study, the output is represented by extremITLLaMa, a fine-tuning on the EVALITA task (Lai et al., 2023).

Figure 2: Distribution of metastasis outcome overall the data, distinguished by set type.

## 2.3 Ensemble

Ensemble is an approach widely used to improve model performance in medical applications, especially in the case of raw data (Nilashi et al., 2022; Doppala et al., 2022; Dutta et al., 2022) and images (Khamparia et al., 2020; Tasci et al., 2021). However, recent works have applied these techniques to the domain of natural language processing, (Yang et al., 2023; Abdennour et al., 2023; Chen et al., 2023; Zhou et al., 2023). In our work, we adopt the approach of Zhou et al. (2023) which combined the BERT predictions with the generated tokens of a large language model to obtain the final ensemble output. We also compared this with the average voting approach of Dutta et al. (2022), where the predicted probability of a class, is a weighted average over all the models.

## 3 Method

Our methodology involves 1) the selection of the data corpus for fine-tuning, 2) The summarization of EHRs to obtain shorter texts, 3) the instruction tuning of an existing large language model on the metastatic classification task and 4) the ensemble of the obtained instruction-tuned model with a BERT-based model fine-tuned on the same task.

### 3.1 Data Corpus

In this study, we used EHRs from a data mart consisting of a collection of structured and textual data referencing patients diagnosed with Breast Cancer and being treated at the Italian Gemelli Hospital of Rome.

We selected all the data sources for extracting information relating to tumour metastasis. Guided by a team of physicians, we chose data on clinical diaries, medical histories, and radio-diagnostic reports, because these texts typically contain past and current information about the patient's health status and examination results. We extract all the relevant EHRs for this study from the Gemelli Breast data mart (Marazzi et al., 2021).

### 3.2 Text Summarization

The EHR length distribution was highly varied and a large portion of the data would risk not being fully processed, due to limits in maximum number of tokens allowed by many large language models.

Additionally, text semantics of clinical reports can be very complex, with relevant information (in this case, the presence or absence of metastasis) not always explicitly reported.

For this reason, we decided to use text summarization methodologies to include data with a reasonable range of tokens, written in a simpler form.

For this purpose, we chose to use Mixtral 8x7B, a pretrained generative Sparse Mixture of Experts language model (Jiang et al., 2024), which outperforms Llama 2 70B (Touvron et al., 2023b) on many benchmarks. To safeguard the confidentiality of clinical reports, we chose to employ locally executable models like Mixtral, thereby excluding the use of GPT (Achiam et al., 2023).

Finally, we formulated an Italian prompt meant to generate a summary of a few words of the input report, retaining all the information relevant to metastasis. We also provided a list of synonymous terminologies as instruction to the model, ensuring a more accurate topic detection. The final prompt was written as follows: *Dato il seguente referto, restituisci una sintesi coincisa in lingua italiana di poche parole, mantenendo tutte le informazioni inerenti a metastasi, lesioni, noduli, attività metabolica o staging:* `{EHR Text}`.

For the implementation of the Mixtral model, we leveraged the Ollama Python library[1].

### 3.3 Instruction Tuning

During the instruction tuning phase, we leveraged the Camoscio language model proposed by Santilli and Rodolà (2023), who fine-tuned the smallest

---

[1] https://github.com/ollama/ollama-python

version of Llama (Touvron et al., 2023a) on the Italian translation of Alpaca instruction-tuning dataset (Taori et al., 2023), using the LoRA technique (Hu et al., 2022). Then, following the methodology of Hromei et al. (2023), we merged the Camoscio adapter[2] to the original Llama model and fine-tuned it on our Italian classification task.

The dataset we used has the `instruction`, `input` and `output` fields, where `input` contains the summarized EHRs, the `output` is the binary information about the presence or absence of metastasis, and the `instruction` is written as follows: *Dato il seguente referto medico in italiano, indica con 1 presenza di metastasi e con 0 assenza di metastasi.*

These fields are then put together, for generating the final prompt; we used the same prompter template of Camoscio.

### 3.4 Ensemble

In order to enhance the final classification performance, we adopted an ensemble approach by combining our instruction-tuned LLM, with the BERT-based model having the best performance among our experiments.

Our approach takes inspiration from Zhou et al. (2023), where the final ensemble prediction corresponds to the one with the highest confidence among the two models, as shown in Equation 1.

$$pred_{ENS} = \begin{cases} pred_{LLM} & \text{if } prob_{LLM} > prob_{BERT} \\ pred_{BERT} & \text{if } prob_{LLM} < prob_{BERT} \end{cases} \quad (1)$$

For the BERT-based model, the confidence $prob_{BERT}$ is the prediction probability related to the predicted class. While for our instruction-tuned model, we considered $p_{LLM}$ as the predicted probability of the generated tokens. Thus, the final ensemble prediction corresponds to the most confident prediction produced by either the two models.

To compare different methods, we also applied a further ensemble technique, the average approach used by Dutta et al. (2022). In this approach, given $M = 2$ models and C=2 classes, we considered: the model output $Y_j \in \mathbb{R}^C$ for each $j^{th}$-model, and the confidence values $P_i \in \mathbb{R}^M$ for each $i^{th}$ class with $i \in \{0, 1\}$. So, the final ensemble confidence for a given class $k$ is defined as a weighted combination of all the models:

Figure 3: Distribution of Llama tokens for the EHR data.

$$P_k^{ens} = \frac{\sum\limits_{j=1}^{2} P_{kj} \times W_j}{\sum\limits_{i=0}^{1} \sum\limits_{j=1}^{2} P_{ij} \times W_j} \quad (2)$$

In the above equation, $W_j$ is the weight of the $j^{th}$ classifier. Once we have the output $Y \in \mathbb{R}^C$, which contains the confidence values $P_i^{ens} \in [0, 1]$ computed on the unseen data $X$, the final prediction will be the $i$-class, such that: $\arg\max_i Y(X)$.

## 4 Experiments

We started by generating the instruction-tuned model on the metastatic classification task, using the summarized EHRs as training data.

We then compared the performance of the instruction-tuned model with several baseline methods, including BERT-based approaches fine-tuned on our classification task and large language models implemented in a zero-shot environment.

Finally, we applied the ensemble techniques between the instruction-tuned model and the BERT-based fine-tuned model with the highest performance in order to obtain the final LlamaMTS classifier.

| Model | Precision | Recall | Accuracy | AUC | F-Score |
|---|---|---|---|---|---|
| *Zero-Shot LLM* | | | | | |
| Mixtral 8x7B | 92,3 | 71,6 | 72,6 | **74** | 80,6 |
| LLaMa2 7B | 79,8 | 1 | 79,8 | 50 | **88,7** |
| Camoscio | 79,8 | 1 | 79,8 | 50 | 88,7 |
| *BERT-Based Fine-Tuning* | | | | | |
| dbmdz BERT | 84,5 | 89,5 | 78,6 | 62,4 | 86,9 |
| BioBIT | 84,8 | 86,6 | 76,7 | 62,4 | 85,6 |
| MedBIT | 88,7 | 88 | 81,5 | **72** | **88,4** |
| MedBIT-r3-plus | 86,3 | 85 | 77,4 | 66 | 85,7 |
| mBERT-Ita | 85 | 88 | 78 | 63,1 | 86,4 |
| DistilBERT | 84,1 | 87,3 | 76,8 | 61,3 | 85,7 |
| RoBERTa-Ita | 79,5 | 98,5 | 78,6 | 49,3 | 88 |
| BERT-Tiny-Ita | 79,5 | 98,5 | 78,6 | 49,3 | 88 |
| *Instruction Tuning* | | | | | |
| Instruction-Tuned LlamaMTS | 80,2 | 1 | 80,4 | **51,5** | **89** |

Table 1: Results of the intruction-tuned LlamaMTS, compared with the zero-shot large language models and the BERT-based fine-tuning experiments, on the metastatic classification task.

## 4.1 Data and Summarization

Starting from our selected data corpus, we focused on a subsample of 1168 EHRs, randomly selected from three different data sources, clinical diaries, medical histories, and radio-diagnostic reports.

A total of 168 EHRs (14% of the available data) were annotated by a team of physicians and used as the gold standard for the final evaluation. In contrast, the remaining 1000 EHRs (86% of the overall data) were used as a training set for the model fine-tuning. As shown in Figure 2, the 80% of gold standards (which corresponds to 134 of the 168 EHRs) were positive to the presence of metastasis, while training set had the 65% of positive-labeled samples (that are 647 of the overall 1000 train reports).

We then analyzed the number of tokens in the final texts, using the model tokenizer. Figure 3 shows that the original EHR data has a median of 938 tokens, with first and third quartiles equal to 577 and 1351 respectively and with a maximum value that achieves 4453 tokens.

Considering the maximum number of tokens supported by Llama (2048) (Touvron et al., 2023a), we adopted approaches to reduce the size of the input texts used in our instruction-tuning environment. For this reason, we opted for the text summarization approach using Mixtral 8 x7B (Jiang et al., 2024), which returned a summarized version of the original data, whose tokens' distribution has a median of 301.5, with a first and third quartiles

respectively equal to 255 and 360 tokens (as shown in Figure 3). For privacy reasons, we do not report practical examples of summaries, but we provide summary metrics.

## 4.2 Instruction-Tuned Model

Our first experiment concerns the instruction tuning of the smallest version of Llama (Touvron et al., 2023a) through the Italian adapter of Santilli and Rodolà (2023). Following the Camoscio repository[3], we set up the fine-tuning by first preparing the input base model. We then merged the adapter checkpoints with the original Llama model and then selected 10 epochs for training, using the 1000 summarized clinical texts described in the above paragraph as inputs. We also set the cutoff length at the maximum value supported by Llama, i.e. 2048 tokens.

In the inference phase, we forced the maximum number of generated tokens to 1. We also prefixed the generation of tokens in order to output binary values for classification.

The resulting model represents our instruction-tuned LlamaMTS that will be ensembled with the best-performing BERT-based model to create the final LlamaMTS classifier.

## 4.3 BERT-Based Fine-Tuning

As baseline experiments, we considered several BERT models available on Hugging Face (Wolf

---
[3]https://github.com/teelinsan/camoscio

| Model | Precision | Recall | Accuracy | AUC | F-Score |
|---|---|---|---|---|---|
| Ensemble *Max Method* | 88,1 | 88,8 | 81,5 | 72,3 | 88,8 |
| Ensemble *AUC-Weighted Avg* | 88,7 | 88 | 81,5 | 72 | 88,4 |
| Ensemble *F1-Weighted Avg* | 88,8 | 88,8 | **82,1** | **72,3** | **88,8** |

Table 2: Results of the ensemble between the instruction-tuned LlamaMTS and the best BERT-based model. The third approach, about the F1-Weighted Average, represents our final LlamaMTS classifier.

et al., 2020) for the Italian language, fine-tuning them on our classification task. The fine-tuning was performed for 10 epochs and the models we chose are pre-trained in the Italian language.

We focused on the work of Buonocore et al. (2023), using their three models (BioBit[4], MedBit[5] and MedBIT-r3-plus[6]), which are different versions of pre-trainings on Italian clinical texts.

Additionally, we explored the work of Abdaoui et al. (2020), fine-tuning their multilingual models[7], pre-trained on a reduced number of parameters.

Finally, we applied other available models trained in the Italian language[8], for further comparisons.

## 4.4 Zero-Shot LLM

As additional baselines, we considered the classification capability of conversational large language models, forcing the answers to be binary values (meaning presence or absence of metastasis). We chose the two best-performing open-source models, Llama2 (Touvron et al., 2023b) and Mixtral (Jiang et al., 2024), using the Ollama Python library[9], with a prompt in the Italian language. The prompt asks to return an integer number for the given task, where the task is to output a binary value indicating the presence or absence of metastasis in the given text. The final prompt was written as follows: *'Per il seguente task, restituisci solo un numero come risposta, senza ulteriore testo. Dato il seguente referto, rispondi con "1" se è indicata presenza di metastasi, altrimenti rispondi con "0":* {EHR Text}.

Whenever other strings are returned in addition to the binary output, then a regex search of the desired values is performed on the generated response, to produce the appropriate binary value.

Moreover, to show the advantage of performing the Llama instruction-tuning, we also applied the Camoscio checkpoints on the same metastatic classification task, with the same inference configuration previously discussed for the instruction-tuned LlamaMTS in subsection 4.2. We chose to focus just on Llama2, as it was the only version available in the Ollama library.

## 4.5 Ensembling Models

The instruction-tuned LlamaMTS was then combined with the best-performing BERT-based fine-tuned model, to achieve improvements in the final performance metrics. We implemented two different ensemble approaches, as described in subsection 3.4, and considered the ensemble with the best performances as our final LlamaMTS classifier.

In the ensemble experiments, we didn't consider the LLM-based models, because we couldn't compute the corresponding predicted probabilities. For this reason, these models are only used as a baseline benchmark, for a first comparison of the results.

The ensemble results consist of three experiments, where the first one leverages on the approach described by Equation 1, while the second and the third implementations are based on Equation 2, using AUC and F-Score as weights respectively.

## 4.6 Results and Discussion

Results are measured through the Python Scikit-Learn package (Pedregosa et al., 2011) by computing the typical scores for classification tasks: Precision, Recall, Accuracy, F-Score, and AUC. For the evaluation of the models' performances, we focus on the F-Score and on the AUC metrics, which are typically preferred to Accuracy when the test set is not perfectly balanced among classes. In our case, gold standards present the 80% of positive metastatic samples overall the 168 EHRs.

Table 1 shows that our instruction-tuned LlamaMTS presents the best performances in terms of F-Score, which is 89%. In particular, it presents good sensitivity, that is approximately 100%, and

---

[4]IVN-RIN/bioBIT

[5]IVN-RIN/medBIT

[6]IVN-RIN/medBIT-r3-plus

[7]Geotrend/bert-base-it-cased, Geotrend/distilbert-base-it-cased

[8]osiria/roberta-base-italian, mascIT/bert-tiny-ita

[9]https://github.com/ollama/ollama-python

a precision of over the 80%. However, we got an AUC of 51.5%, which is lower when compared with the other models. As far as computational resources are concerned, the Llama instruction-tuning spent about 6h 57m 47s, by using an Nvidia RTX 5000 Graphics Processing Unit (GPU) and 16GB of Random Access Memory (RAM).

Among the BERT-based fine-tuned models, MedBIT shows the best metrics in terms of both AUC and F-Score, which are equal to 72% and 88.4% respectively. All the BERT-based experiments present F-Scores over 85%, but an AUC that ranges between 49% and 72%.

With the zero-shot learning of generative large language models, we got the highest results in terms of AUC with Mixtral (74%). Llama2 does not perform well in terms of AUC, that is 50%, though it has a higher F-Score when compared to Mixtral, with a value of 88.7%. Moreover, Llama2 and Camoscio present identical results: this suggests that the adaptation of Llama1 to Italian in Camoscio does not yield superior results compared to the advancements achieved by Llama2, which involved pre-training on a larger Italian corpus.

Additionally, Table 2 shows the results for the three ensemble experiments, performed combining the instruction-tuned LlamaMTS with the fine-tuned MedBit. The approach based on the selection of the highest confident prediction, and the average approach weighted by the F-Score, present the best performances, both having AUC and F-Score equal to 72.3% and 88.8%. Moreover, the F1-average approach has also a higher accuracy of 82.1% (if compared to the 81.5% of the first technique). Then this last method returns the final LlamaMTS classifier, with an AUC that is higher if compared to the instruction-tuned model and MedBit, and with an F-Score that is halfway between the values obtained from the two ensembled models.

## 5 Conclusions

The instruction-tuning allowed us to specialize an existing large language model on a medical classification task in an optimized fine-tuning environment, using the LoRA approach. Our study shows that LlamaMTS, which is a fne-tuned LLM using LoRA, has higher performance metrics when compared to the base model Camoscio and to other existing approaches that involve conversational LLMs and BERT-Based checkpoints (Table 1). Indeed, the instruction-tuned classifier tends to iden-

tify well all the existing positives, even if with low performances in distinguishing the negative samples. This is reflected in the high F-score of 89% and the low 51.5% AUC. We then applied the ensemble technique, combining the classification capability of the instruction-tuned model, with the best-performing BERT-based fine-tuned model. Thus we obtained our final LlamaMTS classifier, which outperforms both the models in terms of AUC, achieving a value of 72.3%, and with an F-Score of 88.8%, close to that of instruction-tuned model.

With this work, we extended advanced NLP techniques on clinical EHR data, automating processes through the usage of powerful language models, trained in the Italian language, on a specific classification task, for the extraction of the tumor metastasis information from EHRs. We proposed an approach that is easily portable to other kinds of outcomes, for extracting information not necessarily available in a structured format, from textual EHRs. Furthermore, the instruction-tuning approach enables fine-tuning large language models in reasonable time frames, leveraging mid-range computational resources.

## Limitations

While our study presents promising results for metastasis classification in Breast cancer patients, several limitations may be investigated in future research. These include the application of the model to new outcomes beyond metastasis and its adaptation to both binary and multi-classification tasks. Additionally, new work could be focused on testing the portability of the model by evaluating its performance on EHRs from new hospitals. Furthermore, improvements in model performance could be explored through extended fine-tuning on additional epochs and training data.

## Ethics Statement

For this study, the use of electronic health records was essential for training and testing our new technology. However, these data contain sensitive patient information and it was fundamental adhering to strict privacy and confidentiality guidelines. To this purpose, the dataset used in this paper was fully de-identified and we received approval from our institution to conduct the presented research. Approval protocol number from the relevant Ethics Committee can be provided on request.

## Acknowledgements

## References

Amine Abdaoui, Camille Pradel, and Grégoire Sigel. 2020. Load what you need: Smaller versions of multilingual bert. *arXiv preprint arXiv:2010.05609*.

Ghada Ben Abdennour, Karim Gasmi, and Ridha Ejbali. 2023. Ensemble learning model for medical text classification. In *International Conference on Web Information Systems Engineering*, pages 3–12. Springer.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Andrea Bacciu, Giovanni Trappolini, Andrea Santilli, Emanuele Rodolà, Fabrizio Silvestri, et al. 2023. Fauno: The italian large language model that will leave you senza parole! In *IIR 2023*. Pisa.

Tommaso Mario Buonocore, Claudio Crema, Alberto Redolfi, Riccardo Bellazzi, and Enea Parimbelli. 2023. Localizing in-domain adaptation of transformer-based biomedical language models. *Journal of Biomedical Informatics*, 144:104431.

Chao-Yi Chen, Kao-Yuan Tien, Yuan-Hao Cheng, and Lung-Hao Lee. 2023. Ncuee-nlp at semeval-2023 task 7: Ensemble biomedical linkbert transformers in multi-evidence natural language inference for clinical trial data. In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 776–781.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

Lorenzo De Mattei, Michele Cafagna, Felice Dell'Orletta, Malvina Nissim, and Marco Guerini. 2020. Geppetto carves italian into a language model. *arXiv preprint arXiv:2004.14253*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Bhanu Prakash Doppala, Debnath Bhattacharyya, Midhunchakkaravarthy Janarthanan, Namkyun Baik, et al. 2022. A reliable machine intelligence model for accurate identification of cardiovascular diseases using ensemble techniques. *Journal of Healthcare Engineering*, 2022.

Aishwariya Dutta, Md Kamrul Hasan, Mohiuddin Ahmad, Md Abdul Awal, Md Akhtarul Islam, Mehedi Masud, and Hossam Meshref. 2022. Early prediction of diabetes using an ensemble of machine learning models. *International Journal of Environmental Research and Public Health*, 19(19):12378.

CD Hromei, D Croce, V Basile, R Basili, et al. 2023. Extremita at evalita 2023: Multi-task sustainable scaling to large language models at its extreme. In *CEUR WORKSHOP PROCEEDINGS*, volume 3473, pages 1–9.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Aditya Khamparia, Aman Singh, Divya Anand, Deepak Gupta, Ashish Khanna, N Arun Kumar, and Joseph Tan. 2020. A novel deep learning-based multi-model ensemble method for the prediction of neuromuscular disorders. *Neural computing and applications*, 32:11083–11095.

Moreno La Quatra and Luca Cagliero. 2022. Bart-it: An efficient sequence-to-sequence model for italian text summarization. *Future Internet*, 15(1):15.

Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. *arXiv preprint arXiv:2402.10373*.

Mirko Lai, Stefano Menini, Marco Polignano, Valentina Russo, Rachele Sprugnoli, Giulia Venturi, et al. 2023. Evalita 2023: Overview of the 8th evaluation campaign of natural language processing and speech tools for italian. In *Proceedings of the Eighth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2023), CEUR. org, Parma, Italy*.

Celio Larcher, Marcos Piau, Paulo Finardi, Pedro Gengo, Piero Esposito, and Vinicius Caridá. 2023. Cabrita: closing the gap for foreign languages. *arXiv preprint arXiv:2308.11878*.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461.*

Fei Liu, Xi Lin, Qingfu Zhang, Xialiang Tong, and Mingxuan Yuan. 2024. Multi-task learning for routing problem with cross-problem zero-shot generalization. *arXiv preprint arXiv:2402.16891.*

Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *International Conference on Learning Representations.*

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics,* 23(6):bbac409.

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.

Fabio Marazzi, Luca Tagliaferri, Valeria Masiello, Francesca Moschella, Giuseppe Ferdinando Colloca, Barbara Corvari, Alejandro Martin Sanchez, Nikola Dino Capocchiano, Roberta Pastorino, Chiara Iacomini, et al. 2021. Generator breast datamart—the novel breast cancer data discovery system for research and monitoring: Preliminary results and future perspectives. *Journal of Personalized Medicine,* 11(2):65.

Michael. 2023. Stambecco: Italian instruction-following llama model. https://github.com/mchl-labs/stambecco.

Mehrbakhsh Nilashi, Rabab Ali Abumalloh, Behrouz Minaei-Bidgoli, Sarminah Samad, Muhammed Yousoof Ismail, Ashwaq Alhargan, and Waleed Abdu Zogaan. 2022. Predicting parkinson's disease progression: evaluation of ensemble methods in machine learning. *Journal of healthcare engineering,* 2022.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research,* 12:2825–2830.

Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277.*

Marco Polignano, Pierpaolo Basile, Marco de Gemmis, Giovanni Semeraro, and Valerio Basile. 2019. AlBERTo: Italian BERT Language Understanding Model for NLP Challenging Tasks Based on Tweets.

In *Proceedings of the Sixth Italian Conference on Computational Linguistics (CLiC-it 2019),* volume 2481. CEUR.

Andrea Santilli and Emanuele Rodolà. 2023. Camoscio: An italian instruction-tuned llama. *arXiv preprint arXiv:2307.16456.*

Gabriele Sarti and Malvina Nissim. 2022. It5: Large-scale text-to-text pretraining for italian language understanding and generation. *arXiv preprint arXiv:2203.03759.*

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.

Erdal Tasci, Caner Uluturk, and Aybars Ugur. 2021. A voting-based ensemble deep learning method focusing on image augmentation and preprocessing variations for tuberculosis detection. *Neural Computing and Applications,* 33(22):15541–15555.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971.*

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288.*

E. J. Wang. 2023. Alpaca-lora. https://github.com/tloen/alpaca-lora.

Xidong Wang, Nuo Chen, Junyin Chen, Yan Hu, Yidong Wang, Xiangbo Wu, Anningzhe Gao, Xiang Wan, Haizhou Li, and Benyou Wang. 2024. Apollo: Lightweight multilingual medical llms towards democratizing medical ai to 6b people. *arXiv preprint arXiv:2403.03640.*

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560.*

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652.*

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 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.

Han Yang, Mingchen Li, Yongkang Xiao, Huixue Zhou, Rui Zhang, and Qian Fang. 2023. One llm is not enough: Harnessing the power of ensemble learning for medical question answering. *medRxiv*, pages 2023–12.

Weipeng Zhou, Dmitriy Dligach, Majid Afshar, Yanjun Gao, and Timothy A Miller. 2023. Improving the transferability of clinical note section classification models with bert and large language model ensembles. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2023, page 125. NIH Public Access.

## A Implementation Details

The LlamaMTS model is trained with the LoRA Parameter-efficient Finetuning technique (Hu et al., 2022), using the Hugging Face Transformers and PEFT libraries (Wolf et al., 2020; Mangrulkar et al., 2022) and the Camoscio repository[10]. Specifically, our model is trained for 10 epochs on a desktop GPU Nvidia RTX 5000 Graphics Processing with 16GB of RAM, on a machine with Ubuntu 20.04.3 LTS. Training is implemented with batches of dimension 8 and gradient accumulation to obtain a final "virtual batch" of 128. The maximum length used for training is 2048 tokens. The learning rate is set to 3 x 10-4 with AdamW (Loshchilov and Hutter, 2018) and a total of 100 warmup steps are performed. We used a lora_r (i.e., the dimensionality of the low-rank update of the matrices) equal to 16. As base model, we merged the Camoscio adapter to the LLaMA 7 billion checkpoint[11]. In the evaluation we limited the max_new_tokens parameters to 1, forcing values to be binary through the prefix_allowed_tokens_fn parameter.

The BERT-based models are fine-tuned by using 10 epochs, 16 batches and a learning rate of 2 x 10-5.

---

[10]https://github.com/teelinsan/camoscio
[11]decapoda-research/llama-7b-hf

# Using Structured Health Information for
# Controlled Generation of Clinical Cases in French

**Hugo Boulanger**[*‡], **Nicolas Hiebel**[*†], **Olivier Ferret**[‡], **Karën Fort**[⋆], **Aurélie Névéol**[†]

[†]Université Paris Saclay, CNRS, LISN, France
[‡]Université Paris-Saclay, CEA, List, F-91120, Palaiseau, France
[⋆]Sorbonne Université / Université de Lorraine, CNRS, Inria, LORIA, France
[†]firstname.lastname@lisn.upsaclay.fr, [‡]firstname.lastname@cea.fr, [⋆]karen.fort@loria.fr

## Abstract

Text generation opens up new prospects for overcoming the lack of open corpora in fields such as healthcare, where data sharing is bound by confidentiality. In this study, we compare the performance of encoder-decoder and decoder-only language models for the controlled generation of clinical cases in French. To do so, we fine-tuned several pre-trained models on French clinical cases for each architecture and generate clinical cases conditioned by patient demographic information (gender and age) and clinical features. Our results suggest that encoder-decoder models are easier to control than decoder-only models, but more costly to train.

## 1 Introduction

The performance of current text generation models makes it difficult for humans to distinguish between natural and synthetic text (Casal and Kessler, 2023), paving the way for a wide range of applications including data augmentation and addressing resource sparsity (Claveau et al., 2021). In this article, we consider the case of reference documents that cannot be shared because of the personal information they contain but are sufficiently generic to mutualize processing resources on a community scale. One way of developing shared processes is to work with synthetic documents that are comparable in content and style to reference documents. We focus on electronic health records, though our methods can be applied to other fields with document-sharing constraints due to privacy.

Creating relevant synthetic documents is not trivial and must take several dimensions into account. As mentioned before, synthetic documents should be comparable to reference documents in terms of style, structure, and content, without leaking personal information that may be contained in the

training corpora. While directly identifying information can be subject to robust upstream de-identification, this does not make documents *anonymous* according to the definition of the General Data Protection Regulation (GDPR). Indeed, de-identification, whether automatic or manual, does not prevent cross-referencing medical information, which can particularly impact privacy for rare diseases.

It is possible to leverage the abilities of current text generation models to generate synthetic documents. However, such models are not as efficient when it comes to specialized domains such as the medical domain, even more so in languages other than English. Thus, the ability to precisely control the generation process is important both for medical consistency and for preserving the privacy of the information contained in real texts.

In this article, we propose a methodology for controlling text generation in terms of content. More specifically, the goal is to condition the generation of medical reports on patient profiles. Following the example of work carried out on the generation of synthetic patient profiles in terms of structured data (Walonoski et al., 2017), these profiles take the form of a set of medical concepts. This approach, which is part of a data-to-text generation problem, has the advantage over a textual priming approach of being able to finely control the information used for conditioning. The latter is implemented by training a neural language model with a set of pairs, each composed of a patient profile in the form of concepts and a reference report corresponding to this profile. Within this framework, the contributions of our paper are as follows:

- a method for controlling the content of medical report generation;

- a method for creating a training set for carrying out this control;

---

[*]These authors contributed equally to this work. The order is alphabetical.

- an implementation of the strategy using language models with two different architectures[1];

- an automatic multidimensional evaluation of synthetic text.

## 2 Related Work

### 2.1 Controlled text generation

Since the advent of the first large language models (LLMs) such as those of the GPT family (Radford et al., 2018), generating text resembling human production seems easy and the problem of generation has evolved to change focus: the aim is no longer simply to generate plausible text but to be able to control more finely what we generate. The texts produced by generative models may be irrelevant, offensive, or even dangerous (Bender et al., 2021). This is why a significant amount of work is being done on generation control. Control can concern several aspects of generation, such as the lexicon or text style (Zhang et al., 2023). Several control methods have been explored, including training a model with examples conditioned according to chosen criteria (Keskar et al., 2019) or modifying the probabilities of output tokens during inference (Kruszewski et al., 2023).

The *data-to-text* (Lin et al., 2023) approaches constrain generation from structured data (graphs, tables, and, in our case, *slots*). The preferred architectures are encoder-decoder models, which can have a variety of internal architectures, combining pre-trained models as encoders and/or decoders. It is also possible to directly fine-tune encoder-decoder models, such as the T5 model (Raffel et al., 2020). Causal language models, such as those using a Transformer (Vaswani et al., 2017) decoder architecture, use the context at the start of a sequence to generate the rest of the sequence.

### 2.2 Biomedical text generation

In the biomedical field, text generation is being explored either to facilitate the work of doctors or to address resource sparsity due to confidentiality issues. This work falls into the second category.

Earlier methods focus on training neural models from scratch. Melamud and Shivade (2019) train an LSTM to generate shareable clinical notes using differential privacy (Dwork et al., 2006),

and Ive et al. (2020) train a Transformer encoder-decoder model to generate synthetic mental health records conditioned by entities automatically extracted from real documents. However, training a model from scratch requires a substantial amount of data that is not available in languages other than English (Névéol et al., 2018).

Several efforts exploit Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) to generate structured data in the medical domain in English (Choi et al., 2017; Abedi et al., 2022; Torfi et al., 2022).

More recently, text generation is being explored to produce reports of discussions between doctors and patients, with the encoder-decoder architecture often being preferred (Eremeev et al., 2023; Ben Abacha et al., 2023; Asada and Miwa, 2023).

For French, Hiebel et al. (2023) fine-tune pre-trained auto-regressive language models to generate clinical cases with no particular constraints and propose a methodology for automatically evaluating the utility of the synthetic texts for a clinical entity recognition task.

## 3 Overall Method

As outlined in the introduction, we cast the task as a data-to-text generation problem where structured health data is used to shape the contents of synthetic text. Of course, finding the conditioning data within the generated texts cannot be the only criterion for evaluating the models: they would only need to reproduce their input to be judged as perfect. This conditioning must therefore be close in nature to the reference documents we wish to emulate.

As mentioned in section 2.1, this double conditioning can be achieved either by fine-tuning the language model used for generation with control elements or by steering the model during inference. We have opted for the former solution, as the latter implies applying elaborate text analysis processes during generation to check compliance with the conditioning, which is costly. The first solution, however, presupposes the availability of training data combining conditioning data and example texts conforming to this conditioning.

To this end, we have adopted a strategy comparable to Peng et al. (2018) for story generation, taken over by Ive et al. (2020) for medical reports, and consisting in automatically extracting the conditioning data from the example texts. This strategy

---

[1] `https://github.com/HugoBoulanger/ClinicalGenerator`

obviously presupposes the availability of text analysis tools capable of extracting this conditioning data from example texts with a sufficiently high level of performance. It therefore requires a close coupling between generation and analysis capabilities but eliminates the need for costly manual annotation. In the present case, we are focusing on medical concepts and are therefore dependent on models for extracting these concepts from medical reports but the genericity of this strategy means that new conditioning elements can easily be taken into account, as long as they can be automatically extracted from example texts.

## 4 Material and Methods

### 4.1 Clinical case corpora in French

The data used for our experiments come from two freely available clinical case corpora. The first corpus is the CAS corpus (Grabar et al., 2018), a corpus of de-identified clinical cases in French[2]. The second corpus is the E3C corpus (Magnini et al., 2020), a multilingual corpus of de-identified clinical cases. Our study only uses the clinical cases in French.

### 4.2 Defining constraints based on a patient profile

Our goal is to generate consistent clinical cases by controlling the generation using clinical elements. We have worked with clinicians to define the salient features of real clinical cases. These features are then used as constraints to generate text. Table 1 shows an example of features that were selected for a clinical case of the E3C corpus. These include patient demographics (age and gender), pathology location, histological information, various signs or symptoms, treatments and procedures performed, lab results, and scores (measures or codes). In line with clinicians' recommendations, we identify around twenty constraints per case, selecting if possible elements from each category with a majority of symptoms, treatments, and procedures. This approach ensures the selection of the salient information from the clinical cases, according to the doctors.

### 4.3 Extracting constraints from documents

Demographic data for the CAS corpus was directly taken from the existing corpus annotations for pa-

tient age and sex. We manually annotated the 1,009 cases from the E3C corpus to obtain equivalent demographic information for this corpus. Other clinical entities (e.g., signs and symptoms, procedures) were obtained by automatically annotating the two corpora consistently using clinical entity recognition models trained on the MERLOT private corpus (Campillos et al., 2018), which contains manual annotations for the entities of interest.

Constraint sets thus include manually annotated demographic information and automatically extracted clinical entities. For each document, we select age and gender when available. When the exact age is not provided, we use the age categories derived from the MeSH (Medical Subject Headings) thesaurus[3] check tags.

Clinical entities are selected from the MERLOT annotation categories that match the categories discussed with the doctors. For each clinical case, we select the ten procedures (*PROC*) and ten symptoms (*DISO*) with the highest tf.idf score. We also select substances (*CHEM*) and measures (*MEAS*). The latters are filtered to retain only informative measures (single digits such as *6* are annotated as *MEAS* but without additional information). Overall, we obtain an average of 26 constraints ($\pm 9.5$) per clinical case.

### 4.4 Text generation models

We compare the performance of two different architectures for the constrained generation of clinical texts using encoder-decoder vs. decoder-only pre-trained Transformer models.

**Encoder-decoder**   This architecture aims to generate text from structured data. In particular, fine-tuning the T5 model has become a standard method for data-to-text tasks. We chose to use the multilingual version of T5, called mT5 (Xue et al., 2021), with one billion parameters as a pre-trained model, and the Small (77 million parameters), Large (780 million parameters), and XL (3 billion parameters) versions of Flan-T5 (Chung et al., 2022) as models fine-tuned with instructions.

**Decoder only**   This architecture aims to generate text from textual prompts. We have chosen several models for this architecture. The Bloom (Scao et al., 2022) model, a generative model trained on several languages, and the Bloomz model, a variant

---

| Type of clinical feature | Sample value |
|---|---|
| Age | 54 |
| Sex | Masculin |
| Localisation | Vessie |
| Histology | adénocarcinome de l'ouraque peu différencié |
| Sign | hématurie |
| Procedure | scanner CT |
| Treatment | chimiothérapie par Méthotrexate-Vinblastine-Endoxan-Cisplatine |
| Score | T III A (selon la classification de Sheldon) |
| Bio | une négativité pour les cytokératines (ck) 7 et 20 |

Table 1: Sample control data based on manual analysis of a clinical case. The source case is shown in Table 2. We show in Appendix A.1 an English version based on the automatic translation of the document (Tables 5 and 6).

specially trained to perform different tasks (translation, automatic summarization, etc.). For each of these two models, we consider two versions in terms of size: one billion and seven billion parameters.

## 5 Experiments

### 5.1 Structured data representation

The use of these generative models requires the conversion of structured data into text format. We have chosen to linearize the inputs differently for the encoder-decoder models and the decoder-only models. For the encoder-decoder models, a special token representing the entity type is added before each entity. We separate demographic information (age, sex) from medical constraints (symptom, procedure, etc.) with a special token *contraintes* (*constraints*). For decoder-only models, no special tokens are used. Figure 1 shows an example of data representation for encoder-decoders.

### 5.2 Fine-tuning

The training set used to fine-tune our models comprises 1,424 clinical cases, containing over 500,000 tokens excluding constraints. For fine-tuning, we freeze the weights of the pre-trained model and add LoRA trainable matrices (Hu et al., 2022). The location of the trainable matrices depends on the type of model. For encoder-decoder models, we add LoRA matrices on the *queries* and *values* of the Transformer layers and the model head. For decoder-only models, LoRA matrices are added to the linear layers of the models. Special tokens are added to the embeddings via randomly initialized vectors. The processing of word embeddings varies according to two configurations defined as follows:

**"Frozen" configuration**: embeddings are frozen but we add LoRA matrices to enable adaptation to the task at a low memory cost.

**"Unfrozen" configuration**: the embeddings are unfrozen, to enable adaptation to the task, but at a higher cost.

We show the total number of parameters and the number of trainable parameters for each model in Table 7 in Appendix A.2.

### 5.3 Automatically generating clinical cases

Our test set consists of 156 clinical cases and their constraints. The constraints are given as input to the generative models and the real clinical cases are used as a reference when computing evaluation metrics. Decoding is performed using a beam search with five beams. We use sampling with a top-p of 0.9, a temperature of 1, and a repetition penalty of 3. Using sampling means that the same model might generate different texts from the same input. We run five generations for each test example to account for this variability.

### 5.4 Evaluation metrics

Automatic evaluation of text generation is notoriously difficult (Novikova et al., 2017). Numerous metrics exist to measure different aspects of text generation (Frisoni et al., 2022). Our metric selection aims to cover several dimensions of evaluation.

**Fit to constraints - *Accuracy*** This measure is used to assess the model's ability to implement the constraints. We calculate the proportion of constraints respected in generated texts in relation to the total number of constraints imposed.

**Language quality - *Perplexity*** Perplexity evaluates how well the textual data matches the probability distribution of a language model. We use a

```json
{
  "age": "22",
  "sexe": "masculin",
  "contraintes": [
    [
      "déhiscence cornéenne",
      "DISO"
    ],
    [
      "réparation chirurgicale",
      "PROC"
    ]
  ]
}
```

<age> 22 <sexe> masculin <contraintes>
<DISO> déhiscence cornéenne
<PROC> réparation chirurgicale

Figure 1: Example of data representation for encoder-decoder architecture (see Figure 2 in Appendix A.1 for its translation).

model specific to French, GPTFR (Simoulin and Crabbé, 2021). For this metric, we want the perplexity obtained on the generated data to be close to the perplexity obtained on the real data (equal to 19.5 for the training corpus).

**Diversity of generated texts - *Self-BLEU*** The Self-BLEU (Zhu et al., 2018) score is the average of the BLEU scores of all the sentences in a corpus. Thus, a redundant corpus will have a high Self-BLEU score while a varied corpus will have a lower score.

**Proximity to natural corpus - *Corpus-BLEU*** Corpus-BLEU (Yu et al., 2017) is a measure of proximity between two corpora and corresponds to the average BLEU score between each sentence in the generated corpus and all sentences in the natural corpus. We calculate Corpus-BLEU by comparing the clinical cases in the test corpus with the generated texts.

**Proximity with the clinical case corresponding to the constraints - *BLEU*** The BLEU (Papineni et al., 2002) score is calculated between the generated text and the actual clinical case from which the constraints originate. It measures proximity to real data in a more specific way than the Corpus-BLEU score.

## 6 Results

### 6.1 Evaluation of synthetic clinical cases

Table 2 shows examples of texts generated from a set of constraints by an encoder-decoder model (Flan-T5-XL frozen) and a decoder-only model (Bloomz 1b1 unfrozen). Table 3 shows the automatic evaluation of clinical cases generated with the different architectures studied. Among our baselines, the simple copy of the conditioning entities (Copy) obtains, as expected, an accuracy of 100 %, but also a very high perplexity. The Corpus

baseline corresponds to a copy of the test corpus in which we have removed the line breaks. This change explains why the BLEU and corpus-BLEU scores are not perfect and, more surprisingly, reduces perplexity from 30.5 to 19.5. The accuracy score, meanwhile, reveals the limitations of our data and accuracy calculation. The majority of these errors concern the sex of the patient, when this is not indicated by the gender agreement of the term "patient" or the use of the qualifier "male" or "female". Other errors are mainly due to rephrasing or errors in constraints.

The results show several trends. The first trend, which was expected but is confirmed by Table 3, is the positive correlation between the size of the models, both for encoder-only and encoder-decoder models, and their results: larger models obtain better results. When comparing encoder-decoder models of equal size (large), a model that has benefited from a training period with instructions, a Flan model, tends to obtain better overall results than a model pre-trained without instructions, especially for the unfrozen configuration. The Flan models also have the advantage of being fine-tuned more quickly for the same size, with a training period of 16 h for Flan-T5-large versus 60 h for mT5-large. As expected, the Flan-T5-XL models were the best-performing of the encoder-decoders tested. They generate more varied texts (Self-BLEU) and have the best accuracy. The texts generated most closely resemble the references (BLEU) and the perplexity values are better than those of the smaller versions of the model. It should be noted that mT5 models achieve lower perplexity —probably because the initial model is multilingual, whereas Flan-T5 models only saw French on translation tasks— and better Corpus-BLEU. Finally, Flan-T5 models are closer to the Corpus baseline than mT5 models in terms of perplexity, which was not *a priori* obvi-

| | |
|---|---|
| **Automatically extracted constraints** | âge: 54 ; sexe: masculin ; contraintes: hématurie isolée, examen tomodensitométrique, masse, 4 cm, adénocarcinome peu différencié, de type III, bilan d' extension, cystoprostatectomie radicale totale, lymphadénectomie iliaque, obturatrice, omphalectomie, entérocystoplastie de substitution, adénocarcinome de l'ouraque peu différencié, très localement mucosécrétant, ulcéré, carcinome transitionnel, grade III, Antigène Carcino-Embryonnaire, Leu-M1, CD 15, cytokératines, épithélium vésical, classification de Sheldon, Méthotrexate, Vinblastine, Endoxan, Cisplatine |
| **Real clinical case** | Un homme de 54 ans a consulté pour hématurie isolée. Une échographie, puis un examen tomod-ensitométrique, démontraient une masse de 4 cm de diamètre, au centre nécrotique, antérieure au dôme vésical, envahissant uniquement la graisse adjacente (Figure 1A.). Les biopsies, faites à la cystoscopie, révèlaient un adénocarcinome peu différencié, avec des zones d'épithélioma transitionnel de type III (Figure 1B). Le bilan d' extension était négatif. Le patient bénéficiait d'une cystoprostatectomie radicale totale avec lymphadénectomie iliaque et obturatrice, associée à une omphalectomie et une entérocystoplastie de substitution selon Studer. L'examen anatomo-pathologique complet de la pièce opératoire, révèlait un adénocarcinome de l'ouraque peu différencié, invasif, très localement mucosécrétant (Figure 1C). Il était ulcéré dans la lumière vésicale, et contenait des zones de carcinome transitionnel de grade III, ainsi que des cellules argentaffines. Deux poches de lymphonoeuds et la graisse environnante, étaient envahis. L'immunohistochimie montrait un positivité des cellules néoplasiques pour l'Antigène Carcino-Embryonnaire (ACE), Leu-M1 (CD 15), et une négativité pour les cytokératines (ck) 7 et 20. L'épithélium vésical était positif pour les ck 7. Il s'agissait d'un stade T III A (selon la classification de Sheldon). Le patient subissait une chimiothérapie par Méthotrexate - Vinblastine - Endoxan -Cisplatine par voie systémique. Avec 35 mois de recul, il va bien. |
| **Synthetic Bloomz 1b1 unfrozen** | Il s'agit d'un patient âgé de 54 ans sans antécédents pathologiques particuliers qui a consulté pour une hématurie isolée. L'examen clinique trouvait un patient en bon état général, apyrétique, avec des signes d'hypertension artérielle et une sensibilité abdominale à la palpation. La biologie était normale. L'échographie abdominale objectivait une masse tissulaire hétérogène du flanc gauche mesurant 4 cm de grand axe. L'IRM montrait une masse tissulaire homogène au niveau du flanc gauche (Figure 1). L'étude histologique de la pièce opératoire confirmait le diagnostic d'adénocarcinome peu différencié de type III selon la classification de Sheldon [1]. <br> Le bilan d'extension par résonance magnétique (RMN) ne retrouvait pas de métastase cérébrale ou pulmonaire. Une cystoprostatectomie radicale totale associée à une lymphadénectomie iliaque et obturatrice ainsi qu'une omphalectomie et entérocystoplastie de substitution ont été réalisées. L'étude anatomopathologique de la pièce opératoire avait confirmé le diagnostic d'adénocarcinome de l'ouraque peu différencié évoluant vers un carcinome transitionnel stade III selon la classification de Sheldon [1]. |
| **Synthetic Flan-T5-XL frozen** | Il s'agit d'un patient âgé de 54 ans, sans antécédents pathologiques particuliers, qui a présenté une hématurie isolée évoluant depuis 2 mois. L'examen tomodensitométrique a objectivé une masse de 4 cm au niveau de la paroi postérieure de l'ouraque (Figure 1). Le diagnostic d'adénocarcinome peu différencié de type III a été retenu. Le bilan d' extension était négatif. Une cystoprostatectomie radicale totale associée à une lymphadénectomie iliaque non obturatrice et une omphalectomie avec entérocystoplastie de substitution a été réalisée. L'étude anatomopathologique de la pièce opératoire a confirmé le diagnostic d'adénocarcinome de l'ouraque peu différencié comportant un aspect très localement mucosécrétant et ulcéré. Cet aspect est compatible avec un carcinome transitionnel de grade III (Antigène Carcino-Embryonnaire: CD34, Leu-M1 et CD 15). Les cellules tumorales exprimaient les cytokératines spécifiques de l'épithélium vésical (classification de Sheldon). Un traitement associant Méthotrexate, Vinblastine, Endoxan et Cisplatine a été débuté. |

Table 2: Sample clinical cases generated by two models based on constraints that were automatically extracted from a reference clinical case. An automatic English translation is shown in Table 6 in Appendix A.

ous since an instructed-based language model is not necessarily the best starting point for training a base text generator. This is particularly true for the Flan-T5-small models, without an evident explanation.

We observe that encoder-decoder models perform better than decoder-only models. Decoder-only models are also more unstable from one generation to another, with large standard deviations in Accuracy, Self-BLEU, and Corpus-BLEU, especially for the smallest models. In terms of perplexity, these models achieve lower scores and thus,

| | Generation method | Accuracy↑ | Perplexity | Self-BLEU-4↓ | Corpus-BLEU-4↑ | BLEU-4↑ |
|---|---|---|---|---|---|---|
| **Baselines** | Copying constraints | 100 | 194.3 | 14.4 | 25.5 | 1.1 |
| | Copying natural corpus | 98.8 | 19.5 | 33.4 | 97.4 | 97.5 |
| | Bloom 1b1 frozen∗ | s/o | 11.5±1.5 | 86.1±0.4 | 64.8±0.4 | s/o |
| | Bloom 1b1 unfrozen∗ | s/o | 10.2±0.9 | 82.9±0.4 | 60.6±0.5 | s/o |
| | Bloom 7b1 frozen∗ | s/o | 8.4±2.8 | 79.3±1.2 | 57.1±0.5 | s/o |
| **Encoder-decoder Encoder-decoder** | mT5-large frozen | 78.0±0.6 | 13.6±0.2 | 53.5±0.5 | 55.8±0.5 | 12.0±0.1 |
| | mT5-large unfrozen | 73.6±0.8 | 13.4±0.2 | 53.8±0.4 | 56.4±0.3 | 10.9±0.2 |
| | Flan-T5-small frozen | 61.6±0.4 | 18.9±0.4 | 49.1±0.4 | 47.1±0.4 | 6.6±0.1 |
| | Flan-T5-small unfrozen | 61.5±0.7 | 17.6±0.5 | 51.2±0.4 | 50.0±1.4 | 7.0±0.2 |
| | Flan-T5-large frozen | 81.5±1.1 | 14.8±0.4 | 52.8±0.4 | 55.3±0.4 | 12.0±0.1 |
| | Flan-T5-large unfrozen | 80.3±1.0 | 15.6±0.5 | 51.9±0.2 | 55.0±0.4 | 11.7±0.2 |
| | Flan-T5-XL frozen | 84.2±0.8 | 14.9±0.2 | 50.2±0.2 | 54.5±0.2 | 12.8±0.1 |
| | Flan-T5-XL unfrozen | 85.3±0.8 | 14.9±0.2 | 49.0±0.1 | 53.8±0.4 | 12.9±0.2 |
| **Decoder** | Bloom 1b1 frozen | 40.5±3.9 | 8.8±0.2 | 62.5±5.8 | 42.3±11.1 | 4.7±1.0 |
| | Bloom 1b1 unfrozen | 29.6±0.9 | 9.3±0.4 | 63.6±4.7 | 50.4±9.7 | 4.0±0.5 |
| | Bloom 7b1 frozen | 43.5±2.5 | 9.9±0.6 | 54.0±2.1 | 47.5±2.0 | 5.8±1.0 |
| | Bloomz 1b1 frozen | 45.4±4.2 | 9.2±0.2 | 61.9±7.6 | 41.8±11.0 | 5.2±1.3 |
| | Bloomz 1b1 unfrozen | 32.1±1.7 | 9.6±0.2 | 65.7±6.0 | 47.0±13.2 | 4.3±0.7 |
| | Bloomz 7b1 frozen | 39.8±3.0 | 9.9±0.2 | 55.0±1.9 | 49.8±1.5 | 5.4±0.4 |

Table 3: Evaluation of synthetic text generated from the constraints of the test set. Baseline models marked with "∗": training and generation without constraints.

deviate from the training corpus. As the model used to calculate perplexity is also a decoder, the common architecture potentially biases the decoders for this metric. On the other hand, decoder training time is much shorter: 10 to 15 minutes for billion-parameter models and 30 minutes for seven-billion-parameter models.

We can also identify some good practices regarding model pre-training and word embedding configuration. Models that have benefited from fine-tuning with instructions perform better overall than models with pre-training on a language modeling task. This is mainly true for accuracy and the BLEU score. We can assume that the type of instructions used for this fine-tuning – more precisely, whether these instructions are directly related or not to text generation tasks – may have an influence on the performance of these models but this analysis is beyond the scope of this article. We can also observe that frozen models perform better than unfrozen models. This observation could be considered surprising since the unfrozen models are supposed to have better adaptation capabilities but their heterogeneity in terms of parameters (LoRA and word embeddings matrices) is perhaps the source of these results.

## 6.2 Environmental impact

| Model | Fine-tuning | Generation | Perplexity | Total |
|---|---|---|---|---|
| mT5-large | 4.84 | 0.5 | 0.01 | 5.35 |
| flan-T5-small | 0.76 | 0.08 | 0.01 | 0.85 |
| flan-T5-large | 1.3 | 0.5 | 0.01 | 1.81 |
| flan-T5-XL | 4.84 | 0.5 | 0.01 | 5.35 |
| Bloom(z) 1b1 | 0.03 | 0.78 | 0.01 | 0.82 |
| Bloom(z) 7b1 | 0.05 | 0.64 | 0.01 | 0.70 |

Table 4: Environmental impact of the final experiments for each model, in kgCO$_2$e. Each line sums the emissions for different associated configurations. The total emissions reach 14.87 kgCO$_2$e.

Table 4 presents the greenhouse gas emissions of the experiments in terms of kgCO$_2$e. The environmental impact is essentially linked to the training of encoder-decoder models, which takes longer and requires more GPUs for larger models. These estimations were computed using the Machine-Learning Impact calculator presented in (Lacoste et al., 2019) with emission values for France (0.101 kgCO$_2$e/kWh) found in (Moro and Lonza, 2018).

## 7 Conclusion

In this study, we generate French clinical cases conditioned on structured clinical data. We com-

pare models with different architectures, encoder-decoder and decoder-only, which we fine-tune on a corpus of clinical cases using LoRA matrices. We propose an evaluation methodology based on a set of automatic measures: accuracy, perplexity, Self-BLEU, Corpus-BLEU, and BLEU. We observe that models with encoder-decoder architecture achieve better results on the task of generation from structured data, but with more costly training. Our experiments suggest that the best training strategy is to add LoRA matrices to the word embeddings rather than unfreezing them, although this does lengthen training.

The computing power available in a hospital setting limits the possibility of using larger and/or heavier models. The smallest size encoder-decoder model, Flan-T5-Small (77 million parameters), fits on the smaller Nvidia P6000 GPUs for fine-tuning and inference and obtains better performances than the larger decoder models. Small encoder-decoder models should be used if this type of resource is available for multiple hours. Decoders are more suitable if time on the GPUs is limited. However, it would be necessary to generate several candidates and filter them to compensate for the irregularity of these models.

Quantization might also be a solution for lightening computational loads, provided that quantized models achieve comparable results to their regular counterparts.

## 7.1 Limitations

The set of measures we have put in place gives us a fairly good view of what our models generate. There are, however, limits to using only accuracy, especially as calculated, to describe the fidelity of information transcription. Accuracy here seeks an exact match between the constraints and the text. Any reformulation of the model is therefore discarded, even though it may be correct. Moreover, using this measure alone does not give us any information on potential additions of information or entities by the models. In this study, we have exclusively used automatic metrics for the evaluation of generated texts. It is difficult to manually assess the quality of generated texts without clinical knowledge. Manual evaluation by clinical experts would enable us to estimate the medical consistency of generated texts more reliably. Finally, we have found that generations from the same model can be unstable. Filtering texts to keep the best candidate could improve results (Hiebel et al., 2023).

## 7.2 Ethical Considerations

The clinical documents used for fine-tuning the generation models (E3C and CAS) do not contain personal information. Thus, there is no additional risk of generating sensitive information with our models fine-tuned on those documents. The documents used for training clinical entity recognition models (MERLOT) were de-identified according to a protocol approved by the CNIL (*Commission de l'Informatique et des Libertés*), an independent French administrative regulatory body whose mission is to ensure that data privacy law is applied to the collection, storage, and use of personal data. In this work, we only use the models' annotations on the E3C and CAS corpus.

## Acknowledgments

## References

Masoud Abedi, Lars Hempel, Sina Sadeghi, and Toralf Kirsten. 2022. GAN-Based Approaches for Generating Structured Data in the Medical Domain. *Applied Sciences*, 12(14).

Masaki Asada and Makoto Miwa. 2023. BioNART: A biomedical non-AutoRegressive transformer for natural language generation. In *The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, pages 369–376, Toronto, Canada. Association for Computational Linguistics.

Asma Ben Abacha, Wen-wai Yim, Yadan Fan, and Thomas Lin. 2023. An empirical study of clinical note generation from doctor-patient encounters. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2291–2302, Dubrovnik, Croatia. Association for Computational Linguistics.

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.

Leonardo Campillos, Louise Deléger, Cyril Grouin, Thierry Hamon, Anne-Laure Ligozat, and Aurélie Névéol. 2018. A French clinical corpus with comprehensive semantic annotations: development of the Medical Entity and Relation LIMSI annOtated Text corpus (MERLOT). *Language Resources and Evaluation*, 52(2):571–601.

J. Elliott Casal and Matt Kessler. 2023. Can linguists distinguish between chatgpt/ai and human writing?: A study of research ethics and academic publishing. *Research Methods in Applied Linguistics*, 2(3):100068.

Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. 2017. Generating multi-label discrete patient records using generative adversarial networks. In *Proceedings of the 2nd Machine Learning for Healthcare Conference*, volume 68 of *Proceedings of Machine Learning Research*, pages 286–305. PMLR.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

Vincent Claveau, Antoine Chaffin, and Ewa Kijak. 2021. La génération de textes artificiels en substitution ou en complément de données d'apprentissage. In *TALN 2021 - 28e Conférence sur le Traitement Automatique des Langues Naturelles*, volume 1, pages 37–49, Lille, France. ATALA.

Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. 2006. Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology - EUROCRYPT 2006*, pages 486–503, Berlin, Heidelberg. Springer Berlin Heidelberg.

Maksim Eremeev, Ilya Valmianski, Xavier Amatriain, and Anitha Kannan. 2023. Injecting knowledge into language generation: a case study in auto-charting after-visit care instructions from medical dialogue. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2373–2390, Toronto, Canada. Association for Computational Linguistics.

Giacomo Frisoni, Antonella Carbonaro, Gianluca Moro, Andrea Zammarchi, and Marco Avagnano. 2022. NLG-metricverse: An end-to-end library for evaluating natural language generation. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3465–3479, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.

Natalia Grabar, Vincent Claveau, and Clément Dalloux. 2018. CAS: French corpus with clinical cases. In *Proceedings of the Ninth International Workshop on Health Text Mining and Information Analysis*, pages 122–128, Brussels, Belgium. Association for Computational Linguistics.

Nicolas Hiebel, Olivier Ferret, Karen Fort, and Aurélie Névéol. 2023. Can synthetic text help clinical named entity recognition? a study of electronic health records in French. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2320–2338, Dubrovnik, Croatia. Association for Computational Linguistics.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, Online.

Julia Ive, Natalia Viani, Joyce Kam, Lucia Yin, Somain Verma, Stephen Puntis, Rudolf Cardinal, Angus Roberts, Robert Stewart, and Sumithra Velupillai. 2020. Generation and evaluation of artificial mental health records for natural language processing. *npj Digital Medicine*, 3.

Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: A conditional transformer language model for controllable generation. *CoRR*, abs/1909.05858.

Germán Kruszewski, Jos Rozen, and Marc Dymetman. 2023. disco: a toolkit for distributional control of generative models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 144–160, Toronto, Canada. Association for Computational Linguistics.

Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. 2019. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*.

Yupian Lin, Tong Ruan, Jingping Liu, and Haofen Wang. 2023. A survey on neural data-to-text generation. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–20.

Bernardo Magnini, Begoña Altuna, Alberto Lavelli, Manuela Speranza, and Roberto Zanoli. 2020. The E3C Project: Collection and Annotation of a Multilingual Corpus of Clinical Cases. In *Proceedings of the Seventh Italian Conference on Computational Linguistics, CLiC-it 2020*, volume 2769 of *CEUR Workshop Proceedings*, Bologna, Italy. CEUR-WS.org.

Oren Melamud and Chaitanya Shivade. 2019. Towards Automatic Generation of Shareable Synthetic Clinical Notes Using Neural Language Models. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 35–45, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Alberto Moro and Laura Lonza. 2018. Electricity carbon intensity in european member states: Impacts on ghg emissions of electric vehicles. *Transportation Research Part D: Transport and Environment*, 64:5–14. The contribution of electric vehicles to environmental challenges in transport. WCTRS conference in summer.

Aurélie Névéol, Hercules Dalianis, Sumithra Velupillai, Guergana Savova, and Pierre Zweigenbaum. 2018. Clinical natural language processing in languages other than english: opportunities and challenges. *Journal of Biomedical Semantics*, 9(1):12.

Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. 2017. Why we need new evaluation metrics for NLG. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2241–2252, Copenhagen, Denmark. Association for Computational Linguistics.

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

Nanyun Peng, Marjan Ghazvininejad, Jonathan May, and Kevin Knight. 2018. Towards controllable story generation. In *Proceedings of the First Workshop on Storytelling*, pages 43–49, New Orleans, Louisiana. Association for Computational Linguistics.

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

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. *J. Mach. Learn. Res.*, 21(1).

Teven Le Scao et al. 2022. Bloom: A 176b-parameter open-access multilingual language model.

Antoine Simoulin and Benoit Crabbé. 2021. Un modèle Transformer Génératif Pré-entrainé pour le _____ français. In *Traitement Automatique des Langues Naturelles*, pages 246–255, Lille, France. ATALA.

Amirsina Torfi, Edward A. Fox, and Chandan K. Reddy. 2022. Differentially private synthetic medical data generation using convolutional GANs. *Information Sciences*, 586:485–500.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Jason Walonoski, Mark Kramer, Joseph Nichols, Andre Quina, Chris Moesel, Dylan Hall, Carlton Duffett, Kudakwashe Dube, Thomas Gallagher, and Scott McLachlan. 2017. Synthea: An approach, method, and software mechanism for generating synthetic patients and the synthetic electronic health care record. *Journal of the American Medical Informatics Association*, 25(3):230–238.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: sequence generative adversarial nets with policy gradient. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI'17, page 2852–2858. AAAI Press.

Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023. A survey of controllable text generation using transformer-based pre-trained language models. volume 56, New York, NY, USA. Association for Computing Machinery.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '18, page 1097–1100, New York, NY, USA. Association for Computing Machinery.

# A  Appendix

## A.1  Translation of Tables and Figures

Figure 2 presents the translation of the example of data representation shown in Figure 1.

Tables 5 and 6 present an automatic translation of the natural document with the corresponding constraints and generated samples that were presented in Tables 1 and 6. The automatic translation was done with DeepL[4].

## A.2  Model Sizes

Table 7 present the total number of parameters and the trainable parameters for each model.

---

[4] www.deepl.com

```
"age": "22",
"sex": "male",
"constraints": [
  [
    "corneal dehiscence",
    "DISO"
  ],
  [
    "surgical repair",
    "PROC"
  ]
]
```

<age> 22 <sex> male <constraints>
<DISO> corneal dehiscence
<PROC> surgical repair

Figure 2: Example of data representation for encoder-decoder architecture (translation of Figure 1).

| Type of clinical feature | Sample value |
| --- | --- |
| Age | 54 |
| Sex | Male |
| Localisation | Bladder |
| Histology | poorly differentiated adenocarcinoma of the urachus |
| Sign | hematuria |
| Procedure | CT scan |
| Treatment | methotrexate-vinblastine-endoxan-cisplatin chemotherapy |
| Score | T III A (according to Sheldon's classification) |
| Bio | negative for cytokeratins (ck) 7 and 20 |

Table 5: Sample control data based on manual analysis of a clinical case (translation of Table 1).

| | |
|---|---|
| **Automatically extracted constraints** | age: 54 ; sex: male ; constraints: isolated hematuria, CT scan, mass, 4 cm, poorly differentiated adenocarcinoma, of type III, extension work-up, total radical cystoprostatectomy, iliac and obturator lymphadenectomy, omphalectomy, replacement enterocystoplasty, adenocarcinoma of the urachus, very locally mucosecretory, ulcerated, transitional cell carcinoma, grade III, Carcinoembryonic Antigen, Leu-M1, CD 15, cytokeratins, bladder epithelium, Sheldon's classification, methotrexate, vinblastine, endoxan, cisplatin |
| **Real clinical case** | A 54-year-old man presented with isolated hematuria. An ultrasound scan, followed by a CT scan, showed a mass 4 cm in diameter, with a necrotic center, anterior to the bladder dome, invading only the adjacent fat (Figure 1A.). Cystoscopic biopsies revealed a poorly differentiated adenocarcinoma, with areas of type III transitional epithelioma (Figure 1B). The extension work-up was negative. The patient underwent total radical cystoprostatectomy with iliac and obturator lymphadenectomy, combined with omphalectomy and Studer replacement enterocystoplasty. Complete anatomopathological examination of the surgical specimen revealed a poorly differentiated, invasive, very locally mucosecretory adenocarcinoma of the urachus (Figure 1C). It was ulcerated in the bladder lumen, and contained areas of grade III transitional cell carcinoma, as well as argentaffin cells. Two lymph node pockets and surrounding fat were invaded. Immunohistochemistry showed neoplastic cells to be positive for Carcinoembryonic Antigen (CEA), Leu-M1 (CD 15), and negative for cytokeratins (ck) 7 and 20. The bladder epithelium was positive for ck 7. This was stage T III A (according to Sheldon's classification). The patient was undergoing systemic methotrexate-vinblastine-endoxan-cisplatin chemotherapy. With 35 months' follow-up, he is doing well. |
| **Synthetic Bloomz 1b1 unfrozen** | This is a 54-year-old patient with no particular pathological history who consulted for isolated hematuria. Clinical examination revealed a patient in good general condition, apyretic, with signs of arterial hypertension and abdominal tenderness to palpation. Laboratory tests were normal. Abdominal ultrasound revealed a heterogeneous tissue mass on the left flank measuring 4 cm in long axis. MRI showed a homogeneous tissue mass in the left flank (Figure 1). Histological examination of the operative specimen confirmed the diagnosis of type III poorly differentiated adenocarcinoma according to Sheldon's classification [1]. Magnetic resonance imaging (MRI) did not reveal any brain or lung metastases. A total radical cystoprostatectomy associated with an iliac and obturator lymphadenectomy and a substitution omphalectomy and enterocystoplasty were performed. Anatomopathological study of the surgical specimen confirmed the diagnosis of poorly differentiated adenocarcinoma of the urachus, progressing to stage III transitional cell carcinoma according to Sheldon's classification [1]. |
| **Synthetic Flan-T5-XL frozen** | This was a 54-year-old patient with no specific pathological history who presented with isolated hematuria evolving for 2 months. CT examination revealed a 4 cm mass in the posterior wall of the myringotomy (Figure 1). A diagnosis of type III poorly differentiated adenocarcinoma was made. The extension work-up was negative. Total radical cystoprostatectomy combined with non-obturator iliac lymphadenectomy and omphalectomy with replacement enterocystoplasty was performed. Anatomopathological study of the surgical specimen confirmed the diagnosis of poorly differentiated adenocarcinoma of the urachus, with a very locally mucosecretory and ulcerated appearance. This appearance is compatible with a grade III transitional cell carcinoma (Carcinoembryonic Antigen: CD34, Leu-M1 and CD 15). Tumor cells expressed cytokeratins specific to the bladder epithelium (Sheldon classification). Treatment with Methotrexate, Vinblastine, Endoxan and Cisplatine was initiated. |

Table 6: Sample clinical cases generated by two models based on constraints that were automatically extracted from a reference clinical case (translation of Table 2).

| Model | Total parameters | Trainable parameters | Percentage trainable |
|---|---|---|---|
| mT5-large frozen | 1.2 B | 9.6 M | 0.8% |
| mT5-large unfrozen | 1.5 B | 518 M | 34.5% |
| Flan-T5-small frozen | 78.3 M | 1.3 M | 1.7% |
| Flan-T5-small unfrozen | 94.3 M | 33.8 M | 35.8% |
| Flan-T5-large frozen | 787 M | 4.3 M | 0.5% |
| Flan-T5-large unfrozen | 819 M | 69.7 M | 8.5% |
| Flan-T5-XL frozen | 2.9 B | 7.9 M | 0.3% |
| Flan-T5-XL unfrozen | 2.9 B | 139 M | 4.7% |
| Bloom(z) 1b1 frozen | 1.1 B | 6.7 M | 0.6% |
| Bloom(z) 1b1 unfrozen | 1.5 B | 390 M | 26.8% |
| Bloom(z) 7b1 frozen | 7.1 B | 17.8 M | 0.3% |

Table 7: Parameter count as reported by the PEFT library used for fine-tuning. We report the same numbers for Bloom and Bloomz because the models have the same architecture and the same amount of parameters. Shift of total parameters in unfrozen models are due to tied embeddings being counted twice.

# Large Language Models Provide Human-Level Medical Text Snippet Labeling

**Ibtihel Amara**[1,2*], **Haiyang Yu**[2], **Fan Zhang**[2], **Yuchen Liu**[2],
**Benny Li**[2], **Chang Liu**[2], **Rupesh Kartha**[2], and **Akshay Goel**[2]
[1] McGill University and [2] Google Research

## Abstract

This study evaluates the proficiency of Large Language Models (LLMs) in accurately labeling clinical document excerpts. Our focus is on the assignment of potential or confirmed diagnoses and medical procedures to snippets of medical text sourced from unstructured clinical patient records. We explore how the performance of LLMs compare against human annotators in classifying these excerpts. Employing a few-shot, chain-of-thought prompting approach with the MIMIC-III dataset, Med-PaLM 2 showcases annotation accuracy comparable to human annotators, achieving a notable precision rate of approximately 92% relative to the gold standard labels established by human experts.

## 1 Introduction

Advanced natural language processing (NLP) tools especially generative language models have recently made a big difference in healthcare (Liu et al., 2023; Hu et al., 2023; Singhal et al., 2023; Goel et al., 2023; Tu et al., 2024). One key way NLP is used is to find important medical details, like diagnoses, within a patient's unstructured data. Clinicians can quickly search for medical conditions in these documents, speeding up their understanding of a patient's medical history.

In this work, we focus on identifying both potential and confirmed medical conditions throughout the various text snippets of information found in patients' medical records. Particularly, we establish a *mapping between a large comprehensive list of possible medical condition or procedures queries C and text snippets from clinical documents S.* We visualize the core task in Figure 4 in the Appendix. When establishing a connection between a medical condition or procedure and a snippet of medical



Figure 1: **Labeling Framework.** This consists of four components: (1) **Pre-filtering**; the query list is pre-filtered using a keyword search algorithm. (2) **Text Chunking**; the medical note is divided into smaller text snippets. (3) **Alignment**; The remaining queries are associated with the most relevant text snippets. (4) **LLM Labeling**; the text snippets and queries are sent to a large language model (LLM). The LLM confirms which conditions are truly relevant for each snippet.

text, we do not expect the text to include "supporting" components that are directly related to the condition or procedure. Instead, we anticipate that the labeler (here LLM) recognizes significant medical patterns, medications, and symptoms that point to a potential diagnosis (i.e. medical condition) or medical procedure. A straightforward example of this is as follows:

**Text Snippet:** *"The patient has been taking metformin 2500mg a day since last year."*
**Possible LLM Condition/Procedure Labeling:** *Diabetes and Polycystic Ovary Syndrome (PCOS).*
The rationale behind this labeling is that metformin is a medication commonly used in various medical treatments. Mastering this labeling process contributes to building the foundation for powerful information retrieval, search and summarization systems, which has the potential to revolutionize medical search and ultimately improve healthcare workflow. We summarize our main contributions as follows: (1) We demonstrate that LLMs can be used to identify potential labels (i.e medical conditions or procedures) with medical snippets reducing reliance on human experts. (2) We propose a cost-effective and efficient labeling framework with LLMs, which accelerates the annotation process

---

by reducing expensive LLM calls while preserving high labeling quality.

## 2 Related Work

Our work aligns with the field of Named Entity Recognition (NER) (Doan et al., 2012; Mullenbach et al., 2018; Yang et al., 2019; Goel et al., 2023; Guo et al., 2024; Ferraro et al., 2024). While NER primarily focuses on identifying and categorizing words into predefined entities such as procedure codes, medication codes, organizations, and others, our work takes a different approach. We adopt a unique methodology wherein we meticulously structure clinical documents by segmenting them into coherent and meaningful snippets. Our goal is to establish connections between these snippets and pertinent medical conditions or procedures drawn from a comprehensive list of medical queries. This approach allows us to not only identify potential medical conditions or procedures but also understand the context within the document, which ultimately will be useful for building and training medical search and retrieval systems.

## 3 Methodology

We provide in Figure 1 the general framework of our proposed labeling pipeline.

**Pre-filtering.** The first step in the pipeline involves pre-filtering a comprehensive list using cost-effective filtering strategies. This step aims to reduce the number of expensive calls to the LLM and avoid quality label loss (see Appendix I.1). There are several methods for implementing a pre-filtering step, such as embedding similarity, medical search engines, etc. We encourage researchers to explore other available and easy alternatives. In this work, we employed a keyword search algorithm. This technique expands the input queries (through query expansion) and looks for the matched text in the input document, which we regard as reference snippets. More details can be found in Appendix B.

**Text Chunking.** We broke down the patient's medical record into more manageable and informative text segments (i.e. medical snippets). We performed different chunking strategies (see Appendix D), and settled with a hybrid method involving a sentence-based (3-4 sentences) chunking algorithm with a constraint of 10-70 word tokens (Figure E).

**Alignment.** At this stage, we matched the remaining medical conditions and procedures to the corresponding text snippet. In particular, we opted

for fuzzy matching. This can be considered as a secondary pre-filtering step at the snippet level. In our work, since our pre-filtering step outputs a reference snippet per condition or procedure, we attempted to locate these snippets within the different text chunks we have produced. This way, the condition becomes associated with the chunked snippet.[*]

**LLM labeling.** In the final stage of our framework, we paired the text snippets and their corresponding medical conditions. These pairs are then sent to the LLM using appropriate prompting strategies. The LLM assesses the relevance of the text snippet and medical condition in each pair. If it determines a condition to be relevant, the condition label is included as one of the final labels for that snippet.

## 4 Experimental Setup

**Dataset and Pre-processing.** We used the publicly available de-identified dataset MIMIC-III (Johnson et al., 2016). It is a collection of de-identified medical records and notes of more than 40,000 critical care patients at a large tertiary care hospital. It contains over two million unstructured clinical documents from nurses, physicians, etc. In our work, we randomly sampled 1000 patients and fetched all of their corresponding clinical records. Our pre-processing of the dataset was kept simplistic. We used simple regular expressions to identify formatting inconsistencies, such as extra spaces or tabs, in the clinical documents. We provide basic statistics in Section F about the sampled subset from the MIMIC-III dataset.

**Human Labeling Workflow.** The human labeling process was carried out in three separate rounds. In each round, a different group of medical expert raters was recruited to evaluate a distinct set of medical text snippets paired with a condition. Overall, we had 14 different medical experts as human annotators: 3 experts on the first round, 5 on the second, and 6 on the third round of labeling. Specifically, the raters were given a set of multiple-choice options ("Relevant", "Irrelevant", and "Not sure") and were asked to answer the following question: "Is the following text snippet relevant to the following medical condition/procedure?". The raters were given a random sample of snippets. In total, we collected 14,470 labeled snippet-condition pairs.

---

[*]It is important to note that the inclusion of this component is contingent upon the pre-filtering strategy that is ultimately adopted.

| | Zero-shot | | | | | | Zero-shot CoT | | | | | | Few-shot CoT | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NE | | | WE | | | NE | | | WE | | | NE | | | WE | | |
| LLM Architecture | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| PaLM 2* | **83.55** | 96.55 | **89.58** | **71.79** | 93.37 | **81.17** | 50.49 | **91.07** | 64.97 | 39.82 | **87.25** | 54.68 | 83.75 | **96.28** | 89.58 | 73.96 | **94.25** | 82.88 |
| Med-PaLM 2 | 81.10 | **98.51** | 88.95 | 67.37 | **96.92** | 79.49 | **91.74** | 84.37 | **87.91** | **81.49** | 76.48 | **78.91** | **92.70** | 87.62 | **90.09** | **84.94** | 82.98 | **83.95** |

Table 1: **LLM performance compared to the golden human labels.** Med-PaLM 2 has the highest performance overall. NE: golden labels without rater exclusion; WE: golden labels with rater exclusion. P: precision; R: recall; F1: F1 score. *This is a fine-tuned PaLM 2 model variant for programming tasks.

**LLM Labeling and Prompt Engineering Strategies.** In this study we investigated LLM capabilities using simple prompt engineering techniques to more complex reasoning prompting strategies. We used zero-shot, few-shot, chain-of-thought (CoT) (Wei et al., 2022) , self-consistency CoT (Wang et al., 2022), and chain of verification (CoVe) (Dhuliawala et al., 2023). We assess these strategies on providing accurate labeling on medical snippets with respect to the "golden" labels obtained from human annotators. As for the LLM architectures, we used two different models: PaLM 2 (Anil et al., 2023) and Med-PaLM 2 (Singhal et al., 2023).

## 5 Results

We provide details about the basic statistics on both human labeled data and the sampled data from MIMIC-III in Appendix F.

**Agreement Between Human Raters.** Before relying on human labels, it is essential to assess their reliability and validity, especially when there is no clear or accessible ground truth label. To do this, we start by plotting the response distribution of each rater at each labeling round. Figure 2 exhibits significant variations within the different raters' responses. In round 1, for instance, two raters (raters 1 and 2) demonstrated a tendency to provide answers skewed towards the "irrelevant" category. In contrast, rater 3 maintained a balanced approach, assigning an equal number of responses to both the "irrelevant" and "relevant" categories. During the second round of the labeling process, raters 5, 6, and 8 exhibited a similar pattern of providing more "irrelevant" labels. In contrast, raters 4 and 7 produced more "relevant" responses. In the third round, we observe a similar distribution trend, which is predominantly characterized by a skew towards the "irrelevant" side. In Figure 3, we assess inter-rater reliability using Cohen's Kappa statistics (Viera et al., 2005; McHugh, 2012) and we provide in Appendix H the agreement interpretations. We observe that the level of agreement between raters varies across different rounds. In round 1 of labeling, the agreement ranges from "fair" to "moderate," indicating a practical level of consensus. However, in rounds 2 and 3, substantial variations emerge. In round 2, raters 5 and 6 exhibit a stronger agreement compared to other raters. In the third round, we observe a notable agreement between raters 11 and 12 and a moderate agreement between raters 10 and 11.

**Golden Labels.** Based on these reliability and agreement results, we decide to create *two types of golden labels*: (1) Majority vote with no rater exclusion [NE] and (2) Majority vote with rater exclusion [WE]. Indeed, for the first case, we mainly consider all of the raters' responses. As for the second version of golden labels, we consider only the majority voting of rater responses that are at least in a fair agreement with each other. In this case, we consider the following raters in each of the rounds (i.e. all raters in round 1, raters 5, 6, and 8 in round 2, and raters 10, 11, and 12 for round 3). We also applied a rigorous majority voting strategy. This involved selecting cases where there was a clear and consistent consensus among the raters. For instance, for a particular snippet-condition pair, we designated the snippet as relevant (associating it with the condition) only if all raters agreed that the condition was pertinent to the snippet. In cases where raters disagreed, we deemed the condition as "not sure", and excluded it from the evaluation.

**LLM Performance on the Aggregated raters' labels.** In Table 1, we compare the performance of different LLMs using different prompting strategies. Overall, Med-PaLM 2 achieves the highest precision across the different LLM architectures for each prompting strategy. This is likely because Med-PaLM 2 is specifically trained on medical text, which allows it to provide more precise results. However, when considering the recall metric, PaLM 2 achieves highest recall values, albeit with lower precision. When building a dataset for training medical retrieval systems, it is well preferred to have a good balance between precision and recall. Among the various prompting techniques, we

Figure 2: **Response Distribution of each Raters.** There are clear variations in the distribution of annotations across the raters.



Figure 3: **Cohen Kappa's Inter-rater reliability.**

find that a few-shot CoT approach yields superior overall performance. Specifically, we observe improvements in both precision and recall metrics. Med-PaLM 2 outperformed in Few-shot CoT due to its medical focus. In zero-shot settings, PaLM 2 achieved top precision and F1 scores.

**Beyond Basic Prompts.** In addition to the three aforementioned prompts, we also explored two more prompting strategies and tested them on Med-PaLM 2: (1) self-consistency CoT and (2) Chain of Verification (CoVe). The accuracy of all these 5 prompts are shown in the Table 2. Note that the self-consistency prompting is based on the Few-shot CoT prompt with multiple runs using non-zero temperature (T=0.5). Although the ensemble result slightly outperforms the single run with T=0 (few-shot CoT), it requires multiple runs (three in our case), which substantially increases the time expenditure, hence we used the few-shot CoT for our final labeling task. Similarly, utilizing the CoVe prompt entails multiple rounds of verification to attain the final label. Each round demands distinct LLM invocations, rendering this method expensive.

**Time Efficiency Comparison.** On average, human raters took anywhere between 65 and 595 seconds (approximately 10 minutes) to review a single snippet, with an average time of 203 seconds. Considering an average of 8 conditions per snippet,

| Prompts | P | R | Acc. | F1 |
|---|---|---|---|---|
| Zero-shot | 78.73 | 97.30 | 89.97 | 87.03 |
| Zero-shot CoT | 92.79 | 72.59 | 88.57 | 81.45 |
| Few-shot CoT | 91.94 | 83.45 | 91.75 | <u>87.49</u> |
| Self-Consistency | 92.63 | 84.43 | 92.29 | **88.34** |
| CoVe | 73.98 | 77.67 | 82.83 | 75.78 |

Table 2: **Med-PaLM 2 performance on the NE dataset.** The highest F1 score is highlighted in bold, and the second-best score is underlined. Self-consistency yields the best performance. However, given that the few-shot prompt is less expensive than the self-consistency prompt, it is still a viable option.

this translates to roughly 24 seconds to review a snippet-condition pair. The latency of LLMs, on the other hand, varies depending on factors such as model architecture, size, inference infrastructure, and prompt strategies. However, on average, their latency is significantly lower than that of human raters.

## 6   Conclusion

We proposed a framework for labeling clinical notes. Our findings suggest that LLMs can produce high-quality medical data labels, which can serve as a valuable dataset for NLP tasks, such as information retrieval systems. These systems can help clinicians to be more efficient in their daily workflow by finding the key information faster and focus on pertinent facts within a clinical note.

# 7 Limitations

This work focused on a specific task: labeling medical conditions within clinical text snippets. While successful in this context, generalizing this approach to other scenarios might face limitations. Our keyword search method could miss relevant conditions not captured by the search algorithm. Additionally, the large language model (LLM) labeling is sensitive to the way it is prompted and requires further exploration to find optimal strategies for different use cases. Furthermore, the sentence-based chunking algorithm, while effective here, is specifically designed for the MIMIC-III dataset and may need adjustments for broader application. Finally, even human raters showed significant disagreement on labeling, highlighting the challenges posed by limited context in snippets and the inherent uncertainties within the medical domain, particularly when associating conditions with diverse symptoms. These limitations underscore the need for further research to improve generalizability and robustness when applying this type of system to broader medical text analysis tasks.

# 8 Ethical Statement

Labels created by LLMs might reflect biases inherent in the LLMs themselves. To some extent, these biases can be reduced by diversifying the LLMs, as this approach encourages the generation of more robust labels. However, even after implementing this strategy, biases may still persist. In the medical context specifically, additional alignment intervention methods can be utilized to modify the behavior of the LLM, presenting a potential solution to this challenge.

# References

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:2309.11495*.

Son Doan, Nigel Collier, Hua Xu, Pham Hoang Duy, and Tu Minh Phuong. 2012. Recognition of medication information from discharge summaries using ensembles of classifiers. *BMC medical informatics and decision making*, 12:1–10.

Antonino Ferraro, Antonio Galli, Valerio La Gatta, Mario Minocchi, Vincenzo Moscato, and Marco Postiglione. 2024. Few shot ner on augmented unstructured text from cardiology records. In *International Conference on Emerging Internet, Data & Web Technologies*, pages 1–12. Springer.

Akshay Goel, Almog Gueta, Omry Gilon, Chang Liu, Sofia Erell, Lan Huong Nguyen, Xiaohong Hao, Bolous Jaber, Shashir Reddy, Rupesh Kartha, et al. 2023. Llms accelerate annotation for medical information extraction. In *Machine Learning for Health (ML4H)*, pages 82–100. PMLR.

Yuting Guo, Yao Ge, and Abeed Sarker. 2024. Detection of medication mentions and medication change events in clinical notes using transformer-based models. *Studies in Health Technology and Informatics*, 310:685–689.

Yan Hu, Iqra Ameer, Xu Zuo, Xueqing Peng, Yujia Zhou, Zehan Li, Yiming Li, Jianfu Li, Xiaoqian Jiang, and Hua Xu. 2023. Zero-shot clinical entity recognition using chatgpt. *arXiv preprint arXiv:2303.16416*.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Liwei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

Zhengliang Liu, Xiaowei Yu, Lu Zhang, Zihao Wu, Chao Cao, Haixing Dai, Lin Zhao, Wei Liu, Dinggang Shen, Quanzheng Li, et al. 2023. Deid-gpt: Zero-shot medical text de-identification by gpt-4. *arXiv preprint arXiv:2303.11032*.

Mary L McHugh. 2012. Interrater reliability: the kappa statistic. *Biochemia medica*, 22(3):276–282.

James Mullenbach, Sarah Wiegreffe, Jon Duke, Jimeng Sun, and Jacob Eisenstein. 2018. Explainable prediction of medical codes from clinical text. *arXiv preprint arXiv:1802.05695*.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*.

Tao Tu, Anil Palepu, Mike Schaekermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, et al. 2024. Towards conversational diagnostic ai. *arXiv preprint arXiv:2401.05654*.

Anthony J Viera, Joanne M Garrett, et al. 2005. Understanding interobserver agreement: the kappa statistic. *Fam med*, 37(5):360–363.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.

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.

Xi Yang, Jiang Bian, Yan Gong, William R Hogan, and Yonghui Wu. 2019. Madex: a system for detecting medications, adverse drug events, and their relations from clinical notes. *Drug safety*, 42:123–133.

Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, et al. 2024. Bench: Extending long context evaluation beyond 100k tokens. *arXiv preprint arXiv:2402.13718*.

## A Visualization of Medical Note Labeling Task.

The goal is to categorize and classify each medical note text snippet into potential conditions. This pairing of text snippets and conditions can be highly valuable for training dense retrieval systems for medical notes pertaining to specific patients.



Figure 4: Medical Note Labeling Task

## B An Example of Input and Output of the Implemented Mixer Search Algorithm for a Medical Note.

We used a keyword mixer search algorithm. This technique expands the input queries (via query expansion) and identifies their connections and locations within the input document. By positioning the keywords in the input document (25 tokens as the context with the searched keyword in the center), the algorithm generates reference sentences. Ultimately, the most representative reference sentence is given in relation to the input query (i.e. medical conditions/procedures). We illustrate the behavior of the mixer search algorithm as a technique for pre-filtering unlikely conditions from a medical note. Given a single query condition and the patient's clinical note, the algorithm identifies the most relevant text snippet from the document that is likely to be associated with the condition.

> **Input:**
> query: "coughing"
> note: (note_id, the medical text)
> **Output: reference text from the medical note**
> "... Secretions: produced bloody and yellowish sputum with productive **cough** which was cleared with Yankauer and tracheal suction. Also of note ..."

## C Identifying Sentence Boundaries

To identify sentence boundaries in the medical notes within MIMIC-III, we use regular expressions after some simple pre-processing as described in the experimental setup section. Regular expressions provide a flexible and efficient way to capture full sentences. They allow us to define patterns that match specific sentence-ending punctuation marks, such as periods (.), exclamation marks (!), and question marks (?). Additionally, regular expressions can be used to handle more complex cases, such as sentences that end with abbreviations or quotations.

## D Note Chunking/Segmentation Strategies.

We implemented several ways of text chunking to split each medical note properly:

**(1) Sentence-base (SB) segmentation:** The medical note is fragmented according to a collection of one or more sentences. We divide the document into $n$ non-overlapping sentences without regard for the notes' structure and sectioning.

**(2) Word-base (WB) segmentation:** The medical note is fragmented according to a collection of one or more word tokens. One thing to note is that this word base does not consider the cut offs. In other words, it would take the number of words given in the input regardless of it being an incomplete or full sentence. For our use case, snippets would be more readable (contextually and grammatically correct) for later human and LLM labeling.

**(3) Sentence-word fusion (SWF) segmentation:** One major thing we noticed during the execution of our algorithms is that we were getting a lot of very short sentences. To mitigate this, we implemented a hybrid version of the snipping algorithms above. We considered a sentence-based text segmentation, with a constraint on the number of words admissible for each segment via a range threshold. In this work, we chose a balance of 3-4 sentences with the constraint of 10- 70 word tokens.

## E Distribution of the Number of Tokens for Different Chunking Algorithms.

The MIMIC-III dataset was used to extract medical notes, and the chunking algorithm was then used to obtain the distribution of token counts. Naively chunking (into four sentences) resulted in very short sentences, mainly due to the formatting of MIMIC-III and the simple pre-processing done on

Figure 5: **Distribution of token count using different chunking strategies.** Top left: Sentence-based segmentation (4 sentences per snippet). Top right: Sentence-based segmentation with a token count constraint of 10-70. Bottom left: Sentence-based segmentation with a token count of 20-60. Bottom right: Sentence-based segmentation with a token count constraint of 30-50.

the notes. To address this, we opted for a sentence-based constraint on token count, resulting in improved snippets. A 10-70 constraint was chosen as it captures an appropriate amount of atomical (singleton) information, while larger constraints could lead to more extensive snippets with more information.

## F Statistics on Human Labeled Data

There are totally 46,146 patients and 2,083,159 notes in the MIMIC3 dataset. We collected 499 medical conditions as queries and sampled 1,000 patients randomly to generate the labeled data for future model training. For the evaluation purpose, we launched three runs of human evaluation: the first run randomly sampled 100 chunked note snippets across patients and notes, the second and third runs sampled 5 patients each and totally 338 notes and 1,048 note snippets. We asked at least three medical expertise to evaluate the data independently in each human evaluation run, and at the end we had 14 independent raters working on 1,079 note snippets and 14,470 snippet-condition pairs. Due to the raters' availability, 9,812 of the snippet-condition pairs were evaluated by three raters independently, 896 of them were evaluated by two raters, and the left 3,762 pairs were evaluated by only one rater.

The basic statistics of the note snippets and condition queries are shown in Figure 3. Because of the settings of our chunking algorithm, most of the snippets have reasonable length (around 60 tokens). Most of the condition queries are single words or

short phrases with 2 to 3 tokens. The keyword mixer search algorithm efficiently narrows the conditions for each snippet: on average, each snippet has about 13 relevant conditions (compare with the full list of 499 conditions), which will largely reduce the time cost of LLM labeling. About half of these pre-filtered snippet-condition pairs were further labeled as true relevant pairs, according to the majority voting of human raters.



Figure 6: **Basic Statistics of the Note Snippets and Condition Queries.** a) Distribution of snippets over length (token counts); b) Distribution of condition queries over length (token counts); c) Counts of relevant conditions of each snippet (green: based on the search engine pre-filter results; blue: majority voting from human raters); d) Counts of relevant snippets of each condition.

## G Prompting strategies

### G.1 Zero shot

You are an expert medical assistant. Your task is to give an answer of Yes/No for the relevance between a snippet and condition pair. A snippet is relevant to a condition if it includes information about the symptoms, assessments, labs, vitals, medications, procedures, or past medical history of a patient that is relevant to the given condition.

### G.2 Zero-shot CoT

You are a clinical specialist. You will be given a medical note snippet (S) and a medical condition or procedure (C). Your task is to mark whether the snippet S mentions meaningful information for C to you. Mark the answer with a binary number (0 or 1). A score of 0 indicates that the snippet does not contain meaningful content to the condition, while a score of 1 indicates that the snippet contains meaningful content. Walk me through your thoughts.

If C is a condition, snippet S contains meaningful information for C if it satisfies one of the following criterias:

(1) The snippet contains description of the condition (including explicit denial of the condition).
(2) The snippet contains description of a common cause to the condition.
(3) The snippet contains description of symptom(s) that are strongly correlated with the condition.
(4) The snippet contains description of findings that could suggest the condition (including findings that can rule out this condition).

If C is a procedure, snippet S contains meaningful information for C if it contains description of the procedure C.
S: snippet.
C: condition.
A:

### G.3 Few-shot CoT

You are an experienced clinician. You will be given a medical note snippet (S) and a medical condition or procedure (C).
Your task is to decide whether the snippet mentions useful information to a clinician for understanding the condition or procedure.
Think step by step without hallucination and provide a final Yes/No answer.

If C is a [condition], snippet S contains useful information for C if it satisfies one of the following criteria:
(1) The snippet contains information that clearly certifies or excludes C.
(2) The snippet contains highly specific information for C (symptoms, signs, or test values).

If C is a [procedure], snippet S contains useful information for C if it contains one of the following criteria.
(1) The snippet contains information that clearly certifies or excludes C.
(2) The snippet mentions clinical conditions that are highly specific to C.

Example1: C: foot pain S: ros: the patient denies any fevers, chills, weight change, nausea, vomiting, abdominal pain, diarrhea, constipation, melena, hematochezia, chest pain, shortness of breath, orthopnea, pnd, lower extremity edema, cough, urinary frequency, urgency, dysuria, lightheadedness, gait unsteadiness, focal weakness, vision changes, headache, rash, or skin changes. A:

Step 1. C (foot pain) is a common [condition] that refers to pain in the foot (lower extremity).
Step 2. Is there an explicit positive/negative signal of C in S? : No, S contains multiple negative symptoms as part of a ROS but does not contain any features related to foot pain.
Thus the answer is No.

Example2
....
ExampleN

C: condition
S: snippet
A: """

### G.4 Chain-of-Verification CoVe

**BASELINE PROMPT** = You are a medical specialist/clinician. You will be given a medical note snippet (S) and a condition/procedure (C).
Your task is to answer the below question (Q) correctly and concisely with a Yes/No answer then provide your explanation and thoughts.
Q: Does the snippet (S) directly or indirectly relate to the condition or procedure (C)?
A direct relationship is when the snippet (S) contains a description of the condition/procedure (C) or perhaps a common cause to the condition/procedure (C).
An indirect relationship is when the snippet (S) contains description of symptoms that are strongly correlated with the condition/procedure (C) or findings that could suggest the condition/procedure (C). Provide clear step by step explanations and thoughts.
S: snippet
C: condition
Answer:
**VERIFICATION QUESTIONS** = You are a medical expert. You will be given a medical note snippet (S), a condition (C ), a question (Q), and a baseline response (BR) coming from another clinician.
Your goal is to generate three verification questions that relate to both (S) and (C ). These verification questions should give a clearer guidance on how to get factual answers based on the (Q) and (BR). They are meant for verifying the factual accuracy in the baseline response (BR). The verification questions must show consistency with (Q), (BR), (S), and (C ).
S: snippet

C: condition
Q: Does the snippet (S) directly or indirectly relate to the condition or procedure (C )?
BR:baseline response
Verification Questions:
**EXECUTE PLAN PROMPT** = You are a medical expert. You will be given a medical note snippet (S), a condition (C ), some verification questions (VQ) to answer as a second opinion expert.
Your task is to provide answers to the verification questions (VQ) as correctly as possible based on the given snippet (S) and condition (C ). The verification questions (VQ) could be tricky as well, so think step by step and answer them correctly.
S: snippet
C: condition
VQ: verification questions
Answer:
**REFINEMENT** = You are a medical expert. You will be given a medical note snippet (S), a condition (C ), a medical question (MQ), a baseline response (BR), some verification questions (VQ) related to all the above, and their corresponding verification answers (VA) provided by another medical assistant.
S: snippet
C: condition
MQ: Does the snippet (S) directly or indirectly relate to the condition or procedure (C )? A direct relationship is when the snippet (S) contains a description of the condition/procedure (C ) or perhaps a common cause to the condition/procedure (C ).
An indirect relationship is when the snippet (S) contains description of symptoms that are strongly correlated with the condition/procedure (C ) or findings that could suggest the condition/procedure (C ).
BR: baseline response
VQ: verification questions
VA: verification answer
Your task is to analyze all of the above information and provide a refined [Yes/No] answer to the medical question (MQ). You must answer with a [Yes/No] response.
Make sure to provide clear explanations, a good walk through of your thoughts based on the information in (S), (C ), (MQ), (BR), (VQ), and (VA).
Answer:

## H  Interpretation of Cohen Kappa's statistics.

Table 3 provides detailed breakdown for interpreting the cohen Kappa value.

| Kappa Values | Agreement |
|---|---|
| <0 | Less than chance agreement |
| 0.01 - 0.20 | Slight agreement |
| 0.21 - 0.40 | Fair agreement |
| 0.41 - 0.60 | Moderate agreement |
| 0.61 - 0.80 | Substantial agreement |
| 0.81 - 0.99 | Almost perfect agreement |

Table 3: Interpretation of Kappa statistics (Viera et al., 2005)

## I  Frequently Asked Questions

### I.1  Why was the LLM not given a list of medical conditions to choose from when labeling a medical text snippet?

Research has shown that LLM performance is correlated with the number of tokens provided in the context (Zhang et al., 2024). Therefore, it is not sensible to use a voluminous and comprehensive list of medical conditions and provide it to the LLM for selection. An alternative and better strategy would be to provide the LLM with medical snippet-condition pairs and ask it to determine the relevance of each pair, which is the strategy used in this work.

Although this approach can reach high accuracy, it presents challenges too: as performing multiple inferences on the LLM can be computationally expensive and may result in long labeling times if resources are limited. For example, to label millions of snippets with associated thousands of conditions, the time complexity would be in the order of $O(10^8)$ or $O(10^9)$, and since LLM inference usually is slow (in seconds) thus the time cost will be in the order of $O(10^3)$ or $O(10^4)$ days. Thus, we need a fast condition filter before sending the data to LLM.

## J  The Comprehensive List of Conditions used in this study.

Our study considered the 499 most prevalent, frequently encountered and queried medical conditions and procedures in medical notes. While we only provide 20 examples below, more detailed information is available upon request: amputation,

anemia, angioedema urticaria, angioplasty of blood vessel, burn, cardiac abscess, cardiac arrest, corneal disease, cough, covid 19, flank pain, foot pain, fracture, fracture fixation, insulin resistance, lung malignancy, ovarian abscess, ophthalmologic procedure, oropharyngeal infection, pancreatitis, etc

# Conversational Topic Recommendation in Counseling and Psychotherapy with Decision Transformer and Large Language Models

**Aylin Gunal**
University of Michigan
Ann Arbor, MI
gunala@umich.edu

**Baihan Lin**
Icahn School of Medicine at Mount Sinai
New York, NY
baihan.lin@mssm.edu

**Djallel Bouneffouf**
IBM Research
Yorktown Heights, NY
djallel.bouneffouf@ibm.com

## Abstract

Given the increasing demand for mental health assistance, artificial intelligence (AI), particularly large language models (LLMs), may be valuable for integration into automated clinical support systems. In this work, we leverage a decision transformer architecture for topic recommendation in counseling conversations between patients and mental health professionals. The architecture is utilized for offline reinforcement learning, and we extract states (dialogue turn embeddings), actions (conversation topics), and rewards (scores measuring the alignment between patient and therapist) from previous turns within a conversation to train a decision transformer model. We demonstrate an improvement over baseline reinforcement learning methods, and propose a novel system of utilizing our model's output as synthetic labels for fine-tuning a large language model for the same task. Although our implementation based on LLaMA-2 7B has mixed results, future work can undoubtedly build on the design.

## 1 Introduction

In recent years, there has been a notable uptick in the number of people seeking professional help for mental health concerns, but the available pool of mental health professionals remains small in comparison. To address this need, automated AI-based tools and methods for counseling have been explored and engineered, ranging from systems for training junior mental health counselors (Min et al., 2022; Demasi et al., 2019) to AI-in-the-loop chatbots (Sharma et al., 2022). With the dramatic rise in popularity and accessibility of large language models (LLMs), it's expected that LLMs will play a significant role in the intersection of computing and mental health research, as well.

In our prior work (Lin et al., 2023b), we introduced the SupervisorBot, a reinforcement learning (RL)-based topic recommendation system in counseling conversations. This proves to be a useful

tool for clinicians during their psychotherapy sessions, where the system recommends what topics to discuss next given what has been discussed so far, as well as what works best in the past in terms of the patient outcomes. In this work, we improve upon this meaningful task by introducing the Decision Transformer (Chen et al., 2021), a transformer model designed for reinforcement learning (RL), into the recommendation pipeline, demonstrating better performance than other RL methods. We also explore the potential combination of Decision Transformer with LLMs, by generating labels for unseen transcript data using the pre-trained Decision Transformer model, and feeding the synthetically annotated data to fine-tuning a LLM. Our primary contribution is demonstrating that in the task of topic recommendation, Decision Transformer outperforms baseline RL methods; if such a system were to go through the process of user testing, the Decision Transformer—or models building on or improving Decision Transformer—can be utilized as the backbone for the recommendation module.

We first describe how we implement the preprocessing of the therapy conversation dataset, and how this is fed into the Decision Transformer model. We then describe how we use a portion of the dataset to train the Decision Transformer model, and that trained model's predicted labels are used as input to a large language model to train for the same task of topic recommendation.

## 2 Related Work

Decision Transformer was introduced as a transformer-based architecture to abstract the process of offline reinforcement learning, and has been used successfully in various NLP tasks including natural language understanding (Zhang et al., 2022; Bucker et al., 2023), navigating text-based games (Putterman et al., 2021), and generative language modeling (Memisevic et al., 2022). The Decision Transformer architecture has also been effectively

applied to the clinical domain to generate treatment recommendations based on patient history (Lee et al., 2023). In this work, we effectively apply the Decision Transformer architecture to the mental health domain in a dialogue recommendation task and improve on performances with older reinforcement learning methods.

The improvement of AI-in-the-loop tools to support humans in tasks has typically focused on human feedback, although more recent work has explored the potential for AI tools to improve themselves through a number of methods. (Saunders et al., 2022) demonstrates that a generative language model can improve its own outputs through fine-tuning on its own generations, and that the improvements are more significant as the model size increases. Generative models also have the advantage of being able to generate improvements to their own outputs (Zelikman et al., 2023).

In addition to models improving themselves, ensemble methods in which one model serves some intermediary purpose within the pipeline—e.g. data generation or filtering for input to another model—can be used for conversational modeling tasks as well (Huang et al., 2023). (Stiennon et al., 2020) uses an intermediary model's output as a reward function for another model, outperforming sole supervised learning from the source dataset. In this work, we explore a potential pipeline in which one model's output is used as synthetic data to train a language model for the task of topic recommendation. We consider the idea of AI supplementing a typical reinforcement learning with human feedback (RLHF) process by experimenting with how AI may be able to augment feedback, which can have significant implications given the lack of publicly available mental health dialogue data, let alone annotated data.

## 3 Architecture

In the following sections, we describe the architecture of our system in detail (see Fig. 1).

### 3.1 Decision Transformer

We re-implement the recommendation system pipeline as described in our original paper, (Lin et al., 2023b). This system is designed to provide real-time feedback in the form of next-topic recommendation for mental health counselors in session with patients, using reinforcement learning methods to learn and to recommend the next topic (the

*action* taken by the counselor) to move on from the current segment of dialogue (the current *state*). *Rewards* are calculated using working alliance inventory (WAI) (Horvath and Greenberg, 1989), a score from a survey of questions to determine how aligned a counselor is with their patient within a session. WAI is determined by computing similarity between inventory items and segments of dialogue (Lin et al., 2023a), and inventory items fall under three different categories: Task, Bond, and Goal. We include an aggregate WAI score, referred to as Full.

The original SupervisorBot paper evaluates the system's performance on three baseline RL algorithms: DDPG (Lillicrap et al., 2015), TD3 (Fujimoto et al., 2018), and BCQ (Fujimoto et al., 2019).

We use the Alex Street dataset [1], a dataset composed of counseling session transcripts for patients suffering from depression, anxiety, suicidal thoughts, and schizophrenia. The Alex Street dataset is preprocessed and segmented into turn-pairs, which are then embedded using Word2Vec (Mikolov et al., 2013). We use embedded topic modeling (Dieng et al., 2020) to extract 8 topics from the corpus—as determined optimal by the motivating paper—and label each turn-pair with the topic it best represents. The WAI scores are computed for each turn-pair. As the original system design is done, turn-pair embeddings represent states, topic labels represent actions, and associated WAI scores represent rewards. These items are fed as input into the Decision Transformer in the form of tuples of $(r_t, s_t, a_t)$.

We defer to the original Decision Transformer paper for architecture details. Our model contains a single-head, 3-layer attention mechanism, and we use a context window of 20 for the baseline results. Pearson's correlation between the model predicted actions and real actions taken is used for evaluation for all experiments. We run experiments 5 times using a 95%/5% train-test split, and take the average result.

### 3.2 LLMs for Recommendation

An attractive property of LLMs is flexibility in usage; at their core they simply model language probability distributions, making their outputs malleable to various tasks in NLP. In this section, we

---

[1]https://alexanderstreet.com/products/counseling-and-psychotherapy-transcripts-series

Figure 1: Architecture of the proposed LLM integration, demonstrating how both gold-standard labels from the dataset as well as synthetic annotations from Decision Transformer output can be used to fine-tune the LLM.

explore LLMs' ability for dialogue classification into labels of different psychiatric conditions to demonstrate their usefulness in various components of RLHF, similar to a diagnostic scenario in clinical setting as in (Lin et al., 2022, 2024). We primarily experiment with LLaMA-2 with 7B parameters (Touvron et al., 2023). We fine-tune LLaMA-2 model for sequence classification and test using the same preprocessing steps and train-test split described in Section 3.1. During the fine-tuning process, we do not use a validation dataset to avoid data leakage into the test set.

We experiment with treating Decision Transformer predictions as synthetic gold standard annotations for the LLM to learn from. We split the full dataset in a 40%/40%/20% split; the first 40% of the Alex Street dataset is used to train Decision Transformer, then Decision Transformer outputs predictions for another 40% of the dataset which the LLM is then fine-tuned with, and ultimately the LLM is evaluated on the final 20% of the dataset. Due to computational constraints, we apply low-rank adaptation (LoRA) (Hu et al., 2021), a parameter efficient fine-tuning method, in order to optimize the fine-tuning process. The LoRA configuration includes an alpha of 16 and dropout rate of .05, and we fine-tune for 1 epoch. We target the LLaMA-2 model's attention layers during training and save the final layer's weights to avoid those scores being randomly initialized for inference. To provide some baseline for comparison, we additionally fine-tune the LLaMA-2 model on the original gold-standard labels.

## 4 Evaluation

### 4.1 Results

The results of the Decision Transformer on a 95%/5% train-test split, reflecting the set-up of the original SupervisorBot paper, are provided in Table 1. We reproduce results for the other RL methods in the original paper for performance on the full-scale rewards; Decision Transformer outperforms these baselines as noted in Table 3. We note that Decision Transformer specifically performs best for all reward scales when trained on the full dataset; among individual diseases, the model performs best on the task, bond, and goal scales for anxiety.

We additionally evaluate whether or not the 20-timestep context is necessary for good performance from the Decision Transformer model, and these results are provided in Table 2. We note that 15 time-steps is optimal for a majority of the reward sclaes, suggesting that the Decision Transformer is better able to make decisions provided a briefer learning history. An advantage of utilizing a transformer-based model for this task is that we are able to investigate its internal structure to understand specifically which historical features—including which time-steps—are significant for inference.

Additionally, we note that the LLaMA-2 model trained on the gold-standard data does not necessarily outperform the Decision Transformer for all reward scales as indicated in Table 4, indicating that the off-the-shelf language model may not be conducive for a reinforcement learning task. LLaMA-2 trained on the Decision Transformer output directly

**Decision Transformer**

| | Depression | Anxiety | Schizophrenia | Suicidal | All |
|---|---|---|---|---|---|
| Full | .176 | .233 | .246 | .213 | .361 |
| Task | **.291** | **.320** | .247 | .231 | .323 |
| Bond | .270 | .314 | .231 | **.239** | .335 |
| Goal | **.291** | .313 | **.249** | .229 | **.375** |

Table 1: Results of Decision Transformer on topic recommendation task, using previous 20 turn-pairs as input. Best results per data subset are in bold.

| | Context Lengths | | | |
|---|---|---|---|---|
| **Rewards** | 5 | 10 | 15 | 20 |
| Full | 0.346 | 0.345 | **0.403** | 0.361 |
| Bond | 0.284 | 0.343 | **0.359** | 0.335 |
| Task | 0.272 | 0.298 | **0.342** | 0.322 |
| Goal | 0.278 | 0.339 | 0.348 | **0.375** |

Table 2: Decision Transformer model performance trained on varying context lengths. Best results per reward scale are in bold.

| | DDPG | BCQ | TD3 |
|---|---|---|---|
| Full | .264 (-.97) | .170 (-1.91) | .286 (-.75) |

Table 3: Baseline RL performance on full-scale rewards on the full dataset, with a comparison to DT performance.

also does not perform particularly well; future work may include modifying the way in which the Decision Transformer synthetic labels are used by a language model. It's possible that prompting the language model may yield better results than treating it as a sequence classifier.

### 4.2 Additional Analysis for Interpretability

We extract the final layer of attention weights from the Decision Transformer models trained on the four reward scales for the three types of inputs: returns, states, and actions. We observe both the attention weights for the individual input types as well as the aggregated and averaged set of weights aross all input types. Due to the auto-regressive nature of Decision Transformer, attention weights

| | Full | Task | Bond | Goal |
|---|---|---|---|---|
| **LLaMA-2 7B + DT** | .148 | .118 | .158 | .115 |
| **LLaMA-2 7B + Gold** | .371 | .259 | .315 | .332 |

Table 4: Results of fine-tuning LLaMA-2 7B on DT output and gold-standard labels.



Figure 2: Normalized attention scores associated with absolute timesteps, *without* padded sequences.



Figure 3: Normalized attention scores associated with relative timesteps.

are assigned to the *timesteps* prior to the recommendation made at a given timestep.

We provide visual analyses of attention scores through normalized aggregate attention scores per timestep for absolute timestep values (Fig. 2) as well as relative timestep values (Fig. 3). We note that the model refines its attention to generally focus on items in earlier positions in given input sequence, both in the case for absolute and relative timesteps. These results, in tandem with the generally higher performances of the model on 15 previous timesteps rather than 20 timesteps, indicate that potentially there is a beginning index to the current context that can be key for the model's inference ability. Future work may include adjusting the context window dynamically, both for training and inference.

## 5   Limitations

Due to limitations of computational resources, experimentation with fine-tuning LLMs is restricted by model size. Future work can build on this work by applying similar experiments on increasing model sizes or non-quantized versions of models, effectively demonstrating (positively or negatively) that performance scales with model size.

## 6   Ethical Considerations

When implementing a topic recommendation system in counseling contexts, ethical considerations are important due to the sensitive nature of digital mental health discussions, as discussed in (Lin, 2022). One of the primary concerns is the potential limitation imposed by a static set of discussion topics. While such a system can streamline the counseling process, it risks limiting the creativity and flexibility of counselors, particularly those in training, and in the long term, inhibit consideration of their own perspectives on how to continue the conversation. This could inadvertently restrict their ability to tailor sessions according to the unique needs of each patient.

This is particularly relevant since the topics pulled are from one specific dataset that covers only four mental health conditions. The training dataset, derived from this limited number of mental health conditions, might not be representative of the broader population or other conditions. This limitation can lead to biased recommendations if not carefully managed. To mitigate this, it is essential to consider a more dynamic approach where the set

of topics can evolve based on ongoing input from practicing counselors and feedback from therapy sessions. This adaptation would help in maintaining the relevance and sensitivity of the recommendations to diverse patient needs. In deployment, we can also imagine that topics are dynamically chosen, or chosen using human feedback; for example, perhaps before the system is put into use, counselors can input their own topics.

In addition to dataset limitations, the calculation of rewards, based on the Working Alliance Inventory (WAI), while rooted in established psychological theory, may benefit from enhancements through reinforcement learning with human feedback (RLHF). Incorporating direct input from users could refine the understanding and alignment of counselor and patient goals, improving the system's effectiveness and ethical alignment.

## 7   Conclusion

In this study, we introduced a Decision-Transformer-based recommendation system which outperforms baseline RL-based methods in counseling topic recommendation, indicating that transformer-based methods may have better performance in general when it comes to modeling conversation direction and alignment. We additionally find that the model performs best for certain reward scales on shorter input sequences, indicating that some exploration of optimal sequence length can be an avenue for future work. Through additional analysis of the attention scores, we additionally find that the model pays more attention to items earlier on in the input sequence.

## References

Arthur Bucker, Luis Figueredo, Sami Haddadin, Ashish Kapoor, Shuang Ma, Sai Vemprala, and Rogerio Bonatti. 2023. Latte: Language trajectory transformer. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7287–7294. IEEE.

Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, P. Abbeel, A. Srinivas, and Igor Mordatch. 2021. Decision transformer: Reinforcement learning via sequence modeling. In *Neural Information Processing Systems*.

Orianna Demasi, Marti A. Hearst, and Benjamin Recht. 2019. Towards augmenting crisis counselor training by improving message retrieval. *Proceedings of the Sixth Workshop on Computational Linguistics and Clinical Psychology*.

Adji B Dieng, Francisco JR Ruiz, and David M Blei. 2020. Topic modeling in embedding spaces. *Transactions of the Association for Computational Linguistics*, 8:439–453.

Scott Fujimoto, Herke Hoof, and David Meger. 2018. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR.

Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, pages 2052–2062. PMLR.

Adam O Horvath and Leslie S. Greenberg. 1989. Development and validation of the working alliance inventory. *Journal of Counseling Psychology*, 36:223–233.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Xin Huang, Kye Min Tan, Richeng Duan, and Bowei Zou. 2023. Ensemble method via ranking model for conversational modeling with subjective knowledge. In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 177–184.

Seunghyun Lee, Da Young Lee, Sujeong Im, Nan Hee Kim, and Sung-Min Park. 2023. Clinical decision transformer: intended treatment recommendation through goal prompting. *arXiv preprint arXiv:2302.00612*.

Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*.

Baihan Lin. 2022. Computational inference in cognitive science: Operational, societal and ethical considerations. *arXiv preprint arXiv:2210.13526*.

Baihan Lin, Djallel Bouneffouf, Yulia Landa, Rachel Jespersen, Cheryl Corcoran, and Guillermo Cecchi. 2024. Compass: Computational mapping of patient-therapist alliance strategies with language modeling. *arXiv preprint arXiv:2402.14701*.

Baihan Lin, Guillermo Cecchi, and Djallel Bouneffouf. 2022. Working alliance transformer for psychotherapy dialogue classification. *arXiv preprint arXiv:2210.15603*.

Baihan Lin, Guillermo Cecchi, and Djallel Bouneffouf. 2023a. Deep annotation of therapeutic working alliance in psychotherapy. In *International workshop on health intelligence*, pages 193–207. Springer.

Baihan Lin, Guillermo Cecchi, and Djallel Bouneffouf. 2023b. Supervisorbot: Nlp-annotated real-time recommendations of psychotherapy treatment strategies with deep reinforcement learning. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pages 7149–7153.

Roland Memisevic, Sunny Panchal, and Mingu Lee. 2022. Decision making as language generation. In *NeurIPS 2022 Foundation Models for Decision Making Workshop*.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Neural Information Processing Systems*.

Do June Min, Verónica Pérez-Rosas, Kenneth Resnicow, and Rada Mihalcea. 2022. Pair: Prompt-aware margin ranking for counselor reflection scoring in motivational interviewing. In *Conference on Empirical Methods in Natural Language Processing*.

Aaron L Putterman, Kevin Lu, Igor Mordatch, and Pieter Abbeel. 2021. Pretraining for language conditioned imitation with transformers.

William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Ouyang Long, Jonathan Ward, and Jan Leike. 2022. Self-critiquing models for assisting human evaluators. *ArXiv*, abs/2206.05802.

Ashish Sharma, Inna Wanyin Lin, Adam S. Miner, David C. Atkins, and Tim Althoff. 2022. Human–ai collaboration enables more empathic conversations in text-based peer-to-peer mental health support. *Nature Machine Intelligence*, 5:46–57.

Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

E. Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. 2023. Self-taught optimizer (stop): Recursively self-improving code generation. *ArXiv*, abs/2310.02304.

Ziqi Zhang, Yile Wang, Yue Zhang, and Donglin Wang. 2022. Can offline reinforcement learning help natural language understanding? *arXiv preprint arXiv:2212.03864*.

# Leveraging Wikidata for Biomedical Entity Linking in a Low-Resource Setting: A Case Study for German

**Faizan E Mustafa**
QUIBIQ GmbH

**Corina Dima**
University of Stuttgart

**Juan G. Diaz Ochoa**
PerMediQ GmbH
QUIBIQ GmbH

**Steffen Staab**
University of Stuttgart
University of Southampton

## Abstract

Biomedical Entity Linking (BEL) is a challenging task for low-resource languages, due to the lack of appropriate resources: datasets, knowledge bases (KBs), and pre-trained models. In this paper, we propose an approach to create a biomedical knowledge base for German BEL using UMLS information from Wikidata, that provides good coverage and can be easily extended to further languages. As a further contribution, we adapt several existing approaches for use in the German BEL setup, and report on their results. The chosen methods include a sparse model using character n-grams, a multilingual biomedical entity linker, and two general-purpose text retrieval models. Our results show that a language-specific KB that provides good coverage leads to most improvement in entity linking performance, irrespective of the used model. The finetuned German BEL model, newly created UMLS$_{Wikidata}$ KB as well as the code to reproduce our results are publicly available[1].

## 1 Introduction

BEL is the task of disambiguating text spans by linking them to a unique identifier in a biomedical knowledge base (French and McInnes, 2023). For instance, the UMLS (Bodenreider, 2004) entity having the Concept Unique Identifier (CUI) C0007765 is usually mentioned in English under the name *cerebellum* but can be mentioned in German using either *Kleinhirn, Cerebellum* or *Zerebellum*. Each *entity* in the KB has an *entity name* and one or multiple *aliases* associated with it, in multiple languages, as shown in Fig. 1. In this work, we refer to such names as *entity mentions*. The task of biomedical entity linking is to recover the unambiguous entity identifier from a KB given either of the names that can be used to refer to an entity. The task can be performed *with context* - where

the name is provided together with the surrounding text, or *without context* - where only the name itself is provided for the disambiguation. In this paper we tackle the problem of BEL for entity mentions without context.



Figure 1: QID for an entity in Wikidata

A wide range of models using ruled-based and deep learning approaches for BEL have been proposed for English, for which many data resources are available (Shi et al., 2023). However, the in-domain BEL datasets, KBs, and models are scarce for low-resource languages. Multilingual biomedical models such as SapBERT (Liu et al., 2021a) have been proposed and evaluated on cross-lingual BEL benchmarks like XL-BEL (Liu et al., 2021b). This benchmark, however, is only intended for evaluation purposes, as it includes only 1,000 samples per language.

Wang et al. (2023) proposed a comprehensive German BEL benchmark, WikiMed-DE-BEL, which has, however, not yet been used for evaluating BEL models. We adapt several models from the literature to BEL on German, and evaluate them on this new benchmark.

A problematic aspect when training a BEL model for German is the lack of a biomedical KB with entity names and descriptions in German. The Unified Medical Language System (UMLS) (Bo-

---

[1]German-Bio-Entity-Linking GitHub Repository

denreider, 2004), the most comprehensive biomedical thesaurus available to date, which is the standard KB in BEL for English, only contains 1.6% entities with German descriptions (Liu et al., 2021b). We propose a solution to this problem by building a German biomedical KB using UMLS information harvested from Wikidata (Vrandečić and Krötzsch, 2014), an approach that leads to better entity coverage and can be extended to further languages.

## 2 Knowledge Bases

### 2.1 UMLS

UMLS (Bodenreider, 2004) is a metathesaurus integrating information from multiple biomedical vocabularies with the aim of improving interoperability. The terminology utilized across vocabularies is standardized by assigning a unique identifier, called the Concept Unique Identifier (CUI) to the same entities modeled in different vocabularies and across multiple languages. The latest UMLS Metathesaurus release, 2023AB, contains approximately 3.36 million concepts and 15.9 million unique concept names from 185 source vocabularies[2].

### 2.2 Wikidata

Wikidata (Vrandečić and Krötzsch, 2014) is a collaborative knowledge base providing the data for many Wikimedia projects, including the multilingual Wikipedia. Wikidata currently consists of more than 100M items that have been edited over 2 billion times by Wikidata users[3]. A defining trait of Wikidata is that it serves as a hub for integrating knowledge from different domains, including the biomedical domain. Wikidata entities can be connected, for example, to the UMLS, to the Disease Ontology or to many other biomedical vocabularies through pre-defined properties.

## 3 BEL Datasets for German

### 3.1 WikiMed-DE-BEL

WikiMed-DE (Wang et al., 2023) is a silver-standard biomedical entity linking dataset for the German language. It was built starting from German Wikipedia articles with hyperlinked text, where the hyperlinks are considered to be entity mentions and are linked to the corresponding Wikidata unique item identifiers (QIDs). The QIDs were then used to assign unique concept IDs from several

biomedical vocabularies including UMLS. The annotations for each article include the article's title, text, QID, biomedical vocabularies concept IDs as well as a list of mentions, each assigned an unique QID as well as biomedical concept IDs. The creators of the WikiMed-DE dataset released a high-quality subset named WikiMed-DE-BEL which we use as a benchmark. WikiMed-DE-BEL includes 53,981 articles from the German Wikipedia. The `train`, `test` and `dev` splits follow the 80/10/10 rule.

We post-process WikiMed-DE-BEL as follows: for each data split, we only keep unique (mention, CUI) pairs. To increase the number of available pairs we create pairs both from the article title and the CUI assigned to the whole article, as well as from the entity mentions inside the article together with their assigned CUI. The `train`, `dev`, and `test` sets contain 42,679, 13,017, and 13,019 unique CUIs and 79,904, 19,561, and 19,203 (CUI, mention) pairs, respectively.

### 3.2 XL-BEL

XL-BEL (Liu et al., 2021b) is a cross-lingual biomedical entity linking evaluation benchmark that covers 10 languages, including German. Entity mentions from Wikipedia articles in the target languages were linked to language-agnostic UMLS CUIs using the methodology proposed by (Vashishth et al., 2021). The dataset samples are (sentence, mention, CUI) triples extracted from these Wikipedia articles. A number of 1,000 samples were retained for each language, making sure that each surface form appears only once in the sampled examples. We use the German subset of XL-BEL for evaluation purposes.

## 4 Models for German BEL

To the best of our knowledge, there are no existing dedicated models for German BEL that are publicly available. We therefore selected several models that could be adapted to German. Because we perform BEL without context, we also report on results obtained using embedding models trained for text retrieval. In this case, the evaluation is based on the nearest neighbour search, using the mention as an input query.

**ScispaCy.** Neumann et al. (2019) introduce ScispaCy, a Python library for biomedical text processing. One of the provided models creates sparse vector representations of the entity names and aliases

---

from the KB by representing them in terms of the TF-IDF scores of character 3-grams which are part of ten or more entities from the given KB. An entity mention is similarly modeled in terms of its character 3-grams and is linked to the KB by retrieving the $k$ nearest neighbours from the KB. This process is language-agnostic and can easily be applied to other languages as long as there is an available KB in the target language.

**SapBERT.** Liu et al. (2021a) propose SapBERT, a pre-training scheme for self-aligning transformer-based representations to KBs based on synonymy relations. The training process takes as input a list of (mention, CUI) pairs from the KB, where the mention could be either the entity name or one of its aliases (see Fig. 1). The authors use a dedicated mining process to discover informative training examples: within a mini-batch they look for triples of the form $(x_a, x_p, x_n)$ where $x_a$ is an *anchor*, a random mention from the mini-batch, $x_p$ is a *positive match* for $x_a$ and $x_n$ is a *negative match* for $x_a$. A positive match has the same CUI as the anchor mention, whereas a negative match has a different CUI than the anchor. For the example in Fig. 1 a triplet could be (*Kleinhirn*, *Cerebellum*, *Gehirn*), where the first two items in the triplet refer to the same CUI, C0007765, whereas *Gehirn*, German for *brain*, refers to a different CUI, C0006104. Each triplet contributes a positive pair and a negative pair towards the training data. The model is then trained using an adapted version of multi-similarity loss (Wang et al., 2019). The goal is to bring the representations of positive pairs closer to each other while pushing the negative pairs far from each other. Each mention is represented using the output [CLS] token resulting from feeding the mention text through the base transformer model.

**M3 Embeddings.** Chen et al. (2024) proposed a multilingual, hybrid text retrieval approach that can model input texts of up to 8192 tokens. A self-knowledge distillation framework is used to jointly learn three retrieval methods (dense, sparse, multi-vector) which reinforce each other. The model can be used for query-based text retrieval in more than 100 languages, including German.

**Jina Embeddings.** Mohr et al. (2024) developed an German-English bilingual model by pre-training a BERT-based language model on bilingual text. The model is then trained as an embedding retrieval model using contrastive learning by fine-tuning on text pairs $(q, p)$ consisting of a query string $q$ and a target string $p$. The evaluation indicate

a considerable improvement in German-English cross-lingual retrieval performance when compared to multilingual models.

## 5 Creating a German Biomedical Knowledge Base

Liu et al. (2021b) report that 69.6% of the names of the UMLS entities in release 2020AA are in English, but only 1.6% are in German. The multilingual UMLS subset they use to evaluate SapBERT, UMLS$_{SapBERT}$, is provided by the SapBERT authors in their GitHub repository[4]. It contains 399,931 entity names or aliases assigned to 62,094 unique CUIs. Most of the names are in English, with only a small fraction being in German. The number of unique (entity, CUI) pairs amounts to 260,633.

We create a large German biomedical KB, UMLS$_{Wikidata}$, by leveraging Wikidata information. We first obtain a list of Wikidata QIDs that are annotated with CUIs by querying Wikidata using the official SPARQL endpoint[5] to fetch items that have the *UMLS CUI* property (P2892). The QIDs are further used to obtain the German label, description and alias(es) using the Python package *qwikidata* [6]. The resulting KB has 599,330 unique CUIs and 671,797 unique (entity name, CUI) pairs, where all the entity names are in German. Table 1 shows the statistics of the two KBs. UMLS$_{Wikidata}$ KB is made publicly available for further use[7].

| | Unique CUIs | Unique (CUI, Entity) Pairs |
|---|---|---|
| UMLS$_{Wikidata}$ | 599,330 | 671,797 |
| UMLS$_{SapBERT}$ | 62,094 | 260,633 |

Table 1: KB Statistics

## 6 Methodology

The first step in the evaluation of each of the selected models is to create vector representations for all KB entities using each model in turn and then store the obtained entity representations in a Faiss index (Johnson et al., 2019) for efficient retrieval.

The linking step for all the models involves first creating a vector representation for the entity mention using the selected model and then finding the $k$ nearest neighbors from the KB by comparing the

---

[4]SapBERT UMLS subset.
[5]Official Wikidata SPARQL endpoint.
[6]https://pypi.org/project/qwikidata/
[7]https://zenodo.org/records/11003203

mention vector to the KB vector representations stored in the corresponding Faiss index using cosine similarity. Mentions are linked to the 5 nearest neighbors for all the models.

We further fine-tune the SapBERT-UMLS model[8], which is already trained on multilingual UMLS pairs, on $UMLS_{Wikidata}$. We use the same procedure as described in Section 4 and train for 5 epochs using a batch size of 256. The fine-tuned model is available on Hugging Face Model Hub[9]

The only hyperparameter of the ScispaCy model is the size of the character n-grams to be used. We use the 3-grams that appear in 10 or more entities in the target KB. We only use the dense representations from the M3 embedding model. For performance reasons, the maximum sequence length of all the embedding models is set to 40 tokens.

## 7 Results

The evaluation metric precision@k ($p@k$) indicates the percentage of samples where the correct entity is found in the top $k$ KB entities predicted by a model. Tables 2 and 3 report the $p@1$ and $p@5$ obtained by the various models when linking against the $UMLS_{Wikidata}$ and the $UMLS_{SapBERT}$, respectively. As a general trend, the sparse, ScispaCy-based n-gram models score lower than the embedding models. The difference is more pronounced when using the $UMLS_{SapBERT}$ KB (in Table 3) because here the descriptions are mostly in English and thus the character 3-grams selected from the KB for the model have less overlap with the German mentions. The Jina embeddings outperform the rest of the embedding models when using the $UMLS_{Wikidata}$ knowledge base.

SapBERT fine-tuned on $UMLS_{Wikidata}$ offers good, consistent performance: it performs on par with the Jina model when using the $UMLS_{Wikidata}$ KB (see Table 2) and outperforms the rest of the models by a large margin, showing a 6 point improvement in $p@1$ score for XL-BEL when using the $UMLS_{SapBERT}$ KB (see Table 3). We hypothesize that this is due to the benefits of fine-tuning on the extra names contained in $UMLS_{Wikidata}$, as it allows the model to learn a better English-German cross-lingual mapping, as many medical terms are common between English and German.

Overall scores are much higher when using the $UMLS_{Wikidata}$ KB instead of the $UMLS_{SapBERT}$

---

| Model | Metrics | WikiMed-DE-BEL | | | XL-BEL DE |
|---|---|---|---|---|---|
| | | Train | Dev | Test | |
| ScispaCy | p@1 | 0.755 | 0.782 | 0.785 | 0.492 |
| using UMLS$_{Wikidata}$ 3-grams | p@5 | 0.824 | 0.847 | 0.851 | 0.590 |
| SapBERT (Liu et al., 2021a) | p@1 | 0.756 | 0.783 | 0.785 | 0.462 |
| | p@5 | 0.822 | 0.846 | 0.850 | 0.568 |
| SapBERT | p@1 | 0.774 | 0.796 | 0.80 | 0.485 |
| fine-tuned on UMLS$_{Wikidata}$ | p@5 | 0.840 | 0.861 | 0.863 | 0.590 |
| M3 embeddings | p@1 | 0.767 | 0.791 | 0.795 | **0.499** |
| | p@5 | 0.836 | 0.857 | 0.860 | 0.604 |
| Jina embeddings | p@1 | **0.777** | **0.803** | **0.805** | 0.495 |
| | p@5 | **0.840** | **0.861** | **0.864** | **0.605** |

Table 2: Results using the $UMLS_{Wikidata}$ KB.

KB because of its larger size and because it provides better coverage for the German entities in the two evaluation datasets. Moreover, the scores for WikiMed-DE-BEL are significantly higher than for XL-BEL when using the $UMLS_{Wikidata}$ KB but the opposite is true when using the $UMLS_{SapBERT}$ KB. The reason for this behaviour is discussed next.

| Model | Metrics | WikiMed-DE-BEL | | | XL-BEL DE |
|---|---|---|---|---|---|
| | | Train | Dev | Test | |
| ScispaCy | p@1 | 0.118 | 0.117 | 0.117 | 0.286 |
| using UMLS$_{SapBERT}$ 3-grams | p@5 | 0.141 | 0.141 | 0.142 | 0.359 |
| SapBERT (Liu et al., 2021a) | p@1 | 0.139 | 0.147 | 0.146 | 0.346 |
| | p@5 | 0.154 | 0.162 | 0.161 | 0.396 |
| SapBERT | p@1 | **0.172** | **0.181** | **0.177** | **0.401** |
| fine-tuned on UMLS$_{Wikidata}$ | p@5 | **0.197** | **0.206** | **0.204** | **0.473** |
| M3 embeddings | p@1 | 0.138 | 0.143 | 0.143 | 0.342 |
| | p@5 | 0.155 | 0.160 | 0.160 | 0.401 |
| Jina embeddings | p@1 | 0.141 | 0.148 | 0.149 | 0.338 |
| | p@5 | 0.158 | 0.166 | 0.165 | 0.394 |

Table 3: Results using the $UMLS_{SapBERT}$ KB.

## 8 KB Coverage

The results obtained for the different dataset/KB combinations are drastically different. The precision is above 0.75 for WikiMed-DE-BEL using the $UMLS_{Wikidata}$ KB, but below 0.20 when using the $UMLS_{SapBERT}$ KB. If a mention's CUI is not present in KB then the model cannot link to it. Therefore, we check the upper limit for the metric scores by calculating the dataset coverage for the two KBs. Table 4 shows, for each dataset, the percentage of dataset CUIs that are present in the KB CUIs. It can be noticed that only 36% of the WikiMed training set CUIs are present in the $UMLS_{SapBERT}$ KB, in contrast to 98% coverage when using the $UMLS_{Wikidata}$ KB.

| KB | WikiMed-DE-BEL | | | XL-BEL DE |
|---|---|---|---|---|
| | Train | Dev | Test | |
| UMLS$_{Wikidata}$ | 98.2% | 97.5% | 97.6% | 81.0% |
| UMLS$_{SapBERT}$ | 36.4% | 36.7% | 37.1% | 99.8% |

Table 4: CUI coverage.

Another problematic setup is when the a particular name of an entity or alias is not present in the

KB, even its CUI is in KB. Therefore, we compute the (mention, CUI) pair coverage by looking at the percentage of (mention, CUI) pairs present in the respective KB. Table 5 shows that the pair coverage for XL-BEL wrt. to $\text{UMLS}_{SapBERT}$ is 11%, whereas for WikiMed-DE-BEL wrt. to $\text{UMLS}_{Wikidata}$ is 56% — which aligns better with the model performance reported in Tables 2 and 3.

| KB | WikiMed-DE-BEL | | | XL-BEL DE |
|---|---|---|---|---|
| | Train | Dev | Test | |
| $\text{UMLS}_{Wikidata}$ | 56.5% | 62.4% | 63.0% | 33.8% |
| $\text{UMLS}_{SapBERT}$ | 6.2% | 5.9% | 6% | 11.8% |

Table 5: (mention, CUI) pairs coverage.

## 9 Conclusion

The unavailability of knowledge bases, datasets and, subsequently, models makes BEL a challenging task for low-resource languages. To this end, we propose an approach to create a KB for German BEL, $\text{UMLS}_{Wikidata}$, using a methodology that can be easily applied to further low-resource languages. We further compare four different models with various representations and trained on different languages. Our results show that creating a dedicated, large-scale knowledge base in the target language leads to the most improvement for doing entity linking in that language, independently of the used model. The best BEL results for German are obtained using the language-specific $\text{UMLS}_{Wikidata}$ knowledge base.

## Acknowledgements

## References

Olivier Bodenreider. 2004. The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic Acids Res.*, 32(Database-Issue):267–270.

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.

Evan French and Bridget T. McInnes. 2023. An overview of biomedical entity linking throughout the years. *Journal of Biomedical Informatics*, 137:104252.

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.

Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. 2021a. Self-alignment pretraining for biomedical entity representations. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4228–4238, Online. Association for Computational Linguistics.

Fangyu Liu, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2021b. Learning domain-specialised representations for cross-lingual biomedical entity linking. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 565–574, Online. Association for Computational Linguistics.

Isabelle Mohr, Markus Krimmel, Saba Sturua, Mohammad Kalim Akram, Andreas Koukounas, Michael Günther, Georgios Mastrapas, Vinit Ravishankar, Joan Fontanals Martínez, Feng Wang, Qi Liu, Ziniu Yu, Jie Fu, Saahil Ognawala, Susana Guzman, Bo Wang, Maximilian Werk, Nan Wang, and Han Xiao. 2024. Multi-Task Contrastive Learning for 8192-Token Bilingual Text Embeddings.

Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. 2019. ScispaCy: Fast and robust models for biomedical natural language processing. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 319–327, Florence, Italy. Association for Computational Linguistics.

Jiyun Shi, Zhimeng Yuan, Wenxuan Guo, Chen Ma, Jiehao Chen, and Meihui Zhang. 2023. Knowledge-graph-enabled biomedical entity linking: a survey. *World Wide Web*, pages 1–30.

Shikhar Vashishth, Denis Newman-Griffis, Rishabh Joshi, Ritam Dutt, and Carolyn P. Rosé. 2021. Improving broad-coverage medical entity linking with semantic type prediction and large-scale datasets. *Journal of Biomedical Informatics*, 121:103880.

Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85.

Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R. Scott. 2019. Multi-Similarity Loss With General Pair Weighting for Deep Metric Learning. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5017–5025.

Yi Wang, Corina Dima, and Steffen Staab. 2023. WikiMed-DE: Constructing a Silver-Standard Dataset for German Biomedical Entity Linking using Wikipedia and Wikidata. In *Proceedings of the Wikidata Workshop 2023 co-located with 22nd International Semantic Web Conference (ISWC 2023)*.

# Revisiting Clinical Outcome Prediction for MIMIC-IV

**Tom Röhr\*, Alexei Figueroa\*, Jens-Michalis Papaioannou\*◇,**
**Conor Fallon\*, Keno Bressem†, Wolfgang Nejdl◇, Alexander Löser\***

\*DATEXIS, Berliner Hochschule für Technik
†Department of Radiology and Nuclear Medicine, German Heart Center Munich
◇L3S, Leibniz University Hannover

{troehr, afigueroa, michalis.papaioannou, cfallon, aloeser}@bht-berlin.de
bressem@dhm.mhn.de
nejdl@L3S.de

## Abstract

Clinical Decision Support Systems assist medical professionals in providing optimal care for patients. A prominent data source used for creating tasks for such systems is the *Medical Information Mart for Intensive Care* (MIMIC). MIMIC contains electronic health records (EHR) gathered in a tertiary hospital in the United States. The majority of past work is based on the third version of MIMIC, although the fourth is the most recent version. This new version, not only introduces more data into MIMIC, but also increases the variety of patients. While MIMIC-III is limited to intensive care units, MIMIC-IV also offers EHRs from the emergency department. In this work, we investigate how to adapt previous work to update clinical outcome prediction for MIMIC-IV. We revisit several established tasks, including prediction of diagnoses, procedures, length-of-stay, and also introduce a novel task: *patient routing prediction*. Furthermore, we quantitatively and qualitatively evaluate all tasks on several bio-medical transformer encoder models. Finally, we provide narratives for future research directions in the clinical outcome prediction domain. We make our source code publicly available to reproduce our experiments, data, and tasks.

## 1 Introduction

Estimating the future clinical state of a patient upon admission to a medical care facility is a task of critical importance. Clinicians must be able to promptly gauge not only the main affliction of patients, but also all the resources needed to streamline their care. A Clinical Decision Support System (CDSS) aids clinicians in a multifaceted way; for instance, they can interact with a clinician in a conversational manner or they can assist in the diagnosis process by offering discrete suggestions. Generative medical assistants, like AMIE (McDuff et al., 2023), enable clinicians to derive diagnostics and treatments by engaging in a conversation with the language model. One way of communicating these findings is to use the International Classification of Diseases (ICD) taxonomy which is also used by medical practitioners to document the admission of a patient, their stay, and release from a medical care facility. While conversational CDSS can provide reasonable answers and may identify important treatment strategies, their suggestions veer substantially from expert suggestions (Benary et al., 2023). Furthermore, validating these suggestions is difficult, given the arbitrarily large output space of decoder-based transformer architectures such as AMIE. However, it is essential for clinicians to validate the predictions of such systems in order to safeguard the well-being of their patients. Given the discrete space of the ICD taxonomy and the necessity of validation, we argue that classification with encoder models is relevant for the clinical outcome prediction domain.

**Clinical Outcome Prediction from Admission Notes.** We revisit the clinical outcome prediction (COP) tasks as defined in van Aken et al. (2021). These tasks are all based on the third version of the *Medical Information Mart for Intensive Care* (MIMIC-III)(Johnson et al., 2016). Therefore, in this work, we refer to these tasks as *COP-III*. Since the publication of *COP-III*, a new version of MIMIC has been released, MIMIC-IV (Johnson et al., 2023). MIMIC-IV supersedes the third version with more patient data from the intensive care units (ICU). Additionally, it includes data from patients admitted to the emergency department (ED). This increase in available data, both in quantity and diversity, renders the tasks of *COP-III* obsolete. We present *COP-IV*, an updated and extended set of 6 clinical outcome prediction tasks based on MIMIC-IV. This includes 3 out of 4 *COP-III* tasks adapted for the MIMIC-IV ICU

and ED splits respectively, as well as a novel *patient routing* task. The patient routing task utilizes the exclusive routing information of MIMIC-IV to predict the first transfer of a patient upon admission. We update the three *COP-III* tasks by adapting the data-processing methods to suit MIMIC-IV. Alongside updating the admission note data, we update the target space from ICD-9 to ICD-10. This provides more relevance for clinicians since ICD-10 is the coding version in use since 2015. We evaluate all *COP-IV* tasks against a selection of open[1] clinical transformer encoder models. Moreover, we compare our results for *COP-IV* and the results of van Aken et al. (2021) for *COP-III* to assess whether the performance for clinical outcome prediction improves with the new data.

**Contributions.** We summarize our contributions as follows:

- We create novel datasets for several outcome prediction tasks, derived from data in both the intensive care unit (ICU) and the emergency department (ED).

- We introduce a novel *patient routing* task, derived from the patient routing information available in the emergency department module of MIMIC-IV. Resulting in 6 tasks overall, with 3 tasks belonging to ICU and ED prediction respectively.

- We benchmark multiple biomedical transformer encoder models on *COP-IV* and present our qualitative and quantitative analysis.

- We present challenges of *COP-IV* and propose future work directions for clinical outcome prediction.

- We release our source code to reproduce our experiments and datasets[2].

## 2 Related Work

**Bio-medical encoders.** In the context of transfer learning, several works explore adapting encoder transformer networks such as BERT (Devlin et al., 2018) into specialized settings.

*BioBERT*(Lee et al., 2019) presents improved performance in bio-medical text mining tasks, by

continuing pre-training a BERT model on full-text and abstracts of research articles from PubMed.

Both *ClinicalBert* and *DischargeBERT* (Alsentzer et al., 2019) further pre-train BioBERT models on full-text notes and discharge notes respectively from the MIMIC-III dataset.

*CORe* (van Aken et al., 2021) reformulates BERT's unsupervised *next-sentence-prediction* pre-training objective as an *admission-discharge-relation*, tasking a BioBERT model to classify whether a sequence coming from an admission-note relates to the discharge section of the same patient.

In contrast to improving a pre-trained BERT or BioBERT model, *PubmedBERT*(Gu et al., 2020) achieves state-of-the-art results on the majority of bio-medical tasks. This encoder is pre-trained from scratch with a domain-specific tokenizer on a corpus based on PubMed.

**Advancements in COP.** Naik et al. (2021) augments a PubmedBERT model with document retrieval from a PubMed knowledge base. Grundmann et al. (2022) and Winter et al. (2022) incorporate additional modalities in the form of support sets of ICD codes from prior admissions, and knowledge graph completion tasks respectively. Papaioannou et al. (2022) present knowledge transfer strategies to improve performance for low-resource clinical text datasets in different languages. They show that incorporating clinical text written in multiple languages can complement clinical knowledge missing in smaller datasets, especially for non-frequent diagnoses. Deznabi et al. (2021) augment the text modality with time-series data to improve predictions for in-hospital mortality. van Aken et al. (2022) enhances a Pubmed-BERT encoder with a prototypical network to not only improve prediction results, but also increase the explainability of predictions.

## 3 COP-IV Tasks

We revisit the task creation process of van Aken et al. (2021) and update it for the MIMIC-IV data.

### 3.1 MIMIC-IV: Data preparation

**Creation of admission notes.** The electronic health records (EHR) available in MIMIC are all associated with medical discharge summaries about the visit of a patient to the hospital. We follow the same pre-processing as in (van Aken et al., 2021), adapted to MIMIC-IV. Hence, we

---

[1]available on https://huggingface.co/
[2]https://github.com/DATEXIS/ClinicalOutcomePrediction-IV

| | mean (words/note) | std (words/note) | mean (sent/note) | std (sent/note) | total notes |
|---|---|---|---|---|---|
| COP-III-ICU | 396.3 | 233.3 | 32.5 | 23.1 | **48,745** |
| COP-IV-ICU | 495.6 | 236.7 | 26.9 | 16.1 | **59,056** |
| COP-IV-ED | 523.9 | 265.2 | 28.5 | 17.5 | **269,573** |

Table 1: *COP-III* vs *COP-IV* admission notes details. *COP-III* is based on MIMIC-III, while *COP-IV* is based on MIMIC-IV. The amount of available notes in the ICU increases. ED is not available in MIMIC-III.

keep specific sections in the discharge summaries that are known at admission time, such as: *Chief complaint*, *(History of) Present illness*, *Medical history*, *Admission medications*, *Allergies*, *Physical exam*, *Family history*, and *Social history*. An admission note acts as an input for all tasks; in Figure 1 we present an example. Table 1 demonstrates a comparison of the statistics of admission notes in *COP-III* and *COP-IV*. We observe that the resulting ICU data for *COP-IV* contains 21% more admission notes compared to *COP-III*. In sharp contrast, *COP-IV* offers an additional 269,573 admission notes in the novel ED split. We also remark that for *COP-IV* the average length of an admission note increases, while the number of sentences decreases.

Additionally, note that the clearest difference between MIMIC-III and MIMIC-IV in terms of style is the anonymization scheme. MIMIC-III follows HIPAA[3] for anonymization and identifiable entities are replaced with random identifiers and an indication of the previous content. In contrast, MIMIC-IV replaces all identifiable markers with three underscores: "___"(Johnson et al., 2023). We follow van Aken et al. (2021) and do not mask the de-identified tokens and consider them as part of the admission note.

**ICD-10 label space.** For the diagnoses and procedure prediction tasks in *COP-III*, the labels are ICD-9 codes. Since MIMIC-IV includes admission notes annotated with ICD-10 codes, for these specific tasks in *COP-IV* we choose to predict only for this newer ICD version. We do this only for the diagnoses and procedures prediction tasks since the remaining tasks are independent of the ICD standard.

### 3.2 Outcome prediction tasks

**Patient routing (PR).** We introduce a novel task to *COP-IV*. We construct this task by lever-

---
[3]Health Insurance Portability and Accountability Act

aging routing information for patients accessible in MIMIC-IV, which details patient transfers between different units within the hospital. In the patient routing task, we predict the first hospital unit a patient is transferred to upon admission to the emergency department. Note that we only focus on the first transfer of a patient out of the emergency department, since we predict at the time of admission. Furthermore, we consolidate the labels for the patient's routing information that refer to the same class but differ in their naming. For instance, there are several specific hospital section labels related to surgical procedures, which we group together into *surgery*. This process results in a total of 18 classes (Table 2), making this a multi-class classification task.

| Patient Routing Prediction | | |
|---|---|---|
| | Classes | Number of Samples |
| COP-IV-ED | 18 | 328,589 |

Table 2: Novel *patient routing* prediction task summary.

**Diagnoses prediction (DIA).** The diagnoses prediction task in *COP-IV* involves mapping admission notes to the ICD-10 coding standard. Similar to van Aken et al. (2021), we don't capture the full granularity of ICD-10, and limit ourselves to three-digit codes. This significantly reduces label scarcity, but still retains a relevant level of detail since the codes are organized hierarchically (Choi et al., 2017). As we show in Table 3, the label space grows in size significantly compared to the old version *COP-III*. We apply a multi-label stratified sampling approach (Sechidis et al., 2011) to split the dataset into train/val/test. This ensures that all codes appear in the training set at least once. Furthermore, we restrict multiple admissions for a single patient to be present in the same split, to prevent potential data leakage during training. Diagnoses prediction is a multi-label classification task.

**Procedures prediction (PRO).** The procedures prediction task in *COP-IV* also involves mapping admission notes to ICD-10. In contrast to the diagnoses prediction task, instead of using only the first 3 digits, we use the first 4 digits. This is due to the differences in hierarchy between the diagnoses and procedure codes in ICD-10. Table 4 contains

Figure 1: Clinical Outcome Prediction: Given an EHR textual description of a patient admission(left) this task involves determining outcomes (right) such as diagnoses, procedures, hospital section, length of stay, and mortality at discharge.

| Diagnoses Outcome Prediction | | | | |
|---|---|---|---|---|
| | **Total** | Train | Test | Val |
| COP-III-ICU | **1,266** | 1,201 | 1,031 | 906 |
| COP-IV-ICU | **1,447** | 1,447 | 943 | 943 |
| COP-IV-ED | **1,617** | 1,617 | 1,207 | 1,198 |

Table 3: Diagnoses code statistics for *COP-III* vs *COP-IV*. Note that the labels in the *COP-III* diagnoses task are ICD-9 codes and in *COP-IV* these are ICD-10 codes. The label space grows significantly for both splits, ED and ICU.

| Procedures Outcome Prediction | | | | |
|---|---|---|---|---|
| | **Total** | Train | Test | Val |
| COP-III-ICU | **711** | 672 | 563 | 476 |
| COP-IV-ICU | **2,956** | 2,956 | 761 | 756 |
| COP-IV-ED | **4,137** | 4,137 | 1,242 | 1,344 |

Table 4: Procedures code statistics for *COP-III* vs *COP-IV*. The label space grows significantly due to the adoption of ICD-10 in *COP-IV*

a summary of the code distributions for the task. Since in the ICD-10 coding standard there are 19 times more procedure codes than ICD-9[4], the total number of codes increases drastically across the ICU and the ED split. We apply the stratified sampling strategy that we use for the diagnoses outcome prediction task. Procedures prediction is also a multi-label classification task.

**Length-of-stay prediction (LOS).** Predicting the length of a patient's stay for a visit is beneficial for medical facilities to allocate resources accordingly. As in (van Aken et al., 2021), the length of an ICU stay is defined as the number of days between the admission and discharge of a patient. Unlike van Aken et al. (2021) we focus specifically on the length of a stay of a patient in the ICU, since factors beyond the state of a patient like occupied beds, medical professionals availability, etc.

could determine the stay. This information is available in MIMIC-IV and we use the same 4 classes as in *COP-III*: *Under 3 days*, *3 to 7 days*, *1 week to 2 weeks*, and *more than 2 weeks*. We validate these modifications to the task with medical professionals and do not create this task for the ED split. As shown in Table 5, the stay of patients considered in *COP-IV* shifts significantly due to the focus of the stay in the ICU. The majority class is now (*Under 3 days*). Length-of-stay prediction is a multi-class classification task.

| Length-of-stay (in days) | | | | |
|---|---|---|---|---|
| | ≤ 3 | > 3 & ≤ 7 | > 7 & ≤ 14 | > 14 |
| COP-III-ICU | 5,596 | 16,134 | 13,391 | 8,488 |
| COP-IV-ICU | 41,285 | 11,840 | 3,986 | 1,945 |

Table 5: *Length-of-stay* prediction task for *COP-III* & *COP-IV*. The length of a stay is measured in days. We observe a shift in the class distribution between version III and IV. This task is not applicable to the ED split.

**In-hospital mortality prediction.** Since a medical professional writes a discharge summary after the visit of a patient, admission sections may contain explicit references to their death. van Aken et al. (2021) applied pattern matching to remove such admission notes. However, in our attempt to replicate this preprocessing method, we found that a rule-based approach to detecting these cases is not reliable. We trained PubMedBERT following this approach; this led to extremely high scores in both AUROC and PR-AUC. Upon closer examination, we still encounter additional patterns (e.g. cessation, passed) that made the decease of a patient explicit. Since we cannot guarantee exhaustive filtering to remove admission notes with such fragments for the MIMIC-IV data, we omit van Aken et al.'s (2021) in-hospital mortality prediction task in *COP-IV*.

## 4 Experiments

We fine-tune all models in all outcome prediction tasks on both MIMIC-IV splits, except for the *LOS* and *PR* tasks. These tasks are exclusive to the ICU and ED split as mentioned in Section 3. We report performance in *AUROC-macro* as well as in *PR-AUC*. In contrast to van Aken et al. (2021), we include PR-AUC as an additional metric.

While the AUROC provides insight into performance for the majority of the patients, the PR-AUC provides a more balanced view, since it emphasizes the performance of labels that are less frequent in the data.

For comparability, we evaluate all *COP-IV* tasks with the encoder models used in (van Aken et al., 2021), namely BioBERT, CORe, ClinicalBERT, and DischargeBERT. Additionally, we extend this evaluation to PubMedBERT. We conduct a Hyper-Parameter-Optimization (HPO) on PubMedBERT for all tasks for the learning rate and warmup steps using *ray* (Liaw et al., 2018) and (Bergstra et al., 2013). We use the resulting hyperparameters in all experiments. We use early stopping on AUROC with a patience of 5 epochs as in van Aken et al. (2021). We keep a consistent batch size of 50 for all tasks and models. For every experiment, we use a single A100 40GB GPU.

## 5 Results

We present all experimental results in Table 6.

**Overall performance.** PubMedBERT outperforms all models across all tasks. BioBERT is the second best performing model, followed by CORe. ClincalBERT and DischargeBERT are the worst performing models.

**Domain-specific tokenizer.** PubMedBERT is the only model in our work that uses a domain-specific tokenizer. We argue that this is one of the reasons why it is the top-performing model across all tasks. Notably, the average tokenized admission note in MIMIC-IV is longer than 512 tokens. Thus exceding the maximum sequence length for BERT-like models. Therefore, the context window that PubMedBERT processes per admission note contains more information on average compared to the other models.

**Pre-training on MIMIC does not bring benefits.** PubMedBERT and BioBERT are pre-trained on PubMed. They have not explicitly seen any MIMIC discharge summaries during the pre-training. In contrast, CORe, ClinicalBERT, and DischargeBERT incorporate MIMIC-III data into their training routine, thus exposing the parameters to specific details, writing style, and anonymization scheme. The results suggest that the models do not benefit from pre-training on MIMIC-III. This is highlighted by the fact that BioBERT has a very similar performance. Thus, reinforcing the idea that the domain-specific tokenizer has a much greater impact on the performance of these tasks.

**Patient routing.** All models achieve high scores for AUROC. In contrast, the results in PR-AUC indicate that all models have difficulties with capturing the hospital units where transfers occur less often. Similar to other tasks, PubMedBert outperforms all other models.

## 5.1 Performance comparison of CORe on MIMIC-III and MIMIC-IV

To validate that our adaptation of the *COP* tasks to the MIMIC-IV dataset is done correctly, we compare the performance of the CORe model on *COP-III* and *COP-IV*. For *COP-III* we use scores from van Aken et al. (2021) and for *COP-IV* we take the results of the CORe[5] model from our evaluation on the respective task in *COP-IV*. We present this comparison in Table 7. Since the ED split was not available in MIMIC-III, we only compare the

---

[5] https://huggingface.co/DATEXIS/ CORe-clinical-outcome-biobert-v1, accessed 28.02.24

| Split | Task | PR | | DIA | | PRO | | LOS | |
|---|---|---|---|---|---|---|---|---|---|
| | Model | AUROC | PR-AUC | AUROC | PR-AUC | AUROC | PR-AUC | AUROC | PR-AUC |
| ED | BioBERT | 93.83 | 59.33 | 85.86 | 14.77 | 92.87 | 19.32 | - | - |
| | CORe | 93.85 | 59.55 | 85.46 | 14.54 | 93.57 | 19.70 | - | - |
| | DischargeBERT | 93.87 | 59.69 | 84.83 | 14.29 | 92.93 | 19.02 | - | - |
| | ClinicalBERT | 93.85 | 59.19 | 84.73 | 14.05 | 93.18 | 18.74 | - | - |
| | PubMedBERT | **94.28** | **61.44** | **86.86** | **17.24** | **93.64** | **21.62** | - | - |
| ICU | BioBERT | - | - | 78.71 | 13.02 | 86.32 | 17.44 | 70.89 | 36.06 |
| | CORe | - | - | 78.06 | 13.05 | 85.38 | 16.10 | 71.39 | 36.49 |
| | DischargeBERT | - | - | 77.76 | 12.30 | 85.25 | 16.01 | 70.00 | 35.70 |
| | ClinicalBERT | - | - | 77.02 | 12.58 | 84.62 | 14.86 | 70.18 | 35.58 |
| | PubMedBERT | - | - | **79.70** | **15.55** | **87.21** | **18.43** | **71.82** | **36.87** |

Table 6: Results of the models for all outcome prediction tasks. Metrics are macro averaged and scores are in %. PubMedBERT is the best performing model for all *COP-IV* tasks. We observe a big gap between AUROC and PR-AUC, signaling the challenges of the long-tail distribution of labels in MIMIC.

| | DIA | PRO | LOS |
|---|---|---|---|
| CORe COP-III | 83.39 | 87.15 | 72.53 |
| CORe COP-IV | 78.06 | 85.38 | 71.39 |

Table 7: Comparison of the CORe model's AUROC-macro performance in *COP-III* as reported in (van Aken et al., 2021) and *COP-IV*. The scores are in %. Given the non-existence of the ED split in version III, we compare ICU only. The tasks in *COP-IV* are more challenging, the pre-training on MIMIC-III does not transfer positively to MIMIC-IV.

tasks that relate to ICU data. This also excludes the patient routing task.

**Diagnoses and procedures outcome prediction.** *COP-III* and *COP-IV* have different label spaces for diagnoses and procedures. We use ICD-10, whereas *COP-III* uses ICD-9. van Aken et al. (2021) reports better performance for both tasks. We argue that this performance gap might be due to the larger code space of ICD-10 compared to ICD-9 (Cartwright, 2013). Additionally, since *COP-IV* uses only ICD-10 codes, we are limited to a fraction of the total amount of summaries available in MIMIC-IV for the ICU split. Roughly 60% of admission notes in this split are annotated with the ICD-9 standard, hence this results in significantly fewer notes for training in *COP-IV* than in *COP-III*.

**Length-of-stay.** The similar scores for the CORe model in *COP-III* and *COP-IV* in Table 7 indicate that the length-of-stay task is still chal-

lenging, despite the modification aimed at focusing on the ICU stay. As previously noted, this leads to a shift of the label distribution, with the majority of patients experiencing shorter stays compared to the *COP-III* task. We argue that this shift in the distribution of the labels could be a factor explaining the lower scores for the task in *COP-IV*. Additional challenges at predicting the length of stay of a patient come from factors such as *employment* or *marital status* which may not be mentioned in a clinical admission note (Khosravizadeh et al., 2016).

## 6 Discussion & Future Work

### 6.1 Multi-label outcome prediction

The performance reported in Table 6, shows that the AUROC and especially the PR-AUC metric for the DIA and PRO tasks have a large room for improvement.

**Critical long-tail.** In Figure 2 we present the label distribution for the complete ED split in MIMIC-IV. It is worth noting that only 100 labels (6% of all labels) are annotated in approximately 67% of the data, whereas the remaining 1,517 labels (94% of all labels) are distributed among the remaining 33% of the samples. We observe the same behavior in the ICU split. We expand the evaluation of PR-AUC of PubMedBERT for class groups depending on their frequency. Figure 3 demonstrates that the model achieves poor PR-AUC performance in the tail of the distribution and improves towards the head. This behavior in PR-AUC emphasizes a weakness of current methods

Figure 2: ICD-10 code distribution for the MIMIC-IV ED split. Each one of the 3 colors indicates 33.3% of total samples highlighting a pronounced long tail.



Figure 3: PR-AUC in % measured on groups of labels depending on their frequency in the data. Performance in the long tail is generally poor while it improves greatly for the more frequent labels.

since the majority of the labels reside in the tail.

**Label-space.** The larger code space in ICD-10 in comparison to ICD-9 further exacerbates the class imbalance present in the multi-label outcome prediction tasks (DIA & PRO).

**Annotation** Moreover, labels in MIMIC exhibit annotation inconsistencies; in practice the most frequent labels are under-annotated (up to 35%) (Searle et al., 2020). Therefore, some correct predictions made by models will conflict with an incomplete ground truth.

## 6.2 Qualitative analysis on Patient routing

For the novel patient routing task, we conduct an additional analysis on diversity and identify potential gaps for different populations. Next, we further discuss the difference in performance that we

observe in hospital care units. In Figure 4 we disaggregate the PR-AUC for variables such as gender and marital status, as well as admission type and care unit.

**Demographic variables.** We observe that predictions for male patients are worse by a significant margin. A possible reason could be the additional amount of time spent by women on average for physical exams and patient questions when visiting a doctor (Tabenkin et al., 2004), thus producing more relevant information during the anamnesis. This may result in richer admission notes for women. The *marital status* shows an impact on widowed patients. The average patient is 78 years old, which is 18 years older when compared to the other categories. Given that the age of patients has an impact on other tasks (van Aken et al., 2021; Khosravizadeh et al., 2016), we argue that it has an impact on patient routing as well. For all other classes, the marital status does not seem to influence the outcome.

**Admission type** PubMedBERT achieves its best performance with admissions that come through physician referrals. Such referrals may contain relevant information to route patients to the corresponding care unit. Walk-ins and Emergency Room (ER) admissions may prioritize immediate care over EHR documentation. Therefore, we argue that in such cases, routing information might be incomplete.

**Performance of care units.** We observe that performance is not directly coupled to the class distribution. In Figure 4 bottom right, we present the PR-AUC for each care unit, sorting them (from left to right) by the number of occurrences in the data. For instance, *psychiatry* (dark green) and *obstetrics* (dark orange), where PR-AUC is significantly above the average, are units that are less present in the data. We argue that for this task performance is determined by the specificity in the admission notes relevant to each care unit and less so by the class frequency. The fact that the *observation* (pink) category is the worst performing reflects the inherent uncertainty of this care unit. We argue that since the symptomatology is not as clear as for other care units (pregnancy in obstetrics), models have more difficulties in routing the patient to the right care unit.

Figure 4: PR-AUC of the patient routing task disaggregated by **Top** demographic variables: Gender and Marital status, **Bottom** Admission type and Care units. A large gap between genders exists. Physician referrals route best. Frequency and marital status of classes are not directly coupled with prediction performance.

## 6.3 Future Work

Our work aims to be a resource for future research in clinical outcome prediction. We propose future work directions as follows:

**ICD code imbalance.** We see a very pronounced room for improvement in PR-AUC performance due to the distribution of the labels in the data. We believe that models designed to tackle this premise are needed since it's an inherent feature in the distribution of real-world clinical data. This could be accomplished with novel architectures beyond transformers, or further strategies to integrate complementary knowledge.

**Label inconsistency.** MIMIC is the best publicly available EHR data and contains annotation deficiencies. We believe that a great effort towards consistent labeling is needed. Potential avenues of data augmentation could come from leveraging generative methods to rephrase and augment existent verified high-quality data.

**Evaluation on other datasets.** Much of the prior research in clinical NLP has centered around MIMIC. However, evaluating on alternative datasets is crucial. We noticed in our *COP-IV* experiments how models did not benefit from pretraining on MIMIC-III. We believe that these signs of overfitting could be mitigated with broader evaluations using clinical text sourced in different clinics, specialties, and languages.

**Multimodal patient representation.** Although most modalities relevant to medical practitioners can be expressed in natural language, there are numerous additional modalities available not only in MIMIC but also in other domain datasets. We believe that enriching the textual representations of transformers with multi-modal data could be beneficial for the outcome prediction tasks.

**Novel outcome prediction tasks.** In practice, outcome prediction consists of a very broad set of possible tasks. Our novel patient routing task is just one example. We expect that additional tasks would provide valuable insights into the strengths and weaknesses of models employed in real-life clinical settings.

## 7 Conclusion

In this work, we introduce *COP-IV*, a clinical outcome prediction set of tasks based on MIMIC-IV, which updates *COP-III*. In addition, we introduce the novel task of *patient routing* at admission time to clinical outcome prediction. We evaluate qualitatively this task for various patient demographics, as well as hospital care units. We explain in detail our preprocessing approach to reproduce the *COP-IV* tasks. Furthermore, we present a comprehensive evaluation of several bio-medical encoder models and discuss their weaknesses, as well as challenges such as the pronounced class imbalance. Moreover, we give relevant insights into data distribution shifts between *COP-III* and *COP-IV*. Lastly, we propose future research directions for clinical outcome prediction. We release our source code to reproduce the data for our benchmark, experiments, and results.

## Acknowledgements

# References

Emily Alsentzer, John R. Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew B. A. McDermott. 2019. Publicly available clinical bert embeddings.

Manuela Benary, Xing David Wang, Max Schmidt, Dominik Soll, Georg Hilfenhaus, Mani Nassir, Christian Sigler, Maren Knödler, Ulrich Keller, Dieter Beule, Ulrich Keilholz, Ulf Leser, and Damian T. Rieke. 2023. Leveraging Large Language Models for Decision Support in Personalized Oncology. *JAMA Network Open*, 6(11):e2343689–e2343689.

J Bergstra, D Yamins, and D D Cox. 2013. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *TProc. of the 30th International Conference on Machine Learning (ICML 2013*.

Donna J Cartwright. 2013. ICD-9-CM to ICD-10-CM codes: What? why? how? *Adv. Wound Care (New Rochelle)*, 2(10):588–592.

Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. 2017. Generating multi-label discrete patient records using generative adversarial networks. In *Proceedings of the 2nd Machine Learning for Healthcare Conference*, volume 68 of *Proceedings of Machine Learning Research*, pages 286–305. PMLR.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Iman Deznabi, Mohit Iyyer, and Madalina Fiterau. 2021. Predicting in-hospital mortality by combining clinical notes with time-series data. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4026–4031, Online. Association for Computational Linguistics.

Paul Grundmann, Tom Oberhauser, Felix Gers, and Alexander Löser. 2022. Attention networks for augmenting clinical text with support sets for diagnosis prediction. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4765–4775, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.

Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J. Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, Li-wei H. Lehman, Leo A. Celi, and Roger G. Mark. 2023. Mimic-iv, a freely accessible electronic health record dataset. *Scientific Data*, 10(1):1.

Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3(1):160035.

Omid Khosravizadeh, Soudabeh Vatankhah, Peivand Bastani, Rohollah Kalhor, Samira Alirezaei, and Farzane Doosty. 2016. Factors affecting length of stay in teaching hospitals of a middle-income country. *Electron. Physician*, 8(10):3042–3047.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E. Gonzalez, and Ion Stoica. 2018. Tune: A research platform for distributed model selection and training.

Daniel McDuff, Mike Schaekermann, 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, Ewa Dominowska, Juraj Gottweis, Joelle Barral, Katherine Chou, Greg S Corrado, Yossi Matias, Jake Sunshine, Alan Karthikesalingam, and Vivek Natarajan. 2023. Towards accurate differential diagnosis with large language models.

Aakanksha Naik, Sravanthi Parasa, Sergey Feldman, Lucy Lu Wang, and Tom Hope. 2021. Literature-augmented clinical outcome prediction. *CoRR*, abs/2111.08374.

Jens-Michalis Papaioannou, Paul Grundmann, Betty van Aken, Athanasios Samaras, Ilias Kyparissidis, George Giannakoulas, Felix Gers, and Alexander Loeser. 2022. Cross-lingual knowledge transfer for clinical phenotyping. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 900–909, Marseille, France. European Language Resources Association.

Thomas Searle, Zina Ibrahim, and Richard Dobson. 2020. Experimental evaluation and development of a silver-standard for the MIMIC-III clinical coding dataset. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*, pages 76–85, Online. Association for Computational Linguistics.

Konstantinos Sechidis, Grigorios Tsoumakas, and Ioannis Vlahavas. 2011. On the stratification of multi-label data. In *Machine Learning and Knowledge Discovery in Databases*, pages 145–158, Berlin, Heidelberg. Springer Berlin Heidelberg.

Hava Tabenkin, Meredith A Goodwin, Stephen J Zyzanski, Kurt C Stange, and Jack H Medalie. 2004. Gender differences in time spent during direct observation of doctor-patient encounters. *J. Womens. Health (Larchmt)*, 13(3):341–349.

Betty van Aken, Jens-Michalis Papaioannou, Manuel Mayrdorfer, Klemens Budde, Felix Gers, and Alexander Loeser. 2021. Clinical outcome prediction from admission notes using self-supervised knowledge integration. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 881–893, Online. Association for Computational Linguistics.

Betty van Aken, Jens-Michalis Papaioannou, Marcel Naik, Georgios Eleftheriadis, Wolfgang Nejdl, Felix Gers, and Alexander Loeser. 2022. This patient looks like that patient: Prototypical networks for interpretable diagnosis prediction from clinical text. 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 172–184, Online only. Association for Computational Linguistics.

Benjamin Winter, Alexei Figueroa Rosero, Alexander Löser, Felix Alexander Gers, and Amy Siu. 2022. KIMERA: injecting domain knowledge into vacant transformer heads. In *LREC*, pages 363–373. European Language Resources Association.

# Can LLMs Correct Physicians, Yet? Investigating Effective Interaction Methods in the Medical Domain

**Burcu Sayin**[1]    **Pasquale Minervini**[2]    **Jacopo Staiano**[1]    **Andrea Passerini**[1]

[1]DISI, University of Trento, `name.surname@unitn.it`

[2]School of Informatics, University of Edinburgh, `p.minervini@ed.ac.uk`

## Abstract

We explore the potential of Large Language Models (LLMs) to assist and potentially correct physicians in medical decision-making tasks. We evaluate several LLMs, including Meditron, Llama2, and Mistral, to analyze the ability of these models to interact effectively with physicians across different scenarios. We consider questions from PubMedQA (Jin et al., 2019) and several tasks, ranging from binary (yes/no) responses to long answer generation, where the answer of the model is produced after an interaction with a physician. Our findings suggest that prompt design significantly influences the downstream accuracy of LLMs and that LLMs can provide valuable feedback to physicians, challenging incorrect diagnoses and contributing to more accurate decision-making. For example, when the physician is accurate 38% of the time, Mistral can produce the correct answer, improving accuracy up to 74% depending on the prompt being used, while Llama2 and Meditron models exhibit greater sensitivity to prompt choice. Our analysis also uncovers the challenges of ensuring that LLM-generated suggestions are pertinent and useful, emphasizing the need for further research in this area.

## 1 Introduction

Recent advancements demonstrate Large Language Models' (LLMs) effectiveness in medical AI applications, notably in diagnosis and clinical support systems (Sutton et al., 2020). Studies reveal their proficiency in answering diverse medical inquiries with high precision (Nori et al., 2023a,b; Tang et al., 2023; Nazary et al., 2024; Dai et al., 2023; Wang et al., 2023; Chen et al., 2023c; Liu et al., 2023; Liévin et al., 2023; Chen et al., 2023a,b), emphasizing the importance of tailored prompt design (Nori et al., 2023b), and advanced prompting techniques for complex tasks (Tang et al., 2023). Despite their potential, there are still challenges in deploying LLMs in the clinical domain (Salvagno et al., 2023;

Azamfirei et al., 2023; Alkaissi and McFarlane, 2023; Ji et al., 2023). Furthermore, existing works evaluate the quality of the standalone LLM, while we are interested in the setting where the LLM is supporting a human decision-maker. In many high-stakes medical scenarios, human experts (e.g., physicians) are responsible for making final decisions, and they can seek assistance from AI agents: understanding how AI systems and experts can interact is essential for ensuring their practical utility and reliability.

We aim to bridge this gap by analyzing the accuracy of LLMs in medical and clinical tasks when interacting with a domain expert (i.e., a physician). For the sake of simplicity, we consider the setting where the LLM is asked to answer a question after a domain expert verbalizes their opinion. We examine whether LLMs avoid challenging expert inputs, potentially affecting response quality. Through empirical tests, we assess LLMs' ability to rectify expert errors while maintaining collaboration, analyzing the impact of expert performance and prompt design on optimizing the performance in clinical decision-making.

Our study presents two main contributions. First, we introduce a binary PubMedQA (Jin et al., 2019) dataset featuring plausible correct and incorrect explanations generated by GPT4. Second, we highlight the importance of prompt design in enhancing LLM interactions with medical experts, showing its influence on LLMs' ability to correct physician errors, explain medical reasoning, adapt to physician input, and ultimately improve LLM performance.

## 2 Methodology

### 2.1 Prompt Design

Our analysis focuses on evaluating LLM performance in medical question-answering tasks with and without a physician answer and/or a corresponding explanation provided in the

Figure 1: Prompt design. The left figure shows the complete prompt template. We start with task instructions; while a summary is provided here as an example, detailed instructions for each use case can be found in Appendix-A. Then, we incorporate the few-shot examples, with their order varying depending on scenarios 1-4. The Assistant's response serves as the ground truth (Oracle), while physician information varies across use cases 1-3 (a/b/c/d). In the baseline case, no information from the physician is provided. Subsequently, we present the test input, where the user provides context and poses a question, followed by information from the physician depending on the use case. On the right side of the figure, detailed information is provided for few-shot example scenarios and use cases.

prompt. Given the well-known LLMs' sensitivity to prompts and the potential impact of the order of few-shot examples on output quality (Bhavya et al., 2022), we explore several in-context learning scenarios and human expert-LLM interactions.

Figure 1 illustrates our prompt template. We first explain the task instructions to the LLM (see Appendix A). Then, we present simulated conversations between the physician and the LLM, which were created by the authors (see Appendix B). The order of few-shot examples varies according to the scenario. This design aims to explore the impact of modifications in user's input and the arrangement of few-shot examples on the responses generated by the LLM. Scenarios 1-4 are structured to exhibit variability in the level of agreement or disagreement between the user and the LLM on 'yes' and 'no' responses. The prompt concludes with the test input, which includes a specific question, the context, and the physician's response.

## 2.2 Use Cases

We focus on binary classification tasks and consider the medical questions with a binary response, investigating the following experimental settings:

**Baseline** A plain question-answering (QA) setting, with no input from the physician.

**Case 1** The physician provides a binary ("yes/no") answer to the prompt question. We examine four distinct cases: (*Case 1a*): The physician is always right; (*Case 1b*): The physician is always wrong; (*Case 1c*): The physician always answers "yes"; (*Case 1d*): The physician always answers "no".

**Case 2** The physician complements the binary answer with a textual explanation. We use the GPT-4 APIs,[1] to generate plausible correct and incorrect explanations for each test example (see Appendix C). We replicate the same scenarios as in Case 1 (a/b/c/d), enriching the prompts with the physician's explanation. For instance, in *Case 2a*, the physician always provides the correct "yes/no" answer and a plausible correct explanation generated by GPT-4. In *Case 2c*, the physician always responds "yes", together with a plausible correct or incorrect explanation generated by GPT-4 depending on whether the correct answer to the question is "yes" or "no".

**Case 3** The physician provides a (binary) correct answer with a certain probability. We simulate physicians with different expertise by varying the probability $p$ of providing a correct answer, with $p \in \{70\%, 75\%, 80\%, 85\%, 90\%, 95\%\}$.

---

[1]Precisely, we used the gpt-4-32k model.

| | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis |
| **1a** | 22 | **84** | 43 | **97** | 96 | 70 | 54 | **91** | 47 | **85** | 83 | 66 |
| **1b** | 79 | 57 | **95** | 7 | 14 | **85** | 51 | 19 | **95** | 37 | 52 | **90** |
| **1c** | 38 | 70 | **75** | 66 | 71 | **80** | 49 | 70 | **77** | 62 | 70 | **82** |
| **1d** | 62 | **71** | 64 | 38 | 40 | **74** | 55 | 40 | **65** | 60 | 65 | **74** |

Table 1: Accuracy (in %) of models in Case 1.

| | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis |
| **1a** | **32** | 7 | 9 | **23** | 6 | 7 | 23 | **26** | 7 | 23 | **26** | 9 |
| **1b** | **32** | 28 | 30 | **20** | 8 | 8 | **23** | 10 | 8 | 23 | 21 | **28** |
| **1c** | **32** | 7 | 9 | **23** | 6 | 16 | 23 | **26** | 16 | 23 | **26** | 25 |
| **1d** | **32** | 28 | 9 | 20 | **28** | 8 | **23** | 10 | 8 | 23 | 21 | **28** |

Table 2: ROUGE-L scores of models in Case 1

# 3 Experimental Setup

We ran an experimental evaluation aimed at answering the following research questions: **Q1**: Can LLMs correct physicians when needed? **Q2**: Can LLMs explain the reasons behind their answers? **Q3**: Can LLMs correct physicians when they provide arguments for their answers? **Q4**: Can LLMs fed with physician answers outperform both themselves and physicians?

## 3.1 Dataset

We use the PubMedQA dataset (Jin et al., 2019), an established biomedical QA dataset sourced from PubMed abstracts. The task is to answer biomedical questions with "yes/no/maybe" considering the given PubMed abstracts. We created a binary version of the task by taking the `pubmed_qa_labeled_fold0_source` subset from the HuggingFace dataset[2], and discarding the (few) "maybe" instances, yielding 445 test examples (62% of class "yes"). We fed this binary dataset as input into GPT-4, asking it to produce plausible correct and incorrect long answers for each question so as to emulate physicians' explanations (Case 2). We made this dataset publicly available[3] and provide further details in Appendix C.

## 3.2 Models & Frameworks

We use Meditron-7B (Med) (Chen et al., 2023a,b), Llama2-7B chat (Ll2) (Touvron et al., 2023), and Mistral-7B-Instruct (Mis) (Jiang et al., 2023) models. We conduct our experiments via Harness Framework (Gao et al., 2023). Our source code is available online.[4]

# 4 Results

**A1: Prompt design affects LLM performance in correcting erroneous physician responses** Table 1 shows the remarkable influence of prompt design on the models' performances: given appropriate instructions and examples, LLMs can effectively correct physicians. For instance, in Case 1d, the physician always responds with "no" while the ground truth distribution of class "no" is just 38%: Mistral achieves significantly higher accuracy, while Llama2 and Meditron exhibit greater sensitivity to prompt changes, displaying improved performance in Scenarios 1 and 4.

**A2: LLMs *could* explain reasons behind their answers** In examining the detailed responses from each model in Case 1, we observed that the quality of Meditron's explanations exhibits minimal sensitivity to the physician's short answer (see Table 2). Llama2 model typically yields lower ROUGE-L scores in cases 1a (the physician is always right) and 1c (the physician always says "yes"). Conversely, the Mistral model consistently delivers better explanations in Scenario 4 for cases b, c, and d. Overall, results show that LLMs are capable of generating plausible explanations when the prompt is constructively framed.

**A3: LLMs exhibit different levels of dependence on physician-provided arguments.** Table 3 reveals that LLMs exhibit a tendency to rely heavily on physicians when they argument their answers, depending on the few-shot samples provided in-context: Meditron achieves 100% accuracy on Case 2a, Scenario 4, where the physician consistently provides the ground truth short answer along with a plausible correct explanation. This indicates Meditron's inclination to prioritize the last examples

| | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis |
| **2a** | 46 | **97** | 96 | 89 | **99** | 95 | 2 | **99** | 87 | **100** | 98 | 95 |
| **2b** | **99** | 0 | 56 | 30 | 0 | **50** | **93** | 2 | 60 | 4 | 5 | **51** |
| **2c** | 39 | 61 | **72** | **79** | 62 | 69 | 40 | 63 | **69** | 63 | 62 | **69** |
| **2d** | 62 | 38 | **78** | 37 | 37 | **76** | 63 | 39 | **86** | 39 | 41 | **79** |

Table 3: Accuracy (in %) of models in Case 2

| | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis |
| **2a** | 14 | 28 | **31** | **37** | 36 | 21 | 37 | **38** | 34 | 26 | **39** | 30 |
| **2b** | 14 | **24** | 19 | 13 | 13 | **22** | 24 | 15 | **26** | 13 | 15 | **40** |
| **2c** | 14 | 11 | **25** | **37** | 11 | 22 | 37 | **38** | 30 | 26 | **34** | 30 |
| **2d** | 14 | **24** | 19 | 13 | 13 | **22** | 24 | 15 | **26** | 13 | 15 | **40** |

Table 4: ROUGE-L scores of models in Case 2

| | Ground Truth | Med | Ll2 | Mis |
|---|---|---|---|---|
| 1a-3 | An association was found between sublingual varices and hypertension. Examining the lateral borders of the tongue is easily done, causes no harm and could be a valuable method for the dental profession to take active part in preventive healthcare. | Sublingual varices are associated with hypertension. | Sublingual varices are significantly associated with hypertension, as demonstrated by the observed odds ratio and p-values. The presence of sublingual varices can serve as a useful indicator for identifying individuals at risk of hypertension. | Sublingual varices are associated with hypertension. The presence of sublingual varices increases the likelihood of hypertension. However, the positive predictive value is relatively low, indicating that not all patients with sublingual varices have hypertension. |
| 2a-3 | An association was found between sublingual varices and hypertension. Examining the lateral borders of the tongue is easily done, causes no harm and could be a valuable method for the dental profession to take active part in preventive healthcare. | The presence of sublingual varices is associated with hypertension. | Sublingual varices are associated with hypertension, as evidenced by the studys findings. The presence of sublingual varices indicates a higher likelihood of hypertension, although their absence does not rule it out. Sublingual varices may serve as a useful indicator of hypertension, but they should not replace standard diagnostic methods. | Sublingual varices are associated with hypertension, but they are not a definitive diagnostic tool for hypertension. |
| 1b-4 | An association was found between sublingual varices and hypertension. Examining the lateral borders of the tongue is easily done, causes no harm and could be a valuable method for the dental profession to take active part in preventive healthcare. | Sublingual varices are not related to hypertension. | Sublingual varices are associated with hypertension, with a positive predictive value of 0.5 and a negative predictive value of 0.80. | Sublingual varices are a reliable indicator of hypertension, with a positive predictive value of 0.5 and a negative predictive value of 0.80. |
| 2b-4 | An association was found between sublingual varices and hypertension. Examining the lateral borders of the tongue is easily done, causes no harm and could be a valuable method for the dental profession to take active part in preventive healthcare. | The study found a correlation between sublingual varices and hypertension, but it is possible to argue that this correlation does not necessarily imply causation. For instance, it could be hypothesized that both conditions are independently associated with another common factor such as age or lifestyle habits, rather than sublingual varices being directly linked to hypertension. This alternative interpretation, while not supported by the data, suggests that the presence of sublingual varices does not necessarily correspond to hypertension. | The study found a correlation between sublingual varices and hypertension, but it is possible that this correlation is due to a third factor rather than a direct causal relationship between the two conditions. | Sublingual varices and hypertension are correlated, but it is unclear whether this association is causal. Other factors, such as age or lifestyle habits, may be responsible for both conditions. |

Table 5: Example of generated texts for each model on different experimental settings.

in the prompt, as evidenced by its superior performance in Scenario 2 and Scenario 4. Conversely, in Case 2b, where the physician consistently offers the opposite of the ground truth short answer and a plausible incorrect explanation, Meditron exhibits better performance in Scenario 1 and Scenario 3. Notably, Meditron learns to contradict the physician in Scenario 1 and Scenario 3 for Case 2c and Case 2d, while it learns to agree with the physician in Scenario 2 and Scenario 4. Another noteworthy observation is that LLama2 tends to over-rely on the physician across all cases and scenarios when the physician provides an argument for their answer. In contrast, Mistral demonstrates a more robust performance than Meditron and LLama2 and appears the least impacted by prompt variations, showcasing over 75% accuracy in Case 2d across every scenario. This suggests its ability to effectively correct physicians when they provide an incorrect answer and an argument.

Table 4 presents the ROUGE-L scores for the models in Case 2, showing that both Llama2 and Mistral generate plausible and more extensive explanations when the prompt includes physician's opinion (see Table 4 and App. D-Table 7). Conversely, Meditron appears to excessively depend on the physician's input, significantly impacting the quality of its explanations. Table 5 illustrates

this with an example question and the extended responses from each model. Meditron tends to alter its explanations in response to the physician's input, while Llama2 and Mistral exhibit greater consistency, offering reasonable explanations regardless of the physician's stance.

**A4: LLMs improve with expert answers but fail to outperform them** Table 6 presents the results for Case 3. Interestingly, the baseline performance of the models remains relatively consistent across different scenarios. Consistent with our observations from Case 1 and Case 2, trends in Case 3 are discernible. Meditron exhibits enhanced performance in Scenario 2 and Scenario 4, yet it surpasses its baseline performance solely in Scenario 2 when the physician achieves an accuracy of over 80%. LLama2 surpasses its baseline in all scenarios when the physician attains an accuracy exceeding 85%. In contrast, Mistral demonstrates poor performance

| | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis |
| **Baseline** | 81 | 80 | **84** | 83 | 81 | 84 | 85 | 79 | 84 | **84** | 79 | **84** |
| **Phy_70** | 40 | 75 | 58 | 70 | 71 | 74 | 55 | 69 | 61 | 71 | 75 | 73 |
| **Phy_75** | 35 | 77 | 58 | 74 | 75 | 74 | 54 | 73 | 61 | 71 | 74 | 72 |
| **Phy_80** | 34 | 80 | 55 | 79 | 80 | 74 | 51 | 76 | 56 | 79 | 78 | 72 |
| **Phy_85** | 28 | 80 | 52 | 85 | 84 | 72 | 53 | 80 | 55 | 80 | 78 | 70 |
| **Phy_90** | 28 | 80 | 49 | 88 | 87 | 71 | 54 | 83 | 52 | 79 | 80 | 69 |
| **Phy_95** | 24 | 82 | 46 | **92** | **92** | 71 | 53 | **87** | 49 | 82 | 81 | 67 |

Table 6: Case 3 - Accuracy of 7B models

in Case 3, being notably influenced by the physician's answer in each scenario. Overall, while these 7B models, when fed with physician answers, show improved performance over their baseline, they do not outperform the physicians themselves. We further investigated if the 70B version of the models fed with physician answers could outperform both alone, obtaining even worse results when employing the same prompts (see App. E-Table 8). This indicates that larger models do not necessarily yield better performance; indeed, Gramopadhye et al. (2024) recently showed how the LLama2-70B model achieved less than 55% accuracy on the MEDQA dataset (Jin et al., 2021), another medical question answering benchmark featuring questions with multiple options. The reasonable hypothesis that prompt modifications might boost the performance of 70B models falls outside the scope of this work.

## 5 Conclusion and Future Work

Our experimental results reveal several key insights. Firstly, prompt design significantly impacts LLM performance, with models demonstrating sensitivity to prompt variations yet effectively correcting erroneous physician responses with appropriate instructions and examples. For instance, Mistral achieved robust accuracy across all scenarios in Case 1d. Secondly, LLMs exhibit the ability to explain their answers under the condition that the prompt used is carefully designed. Thirdly, LLMs tend to rely on physicians when they provide arguments for their answers and are particularly influenced by the order of few shot examples. Meditron is highly affected by prompt variations, while LLama2 tends to over-rely on the physician. Mistral demonstrates robust performance, indicating resilience to prompt variations. Finally, in Case 3, while Meditron and LLama2 surpass their baselines in specific scenarios, Mistral's performance is notably influenced by the physician's answer. Larger 70B models do not guarantee improved performance, highlighting the importance of prompt design and the need for further investigation.

## 6 Limitations

A limitation of our study is the use of GPT-4 to simulate plausible correct and incorrect responses to the questions, to complement the ground-truth ones contained in the PubMedQA dataset. This choice is justified by recent findings (Tan and Jiang, 2023) highlighting the effectiveness of LLMs as generative reasoners capable of modeling user behavior and simulating their opinions/preferences in human-LLM interactions. Nonetheless, real-world experiments involving interactions with physicians should be planned to corroborate and strengthen the results found in this paper.

A second limitation is that this work is not providing solutions to the problems being raised. Indeed, the main goal of the work is raising awareness on the limitations of current open-source LLMs for medical decision support. We hope that these insights will encourage further research aimed to address these limitations.

## References

Hussam Alkaissi and Samy I. McFarlane. 2023. Artificial Hallucinations in ChatGPT: Implications in Scientific Writing. *Cureus*, 15(2):e35179.

Razvan Azamfirei, Sapna R. Kudchadkar, and James Fackler. 2023. Large language models and the perils of their hallucinations. *Critical Care*, 27(1):1–2.

Bhavya Bhavya, Jinjun Xiong, and ChengXiang Zhai. 2022. Analogy generation by prompting large language models: A case study of InstructGPT. In *Proceedings of the 15th International Conference on Natural Language Generation*, pages 298–312, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.

Zeming Chen, Alejandro Hernández-Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023a. Meditron-70b: Scaling medical pretraining for large language models. *arXiv*, abs/2311.16079.

Zeming Chen, Alejandro Hernández-Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023b. Meditron-70b: Scaling medical pretraining for large language models.

Zhiyu Chen, Yujie Lu, and William Wang. 2023c. Empowering psychotherapy with large language models: Cognitive distortion detection through diagnosis of thought prompting. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4295–4304, Singapore. Association for Computational Linguistics.

Haixing Dai, Zheng Liu, Wenxiong Liao, Xiaoke Huang, Zihao Wu, Lin Zhao, Wei Liu, Ninghao Liu, Sheng Li, Dajiang Zhu, Hongmin Cai, Quanzheng Li, Dinggang Shen, Tianming Liu, and Xiang Li. 2023. Chataug: Leveraging chatgpt for text data augmentation. *ArXiv*, abs/2302.13007.

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, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation.

Ojas Gramopadhye, Saeel Sandeep Nachane, Prateek Chanda, Ganesh Ramakrishnan, Kshitij Sharad Jadhav, Yatin Nandwani, Dinesh Raghu, and Sachindra Joshi. 2024. Few shot chain-of-thought driven reasoning to prompt llms for open ended medical question answering. *arXiv*, abs/2403.04890.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):248:1–248:38.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *arXiv*, abs/2310.06825.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research 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 2567–2577, Hong Kong, China. Association for Computational Linguistics.

Xiao Liu, Da Yin, Chen Zhang, Yansong Feng, and Dongyan Zhao. 2023. The magic of IF: Investigating causal reasoning abilities in large language models of code. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9009–9022, Toronto, Canada. Association for Computational Linguistics.

Valentin Liévin, Christoffer Egeberg Hother, and Ole Winther. 2023. Can large language models reason about medical questions? *ArXiv*, abs/2207.08143.

Fatemeh Nazary, Yashar Deldjoo, and Tommaso Di Noia. 2024. ChatGPT-HealthPrompt. Harnessing the power of xai in prompt-based healthcare decision support using ChatGPT. In *Artificial Intelligence. ECAI 2023 International Workshops*, pages 382–397, Cham. Springer Nature Switzerland.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023a. Capabilities of gpt-4 on medical challenge problems. *arXiv*, abs/2303.13375.

Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoifung Poon, Tao Qin, Naoto Usuyama, Chris White, and Eric Horvitz. 2023b. Can generalist foundation models outcompete special-purpose tuning? Case study in medicine. *ArXiv*, abs/2311.16452.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott

Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. *arXiv*, abs/2303.08774.

M. Salvagno, F.S. Taccone, A.G. Gerli, and ChatGPT. 2023. Can artificial intelligence help for scientific writing? *Critical Care*, 27(1):75. Publisher: BioMed Central Ltd.

Reed Taylor Sutton, David Pincock, Daniel C. Baumgart, Daniel C. Sadowski, Richard N Fedorak, and Karen I. Kroeker. 2020. An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ Digital Medicine*, 3(1):17.

Zhaoxuan Tan and Meng Jiang. 2023. User modeling in the era of large language models: Current research and future directions. *arXiv*, abs/2312.11518.

Xiangru Tang, Anni Zou, Zhuosheng Zhang, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. 2023. MedAgents: Large language models as collaborators for zero-shot medical reasoning. *arXiv*, abs/2311.10537.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv*, abs/2307.09288.

Jiaqi Wang, Enze Shi, Sigang Yu, Zihao Wu, Chong Ma, Haixing Dai, Qiushi Yang, Yanqing Kang, Jinru Wu, Huawen Hu, Chenxi Yue, Haiyang Zhang, Yi-Hsueh Liu, Xiang Li, Bao Ge, Dajiang Zhu, Yixuan Yuan, Dinggang Shen, Tianming Liu, and Shu Zhang. 2023. Prompt engineering for healthcare: Methodologies and applications. *ArXiv*, abs/2304.14670.

## A    Prompt Design - Instructions for different experimental settings

---

**Instructions for baseline experiments**

You are a supportive, respectful, and truthful assistant, dedicated to providing assistance in a clinical context. Your responses must adhere to the highest standards of safety, ethics, and professional integrity. They should be free from any form of bias (e.g., racial, gender-based, socio-economic) and avoid promoting harmful, unethical, illegal, or otherwise inappropriate content. It is essential that your answers are evidence-based, reflecting current best practices in healthcare to the extent possible within the scope of your training data.

In scenarios where the input is unclear, incorrect, or lacks factual basis, kindly clarify the confusion or correct the misinformation, prioritizing educational value and accuracy. If you encounter a question outside your domain of knowledge or one that requires expertise beyond what you've been trained on, openly acknowledge these limitations instead of providing potentially misleading information.

In the dialogue that follows, you will engage in simulated conversations with a physician, hereafter referred to as 'User'. The User will present clinical scenarios, including context and a specific question. Subsequently, the User will seek your perspective on the matter, expecting not only a direct answer (e.g., 'yes' or 'no') but also a rationale for your response. As the Assistant, presumed to have expertise in clinical science and medical knowledge for the purpose of this exercise, your task is to provide a reliable answer to the User's question. Please also provide a constructive explanation, backed by evidence or established clinical guidelines whenever possible.

---

**Instructions used in Case 1**

You are a supportive, respectful, and truthful assistant, dedicated to providing assistance in a clinical context. Your responses must adhere to the highest standards of safety, ethics, and professional integrity. They should be free from any form of bias (e.g., racial, gender-based, socio-economic) and avoid promoting harmful, unethical, illegal, or otherwise inappropriate content. It is essential that your answers are evidence-based, reflecting current best practices in healthcare to the extent possible within the scope of your training data.

In scenarios where the input is unclear, incorrect, or lacks factual basis, kindly clarify the confusion or correct the misinformation, prioritizing educational value and accuracy. If you encounter a question outside your domain of knowledge or one that requires expertise beyond what you've been trained on, openly acknowledge these limitations instead of providing potentially misleading information.

In the dialogue that follows, you will engage in simulated conversations with a physician, hereafter referred to as 'User'. The User will present clinical scenarios, including context, a specific question, and his own response to the question. Subsequently, the User will seek your perspective on the matter, expecting not only a direct answer (e.g., 'yes' or 'no') but also a rationale for your response. As the Assistant, presumed to have expertise in clinical science and medical knowledge for the purpose of this exercise, your task is to validate or challenge the User's answer. Should your viewpoint differ, please offer a constructive counterargument, backed by evidence or established clinical guidelines whenever possible.

You are a supportive, respectful, and truthful assistant, dedicated to providing assistance in a clinical context. Your responses must adhere to the highest standards of safety, ethics, and professional integrity. They should be free from any form of bias (e.g., racial, gender-based, socio-economic) and avoid promoting harmful, unethical, illegal, or otherwise inappropriate content. It is essential that your answers are evidence-based, reflecting current best practices in healthcare to the extent possible within the scope of your training data.

In scenarios where the input is unclear, incorrect, or lacks factual basis, kindly clarify the confusion or correct the misinformation, prioritizing educational value and accuracy. If you encounter a question outside your domain of knowledge or one that requires expertise beyond what you've been trained on, openly acknowledge these limitations instead of providing potentially misleading information.

In the dialogue that follows, you will engage in simulated conversations with a physician, hereafter referred to as 'User'. The User will present clinical scenarios, including context, a specific question, and his own response to the question along with an explanation. Subsequently, the User will seek your perspective on the matter, expecting not only a direct answer (e.g., 'yes' or 'no') but also a rationale for your response. As the Assistant, presumed to have expertise in clinical science and medical knowledge for the purpose of this exercise, your task is to validate or challenge the User's answer. Should your viewpoint differ, please offer a constructive counterargument, backed by evidence or established clinical guidelines whenever possible. Please make sure that you generate a JSON object that contains your answer and the corresponding explanation.

## B  Prompt Design - Few shot examples

### User: *Context*: To evaluate the degree to which histologic chorioamnionitis, a frequent finding in placentas submitted for histopathologic evaluation, correlates with clinical indicators of infection in the mother. A retrospective review was performed on 52 cases with a histologic diagnosis of acute chorioamnionitis from 2,051 deliveries at University Hospital, Newark, from January 2003 to July 2003. Third-trimester placentas without histologic chorioamnionitis (n = 52) served as controls. Cases and controls were selected sequentially. Maternal medical records were reviewed for indicators of maternal infection. Histologic chorioamnionitis was significantly associated with the usage of antibiotics (p = 0.0095) and a higher mean white blood cell count (p = 0.018). The presence of 1 or more clinical indicators was significantly associated with the presence of histologic chorioamnionitis (p = 0.019)., *Question*: Does histologic chorioamnionitis correspond to clinical chorioamnionitis?

### Assistant: *Answer*: yes

*Explanation*: Histologic chorioamnionitis is a reliable indicator of infection whether or not it is clinically apparent.

### User: *Context*: Complex regional pain syndrome type I is treated symptomatically. A protective effect of vitamin C (ascorbic acid) has been reported previously. A dose-response study was designed to evaluate its effect in patients with wrist fractures.In a double-blind, prospective, multicenter trial, 416 patients with 427 wrist fractures were randomly allocated to treatment with placebo or treatment with 200, 500, or 1500 mg of vitamin C daily for fifty days. The effect of gender, age, fracture type, and cast-related complaints on the occurrence of complex regional pain syndrome was analyzed.Three hundred and seventeen patients with 328 fractures were randomized to receive vitamin C, and ninety-nine patients with ninety-nine fractures were randomized to receive a placebo. The prevalence of complex regional pain syndrome was 2.4% (eight of 328) in the vitamin C group and 10.1% (ten of ninety-nine) in the placebo group (p=0.002); all of the affected patients were elderly women. Analysis of the different doses of vitamin C showed that the prevalence of complex regional pain syndrome was 4.2% (four of ninety-six) in the 200-mg group (relative risk, 0.41; 95% confidence interval, 0.13 to 1.27), 1.8% (two of 114) in the 500-mg group (relative risk, 0.17; 95% confidence interval, 0.04 to 0.77), and 1.7% (two of 118) in the 1500-mg group (relative risk, 0.17; 95% confidence interval, 0.04 to 0.75). Early cast-related complaints predicted the development of complex regional pain syndrome (relative risk, 5.35; 95% confidence interval, 2.13 to 13.42)., *Question*: Can vitamin C prevent complex regional pain syndrome in patients with wrist fractures?

### Assistant: *Answer*: yes

*Explanation*: Vitamin C reduces the prevalence of complex regional pain syndrome after wrist fractures. A daily dose of 500 mg for fifty days is recommended.

### **User**: *Context*: Cancer of the buccal mucosa is an uncommon and aggressive neoplasm of the oral cavity. Less than 2% of patients treated for cancer of the oral cavity at Roswell Park Cancer Institute (RPCI) from 1971 to 1997 had primary buccal cancers. Because the majority of these patients did not undergo any adjuvant treatment, this group provided us with the opportunity to assess the relationship between margin status and local recurrence for both small (T1-T2) and large (T3-T4) tumors treated with surgery alone.The RPCI tumor registry database reported 104 patients who were treated for buccal carcinoma. A retrospective chart review identified 27 patients who met our criteria for a buccal mucosal primary tumor (epicenter of the mass in the buccal mucosa). There were 13 men and 14 women, ranging in age from 34 to 94 years (mean, 75). Data were collected regarding patient demographics, presenting symptoms, stage, treatment received, and outcome.All patients underwent surgical resection of their primary lesion; 21 (75%) had T1 or T2 tumors. The rate of local recurrence was 56% for the group as a whole. Patients with close or positive margins had a 66% local failure rate as compared with 52% when surgical margins were negative (greater than or equal to 5 mm from the resection margin after tissue fixation; P = ns). Among those in whom negative margins were achieved, patients with T1-T2 disease had a 40% local failure rate with surgical resection alone., *Question*: Cancer of the buccal mucosa: are margins and T-stage accurate predictors of local control?
### **Assistant**: *Answer*: no
*Explanation*: Local excision of T1 and T2 buccal mucosa cancers with pathologically negative margins had a high rate of local recurrence in our series. Low T-stage and negative margins are not adequate predictors of local control. Even early buccal tumors may benefit from adjuvant therapy to enhance local control.


### **User**: *Context*: To be able to adhere to discharge instructions after a visit to the emergency department (ED), patients should understand both the care that they received and their discharge instructions. The objective of this study is to assess, at discharge, patients' comprehension of their ED care and instructions and their awareness of deficiencies in their comprehension.We conducted structured interviews of 140 adult English-speaking patients or their primary caregivers after ED discharge in 2 health systems. Participants rated their subjective understanding of 4 domains: (1) diagnosis and cause; (2) ED care; (3) post-ED care, and (4) return instructions. We assessed patient comprehension as the degree of agreement (concordance) between patients' recall of each of these domains and information obtained from chart review. Two authors scored each case independently and discussed discrepancies before providing a final concordance rating (no concordance, minimal concordance, partial concordance, near concordance, complete concordance).Seventy-eight percent of patients demonstrated deficient comprehension (less than complete concordance) in at least 1 domain; 51% of patients, in 2 or more domains. Greater than a third of these deficiencies (34%) involved patients' understanding of post-ED care, whereas only 15% were for diagnosis and cause. The majority of patients with comprehension deficits failed to perceive them. Patients perceived difficulty with comprehension only 20% of the time when they demonstrated deficient comprehension., *Question*: Patient comprehension of emergency department care and instructions: are patients aware of when they do not understand?
### **Assistant**: *Answer*: no
*Explanation*: Many patients do not understand their ED care or their discharge instructions. Moreover, most patients appear to be unaware of their lack of understanding and report inappropriate confidence in their comprehension and recall.

### **User**: *Context*: To evaluate the degree to which histologic chorioamnionitis, a frequent finding in placentas submitted for histopathologic evaluation, correlates with clinical indicators of infection in the mother. A retrospective review was performed on 52 cases with a histologic diagnosis of acute chorioamnionitis from 2,051 deliveries at University Hospital, Newark, from January 2003 to July 2003. Third-trimester placentas without histologic chorioamnionitis (n = 52) served as controls. Cases and controls were selected sequentially. Maternal medical records were reviewed for indicators of maternal infection. Histologic chorioamnionitis was significantly associated with the usage of antibiotics (p = 0.0095) and a higher mean white blood cell count (p = 0.018). The presence of 1 or more clinical indicators was significantly associated with the presence of histologic chorioamnionitis (p = 0.019)., *Question*: Does histologic chorioamnionitis correspond to clinical chorioamnionitis?, *Answer*: yes

### **Assistant**: *Answer*: yes

*Explanation*: Histologic chorioamnionitis is a reliable indicator of infection whether or not it is clinically apparent.

### **User**: *Context*: Complex regional pain syndrome type I is treated symptomatically. A protective effect of vitamin C (ascorbic acid) has been reported previously. A dose-response study was designed to evaluate its effect in patients with wrist fractures.In a double-blind, prospective, multicenter trial, 416 patients with 427 wrist fractures were randomly allocated to treatment with placebo or treatment with 200, 500, or 1500 mg of vitamin C daily for fifty days. The effect of gender, age, fracture type, and cast-related complaints on the occurrence of complex regional pain syndrome was analyzed.Three hundred and seventeen patients with 328 fractures were randomized to receive vitamin C, and ninety-nine patients with ninety-nine fractures were randomized to receive a placebo. The prevalence of complex regional pain syndrome was 2.4% (eight of 328) in the vitamin C group and 10.1% (ten of ninety-nine) in the placebo group (p=0.002); all of the affected patients were elderly women. Analysis of the different doses of vitamin C showed that the prevalence of complex regional pain syndrome was 4.2% (four of ninety-six) in the 200-mg group (relative risk, 0.41; 95% confidence interval, 0.13 to 1.27), 1.8% (two of 114) in the 500-mg group (relative risk, 0.17; 95% confidence interval, 0.04 to 0.77), and 1.7% (two of 118) in the 1500-mg group (relative risk, 0.17; 95% confidence interval, 0.04 to 0.75). Early cast-related complaints predicted the development of complex regional pain syndrome (relative risk, 5.35; 95% confidence interval, 2.13 to 13.42)., *Question*: Can vitamin C prevent complex regional pain syndrome in patients with wrist fractures?, *Answer*: no

### **Assistant**: *Answer*: yes

*Explanation*: Vitamin C reduces the prevalence of complex regional pain syndrome after wrist fractures. A daily dose of 500 mg for fifty days is recommended.

### User: *Context*: Cancer of the buccal mucosa is an uncommon and aggressive neoplasm of the oral cavity. Less than 2% of patients treated for cancer of the oral cavity at Roswell Park Cancer Institute (RPCI) from 1971 to 1997 had primary buccal cancers. Because the majority of these patients did not undergo any adjuvant treatment, this group provided us with the opportunity to assess the relationship between margin status and local recurrence for both small (T1-T2) and large (T3-T4) tumors treated with surgery alone.The RPCI tumor registry database reported 104 patients who were treated for buccal carcinoma. A retrospective chart review identified 27 patients who met our criteria for a buccal mucosal primary tumor (epicenter of the mass in the buccal mucosa). There were 13 men and 14 women, ranging in age from 34 to 94 years (mean, 75). Data were collected regarding patient demographics, presenting symptoms, stage, treatment received, and outcome.All patients underwent surgical resection of their primary lesion; 21 (75%) had T1 or T2 tumors. The rate of local recurrence was 56% for the group as a whole. Patients with close or positive margins had a 66% local failure rate as compared with 52% when surgical margins were negative (greater than or equal to 5 mm from the resection margin after tissue fixation; P = ns). Among those in whom negative margins were achieved, patients with T1-T2 disease had a 40% local failure rate with surgical resection alone., *Question*: Cancer of the buccal mucosa: are margins and T-stage accurate predictors of local control?, *Answer*: no

### Assistant: *Answer*: no

*Explanation*: Local excision of T1 and T2 buccal mucosa cancers with pathologically negative margins had a high rate of local recurrence in our series. Low T-stage and negative margins are not adequate predictors of local control. Even early buccal tumors may benefit from adjuvant therapy to enhance local control.

### User: *Context*: To be able to adhere to discharge instructions after a visit to the emergency department (ED), patients should understand both the care that they received and their discharge instructions. The objective of this study is to assess, at discharge, patients' comprehension of their ED care and instructions and their awareness of deficiencies in their comprehension.We conducted structured interviews of 140 adult English-speaking patients or their primary caregivers after ED discharge in 2 health systems. Participants rated their subjective understanding of 4 domains: (1) diagnosis and cause; (2) ED care; (3) post-ED care, and (4) return instructions. We assessed patient comprehension as the degree of agreement (concordance) between patients' recall of each of these domains and information obtained from chart review. Two authors scored each case independently and discussed discrepancies before providing a final concordance rating (no concordance, minimal concordance, partial concordance, near concordance, complete concordance).Seventy-eight percent of patients demonstrated deficient comprehension (less than complete concordance) in at least 1 domain; 51% of patients, in 2 or more domains. Greater than a third of these deficiencies (34%) involved patients' understanding of post-ED care, whereas only 15% were for diagnosis and cause. The majority of patients with comprehension deficits failed to perceive them. Patients perceived difficulty with comprehension only 20% of the time when they demonstrated deficient comprehension., *Question*: Patient comprehension of emergency department care and instructions: are patients aware of when they do not understand?, *Answer*: yes

### Assistant: *Answer*: no

*Explanation*: Many patients do not understand their ED care or their discharge instructions. Moreover, most patients appear to be unaware of their lack of understanding and report inappropriate confidence in their comprehension and recall.

### **User**: {"*Context*": To evaluate the degree to which histologic chorioamnionitis, a frequent finding in placentas submitted for histopathologic evaluation, correlates with clinical indicators of infection in the mother. A retrospective review was performed on 52 cases with a histologic diagnosis of acute chorioamnionitis from 2,051 deliveries at University Hospital, Newark, from January 2003 to July 2003. Third-trimester placentas without histologic chorioamnionitis (n = 52) served as controls. Cases and controls were selected sequentially. Maternal medical records were reviewed for indicators of maternal infection. Histologic chorioamnionitis was significantly associated with the usage of antibiotics (p = 0.0095) and a higher mean white blood cell count (p = 0.018). The presence of 1 or more clinical indicators was significantly associated with the presence of histologic chorioamnionitis (p = 0.019)., "*Question*": Does histologic chorioamnionitis correspond to clinical chorioamnionitis?, "*Answer*": yes, "*Explanation*": The presence of histologic chorioamnionitis appears to be significantly associated with clinical indicators of infection in the mother, such as the usage of antibiotics and a higher mean white blood cell count. This suggests a correlation between histologic chorioamnionitis and clinical chorioamnionitis. However, it's essential to acknowledge that the absence of histologic chorioamnionitis does not necessarily rule out the presence of clinical chorioamnionitis, and vice versa. Clinical manifestations of chorioamnionitis can vary, and histologic examination may not always capture the full spectrum of the condition. Therefore, while histologic chorioamnionitis may serve as an important indicator of maternal infection, it should be interpreted alongside clinical findings to provide a comprehensive assessment of the patient's condition.}

### **Assistant**: {"*Answer*": yes, "*Explanation*": Histologic chorioamnionitis is a reliable indicator of infection whether or not it is clinically apparent.}

### **User**: {"*Context*": Complex regional pain syndrome type I is treated symptomatically. A protective effect of vitamin C (ascorbic acid) has been reported previously. A dose-response study was designed to evaluate its effect in patients with wrist fractures.In a double-blind, prospective, multicenter trial, 416 patients with 427 wrist fractures were randomly allocated to treatment with placebo or treatment with 200, 500, or 1500 mg of vitamin C daily for fifty days. The effect of gender, age, fracture type, and cast-related complaints on the occurrence of complex regional pain syndrome was analyzed.Three hundred and seventeen patients with 328 fractures were randomized to receive vitamin C, and ninety-nine patients with ninety-nine fractures were randomized to receive a placebo. The prevalence of complex regional pain syndrome was 2.4% (eight of 328) in the vitamin C group and 10.1% (ten of ninety-nine) in the placebo group (p=0.002); all of the affected patients were elderly women. Analysis of the different doses of vitamin C showed that the prevalence of complex regional pain syndrome was 4.2% (four of ninety-six) in the 200-mg group (relative risk, 0.41; 95% confidence interval, 0.13 to 1.27), 1.8% (two of 114) in the 500-mg group (relative risk, 0.17; 95% confidence interval, 0.04 to 0.77), and 1.7% (two of 118) in the 1500-mg group (relative risk, 0.17; 95% confidence interval, 0.04 to 0.75). Early cast-related complaints predicted the development of complex regional pain syndrome (relative risk, 5.35; 95% confidence interval, 2.13 to 13.42)., "*Question*": Can vitamin C prevent complex regional pain syndrome in patients with wrist fractures?, "*Answer*": no, "*Explanation*": Despite the evidence suggesting a significant reduction in the prevalence of complex regional pain syndrome (CRPS) in wrist fracture patients treated with vitamin C, an alternative interpretation could propose that the observed protective effect of vitamin C might not be directly related to its pharmacological properties. It could be argued, for the sake of discussion, that vitamin C's effect might be a surrogate for other health-promoting behaviors or factors associated with individuals more likely to take supplements, such as overall better health status, more conscientious health practices, or even a placebo effect driven by patients' expectations of treatment benefits. This perspective might suggest that the actual reduction in CRPS incidence among the vitamin C-treated groups could be attributed to these confounding factors rather than the biological action of vitamin C itself. Therefore, the correlation between vitamin C supplementation and decreased CRPS prevalence might not imply a causal relationship, challenging the notion that vitamin C can directly prevent CRPS in patients with wrist fractures.}

### **Assistant**: {"*Answer*": yes, "*Explanation*": Vitamin C reduces the prevalence of complex regional pain syndrome after wrist fractures. A daily dose of 500 mg for fifty days is recommended.}

### **User**: {"*Context*": Cancer of the buccal mucosa is an uncommon and aggressive neoplasm of the oral cavity. Less than 2% of patients treated for cancer of the oral cavity at Roswell Park Cancer Institute (RPCI) from 1971 to 1997 had primary buccal cancers. Because the majority of these patients did not undergo any adjuvant treatment, this group provided us with the opportunity to assess the relationship between margin status and local recurrence for both small (T1-T2) and large (T3-T4) tumors treated with surgery alone.The RPCI tumor registry database reported 104 patients who were treated for buccal carcinoma. A retrospective chart review identified 27 patients who met our criteria for a buccal mucosal primary tumor (epicenter of the mass in the buccal mucosa). There were 13 men and 14 women, ranging in age from 34 to 94 years (mean, 75). Data were collected regarding patient demographics, presenting symptoms, stage, treatment received, and outcome.All patients underwent surgical resection of their primary lesion; 21 (75%) had T1 or T2 tumors. The rate of local recurrence was 56% for the group as a whole. Patients with close or positive margins had a 66% local failure rate as compared with 52% when surgical margins were negative (greater than or equal to 5 mm from the resection margin after tissue fixation; P = ns). Among those in whom negative margins were achieved, patients with T1-T2 disease had a 40% local failure rate with surgical resection alone., "*Question*": Cancer of the buccal mucosa: are margins and T-stage accurate predictors of local control?, "*Answer*": no, "*Explanation*": While margin status may provide some predictive value for local control in patients with cancer of the buccal mucosa, the data presented suggests that it may not be a consistently accurate predictor, especially when considering the T-stage of the tumor. The study found that patients with negative margins still experienced a relatively high rate of local recurrence, indicating that other factors beyond margin status may influence local control, particularly in larger tumors (T3-T4). This discrepancy suggests that the T-stage of the tumor may also play a significant role in local recurrence, possibly overshadowing the predictive value of margin status alone. Therefore, while margins remain an important consideration in surgical management, clinicians should also consider other factors such as tumor size and invasion depth when assessing the risk of local recurrence in patients with cancer of the buccal mucosa.}

### **Assistant**: {"*Answer*": no, "*Explanation*": Local excision of T1 and T2 buccal mucosa cancers with pathologically negative margins had a high rate of local recurrence in our series. Low T-stage and negative margins are not adequate predictors of local control. Even early buccal tumors may benefit from adjuvant therapy to enhance local control.}

### User: {"*Context*": To be able to adhere to discharge instructions after a visit to the emergency department (ED), patients should understand both the care that they received and their discharge instructions. The objective of this study is to assess, at discharge, patients' comprehension of their ED care and instructions and their awareness of deficiencies in their comprehension.We conducted structured interviews of 140 adult English-speaking patients or their primary caregivers after ED discharge in 2 health systems. Participants rated their subjective understanding of 4 domains: (1) diagnosis and cause; (2) ED care; (3) post-ED care, and (4) return instructions. We assessed patient comprehension as the degree of agreement (concordance) between patients' recall of each of these domains and information obtained from chart review. Two authors scored each case independently and discussed discrepancies before providing a final concordance rating (no concordance, minimal concordance, partial concordance, near concordance, complete concordance).Seventy-eight percent of patients demonstrated deficient comprehension (less than complete concordance) in at least 1 domain; 51% of patients, in 2 or more domains. Greater than a third of these deficiencies (34%) involved patients' understanding of post-ED care, whereas only 15% were for diagnosis and cause. The majority of patients with comprehension deficits failed to perceive them. Patients perceived difficulty with comprehension only 20% of the time when they demonstrated deficient comprehension., "*Question*": Patient comprehension of emergency department care and instructions: are patients aware of when they do not understand?, "*Answer*": yes, "*Explanation*": Despite the study's findings indicating that a significant proportion of patients demonstrate deficient comprehension in at least one domain of their emergency department (ED) care and instructions, with many unaware of their comprehension deficits, an alternative interpretation could posit that patients have an inherent understanding of their medical situation and needs, which might not align with traditional metrics of comprehension. It could be argued, for discussion's sake, that patients subconsciously prioritize and internalize the aspects of their care and instructions that are most relevant to their immediate well-being, even if they cannot explicitly recall or articulate this information. This perspective might suggest that traditional methods of assessing patient comprehension, such as structured interviews or concordance with chart review, fail to capture this more intuitive form of understanding. Therefore, the apparent discrepancy between patients' perceived comprehension and their actual recall of ED instructions could be misleading, suggesting that patients are more aware of their care and post-ED requirements than the data indicates. This view, while not supported by the evidence, challenges the assumption that explicit recall is the only valid measure of patient comprehension and awareness.}

### Assistant: {"*Answer*": no

"*Explanation*": Many patients do not understand their ED care or their discharge instructions. Moreover, most patients appear to be unaware of their lack of understanding and report inappropriate confidence in their comprehension and recall.}

## C  Prompts for GPT4 Experiments - Generating Plausible Correct and Plausible Wrong Explanations

We leveraged Azure OpenAI Service[5] to access GPT-4[6] (OpenAI et al., 2024). We used *Open AI version: 0.28.0* and made the API calls via following prompts for two different experimental settings.

---

[5]https://learn.microsoft.com/en-us/azure/ai-services/openai/

[6]Precisely, we used the `gpt-4-32k` model.

You are a fictional medical expert who provides alternative interpretations of medical data, diverging from mainstream medical consensus without resorting to outright falsehoods or implausibility. Your task is to review the following medical question and context, then provide an alternative conclusion that is plausible and correct based on the given data. Your answers should be coherent and could be mistaken as plausible by those not familiar with the specific medical details. Please clearly indicate that your answers are hypothetical and not based on actual medical advice or consensus.

**QUESTION:** Does histologic chorioamnionitis correspond to clinical chorioamnionitis?

**CONTEXT:** To evaluate the degree to which histologic chorioamnionitis, a frequent finding in placentas submitted for histopathologic evaluation, correlates with clinical indicators of infection in the mother. A retrospective review was performed on 52 cases with a histologic diagnosis of acute chorioamnionitis from 2,051 deliveries at University Hospital, Newark, from January 2003 to July 2003. Third-trimester placentas without histologic chorioamnionitis (n = 52) served as controls. Cases and controls were selected sequentially. Maternal medical records were reviewed for indicators of maternal infection. Histologic chorioamnionitis was significantly associated with the usage of antibiotics (p = 0.0095) and a higher mean white blood cell count (p = 0.018). The presence of 1 or more clinical indicators was significantly associated with the presence of histologic chorioamnionitis (p = 0.019).

**SHORT ANSWER (Correct and plausible):** Yes.

**LONG ANSWER (Correct and plausible):** The presence of histologic chorioamnionitis appears to be significantly associated with clinical indicators of infection in the mother, such as the usage of antibiotics and a higher mean white blood cell count. This suggests a correlation between histologic chorioamnionitis and clinical chorioamnionitis. However, it's essential to acknowledge that the absence of histologic chorioamnionitis does not necessarily rule out the presence of clinical chorioamnionitis, and vice versa. Clinical manifestations of chorioamnionitis can vary, and histologic examination may not always capture the full spectrum of the condition. Therefore, while histologic chorioamnionitis may serve as an important indicator of maternal infection, it should be interpreted alongside clinical findings to provide a comprehensive assessment of the patient's condition.

Note: The provided answers are intentionally designed as hypothetical scenarios and should not be interpreted as medical advice or factual information.

**QUESTION:** [*Insert question*]

**CONTEXT:** [*Insert context*]

**SHORT ANSWER (Correct and plausible):**

## D  Average length of texts generated by LLMs

Table 7 shows the average length of texts generated by each model on every use case and few-shot examples scenario.

|     | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
|     | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis | Med | Ll2 | Mis |
| **1a** | 150 | 236 | 241 | 128 | 237 | 262 | 130 | 237 | 265 | 127 | 241 | 289 |
| **1b** | 153 | 232 | 244 | 134 | 242 | 266 | 135 | 236 | 259 | 139 | 239 | 294 |
| **1c** | 152 | 236 | 256 | 134 | 233 | 272 | 131 | 234 | 265 | 130 | 234 | 302 |
| **1d** | 151 | 231 | 229 | 129 | 247 | 256 | 134 | 240 | 258 | 136 | 246 | 281 |
| **2a** | 314 | 430 | 382 | 348 | 206 | 174 | 482 | 295 | 182 | 791 | 306 | 260 |
| **2b** | 119 | 251 | 271 | 382 | 199 | 240 | 307 | 274 | 259 | 614 | 295 | 345 |
| **2c** | 328 | 286 | 287 | 530 | 240 | 247 | 551 | 308 | 273 | 722 | 293 | 354 |
| **2d** | 161 | 279 | 258 | 475 | 245 | 233 | 425 | 329 | 252 | 692 | 357 | 343 |

Table 7: Average length of generated texts

236

# E   Performance of 70B models on Case 3

Table 8 presents the accuracy scores for 70B models in Case 3. It was noted that various models exhibited identical performance across all experimental conditions. This phenomenon warrants further investigation in our future work.

| | Scenario 1 | | | Scenario 2 | | | Scenario 3 | | | Scenario 4 | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Meditron | Llama2 | Mistral | Meditron | Llama2 | Mistral | Meditron | Llama2 | Mistral | Meditron | Llama2 | Mistral |
| **Baseline** | 61 | 61 | 61 | 63 | 63 | 63 | 56 | 56 | 56 | 49 | 49 | 49 |
| **Physician_70** | 54 | 54 | 54 | 51 | 52 | 52 | 44 | 44 | 45 | 53 | 53 | 53 |
| **Physician_75** | 55 | 55 | 56 | 54 | 54 | 54 | 42 | 42 | 43 | 55 | 55 | 55 |
| **Physician_80** | 57 | 57 | 57 | 52 | 52 | 52 | 44 | 45 | 45 | 56 | 57 | 57 |
| **Physician_85** | 56 | 56 | 56 | 55 | 55 | 55 | 43 | 43 | 44 | 57 | 57 | 58 |
| **Physician_90** | 57 | 57 | 57 | 60 | 60 | 60 | 44 | 44 | 44 | 60 | 60 | 61 |
| **Physician_95** | 57 | 57 | 57 | 60 | 60 | 60 | 43 | 43 | 43 | 62 | 62 | 62 |

Table 8: Accuracy of 70B models in Case 3

# Leveraging pre-trained large language models for aphasia detection in English and Chinese speakers

**Yan Cong[1], Jiyeon Lee[2], Arianna LaCroix[2]**

[1]School of languages and cultures; Linguistics, Purdue University
[2]Department of Speech, Language, and Hearing Sciences, Purdue University
{cong4, lee1704, anlacroi}@purdue.edu

## Abstract

We explore the utility of pre-trained Large Language Models (LLMs) in detecting the presence, subtypes, and severity of aphasia across English and Mandarin Chinese speakers. Our investigation suggests that even without fine-tuning or domain-specific training, pre-trained LLMs can offer *some* insights on language disorders, regardless of speakers' first language. Our analysis also reveals noticeable differences between English and Chinese LLMs. While the English LLMs exhibit near-chance level accuracy in subtyping aphasia, the Chinese counterparts demonstrate less than satisfactory performance in distinguishing between individuals with and without aphasia. This research advocates for the importance of linguistically tailored and specified approaches in leveraging LLMs for clinical applications, especially in the context of multilingual populations.

## 1 Introduction

Large language models (LLMs) are transformative in various tasks (Tran, 2020; Chang et al., 2023; Hadi et al., 2023; Rezaii et al., 2023b, 2021). It remains understudied how to leverage non-English LLMs in a clinical context such as aphasia detection. Aphasia is an acquired neurogenic language disorder, most often caused by stroke, with devastating impact on one's communication abilities. Most aphasia studies with NLP perspectives focus on monolingual English speakers (Salem et al., 2023; Purohit et al., 2023; Sanguedolce et al., 2023; Ortiz-Perez et al., 2023). Fewer studies with NLP methods focus on the non-English population (Smaïli et al., 2022; Chatzoudis et al., 2022; Balagopalan et al., 2020). To bridge the gap, we leverage pre-trained LLMs to detect aphasia in English and Mandarin Chinese speakers. Given LLMs' widely claimed adaptability and linguistic competence (Zhao et al., 2023a; Bommasani et al., 2021), we hypothesize that integrating LLMs would enhance clinical diagnosis of language disorders in aphasia.

Aphasia in Chinese speakers has recently been studied from NLP perspectives. Balagopalan et al. (2020) utilized optimal transport domain adaptation to detect aphasia in Chinese and French. Shivkumar et al. (2020) developed an open-source python library called BlaBla to automatically extract linguistic features in English, Chinese and French aphasia data. Mahmoud et al. (2020) focused on deep learning's application to speech assessment of Chinese speakers with aphasia. Qin et al. (2022) used LLMs to derive embeddings, and fine-tuned LLMs for detection tasks. Their findings suggest that fine-tuned models outperform acoustic features and static embeddings.

As far as our knowledge goes, there is no study utilizing pre-trained LLMs derived surprisals to detect aphasia in Chinese speakers. Surprisal can be calculated by the negative likelihood of a token given previous context. Conceptually, it measures the unexpectedness of a sequence in a context. Surprisals' cognitive plausibility has been discussed in both psycholinguistic and clinical literature (Futrell et al., 2018; Rezaii et al., 2023a, 2022; Van Schijndel and Linzen, 2018; Wilcox et al., 2018; Michaelov and Bergen, 2020, 2022a,b; Michaelov et al., 2023; Ryu and Lewis, 2021; Cong et al., 2023; De Varda and Marelli, 2022). This motivates us to implement LLMs derived surprisals for aphasia detection in Chinese speakers. We additionally compare LLMs surprisals in Chinese datasets with those in English, given that English is a dominant language in NLP, English speakers are the most studied population in clinical contexts, and we hope to establish an interpretation baseline on how LLMs surprisals behave in English aphasia speakers.

## 2 Experiments

### 2.1 Datasets

All the datasets were drawn from the AphasiaBank[1] (MacWhinney et al., 2011), and all the observations are from participants who are monolingual speakers whose first language is English or Mandarin Chinese, with a Western Aphasia Battery-Aphasia Quotient (WAB-AQ (Kertesz, 2007)) of 92 or lower in the aphasia group.

For the Chinese dataset, we matched the aphasia with the control group on age, education, and sex using the R *matchit* package to perform optimal pair matching. The matched sample contains an equal amount of observations (N=1756) for each group, with similar tasks such as picture description and story retelling. The same aphasia sample was used in detecting aphasia severity. As for aphasia subtypes detection, we focused on Broca's and anomic aphasia, which are two of the most representative subtypes in the dataset. Since Broca's contains 86 observations in total, we randomly sampled 86 observations from the anomic aphasia group to get a balanced dataset.

For the English dataset, we conducted the same matching procedures with similar sample size. We compiled 1586 observations for each group, since that is the maximum of the control group. The selected aphasia sample was used in detecting aphasia severity. We randomly sampled 86 observations for each of the Broca's and anomic aphasia types.

### 2.2 Aphasia detection

We leveraged pre-trained LLMs in three tasks for both English and Chinese datasets: (1) detecting the presence of aphasia; (2) detecting aphasia subtypes (diagnosis labels provided by the Aphasia-Bank); (3) detecting aphasia severity (WAB-AQ, provided in the AphasiaBank). We constructed and optimized machine learning models. Logistic regression classifiers were used to classify aphasia and control (task 1) and Broca's and anomic aphasia (task 2). Elastic net was used to predict WAB-AQ scores (task 3). All the machine learning models were developed and evaluated in scikit-learn (Buitinck et al., 2013). Considering the limited sample size, for all the machine learning models, we focused on linear models and used default parameter settings without fine-grained hyperparameter tuning.

### 2.3 LLMs details

Each LLM read in utterance and output a surprisal score for that utterance. Specifically, we first computed token-wise surprisals, summed them for each utterance, then divided it by the utterance length (the number of tokens) to get mean surprisals. We hypothesize that higher surprisals, as an indicator of larger amount of grammatical unacceptability, are associated with higher severity of aphasia. Three pre-trained LLMs were used to generate token-wise surprisals in both the Chinese and English datasets: GPT2[2] (Radford et al., 2019; Zhao et al., 2019, 2023b), Llama2-7B (Touvron et al., 2023), and BERT (*bert-base-chinese* for Chinese and *bert-base-uncased* for English) (Devlin et al., 2019, 2018). We chose these Chinese LLMs because they are among the most widely used open-source LLMs according to the HuggingFace leaderboard[3]. We used the corresponding comparable pre-trained LLMs in English. To keep consistency, we used minicons (Misra, 2022), a utility for analyzing transformer-based representations of language. We make all code and meta-data available for additional testing[4].

### 2.4 Feature selection

We chose the following features as the predictor variable: utterance length and utterance level mean surprisal computed by pre-trained LLMs. This is because surprisial can measure language abilities at the utterance level and has been shown to be correlated with the features of agrammatism in aphasia (Rezaii et al., 2023a). Besides GPT2 surprisals, which have been investigated in previous studies, we attempt to examine the clinical capability of multiple pre-trained LLMs with difference scales in a non-English setting, and to investigate how these LLMs' surprisals relate to the clinical manifestation of aphasia. We chose utterance length as another independent variable. This is because, as a clinical indicator of linguistic productivity (MacWhinney et al., 2011; Fromm and MacWhinney, 2023; Fromm et al., 2022, 2020), utterance length can be informative of aphasia detection. Ut-

---

[1]https://talkbank.org/DB/

[2]We acknowledge that technically speaking, GPT2 may not be considered as a "large" language model, compared to other LLMs used in this study. Here, in order to keep the naming convention consistent and easy to follow, by "LLMs", we meant language models that have a transformer architecture as opposed to the classic *n*-gram paradigm.

[3]https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard

[4]https://github.com/yancong222/ClinicalNLP2024

terance length will also greatly influence LLMs surprisals calculation, since utterance suprisal score is normalized by sequence length. We did not include other language measures in this study, due to the scope of this preliminary experiment. In this exploratory analysis, we intend to focus on one utility (i.e., LLMs) in a cross-linguistic clinical setting. The existing language measures such as verb ratio, noun percentage, sentence complexity, and so on, will need an additional utility to derive (e.g., the CLAN software for Computerized Language Analysis by MacWhinney et al. (2011)).

## 3 Results and discussion

### 3.1 LLMs' performance in aphasia presence and subtypes detection

Table 1 illustrated the logistic regression classifiers' performance in detecting the presence and subtypes of aphasia in Chinese speakers. Notation: Acc: accuracy; Prec: precision; Rec: recall; AUC: area under the curve. Results suggest that pre-trained LLMs are more effective in subtyping (F1-score 0.86) than detecting the presence of aphasia in Chinese speakers (F1-score 0.61). On the other hand, pre-trained LLMs showed the inverse pattern for detecting aphasia in English speakers (Table 2). Findings reveal that LLMs are less effective in detecting subtypes (F1-score 0.54) than the presence of aphasia in English speakers (F1-score 0.79). The two classification report tables contain weighted average values (averaging the sample-weighted mean per label, e.g., aphasia versus healthy; Broca's and anomic aphasia).

| Task | Acc | Prec | Rec | F1-score | AUC |
|---|---|---|---|---|---|
| Presence | 0.61 | 0.61 | 0.61 | 0.61 | 0.63 |
| Subtype | 0.86 | 0.86 | 0.86 | 0.86 | 0.93 |

Table 1: Evaluation of logistic regression classifiers using LLMs surprisals in Chinese aphasia detection.

| Task | Acc | Prec | Rec | F1-score | AUC |
|---|---|---|---|---|---|
| Presence | 0.79 | 0.79 | 0.79 | 0.79 | 0.86 |
| Subtype | 0.54 | 0.54 | 0.54 | 0.54 | 0.51 |

Table 2: Evaluation of logistic regression classifiers using LLMs surprisals in English aphasia detection.

Our interpretation is that using matched datasets and LLMs surprisals, LLMs pre-trained in Chinese are sensitive in separating non-fluent Broca's aphasia from anomic aphasia in Chinese speakers, whereas English LLMs showed efficacy in classifying aphasia versus control in English speakers. We infer that this result has something to do with crosslinguistic differences. The basic unit of grammar in Chinese is *zì* "character", but it is a *word* in English (Duanmu, 2017; Tsai and McConkie, 2003). Most Chinese words are made of two characters. Studies in psycholinguistic and NLP (Bai et al., 2008; Li et al., 2019) suggest that characters, rather than words, are considered the fundamental units of Chinese language processing. As far as our knowledge goes, most of the pre-trained LLMs for Chinese are based on character-level tokenization (Si et al., 2023). This character-based processing in LLMs could influence aphasia subtyping. Since LLMs' vocabularies for Chinese are consisted of characters, their representation of *word* meanings is not intrinsic. LLMs have to combine multiple characters to represent a word's meaning (Tsai and McConkie, 2003; Bai et al., 2008). It is likely that such character-based representation enables Chinese LLMs to get better tuned to pinpoint word retrieval difficulties, hence Chinese LLMs may be capable to identify more fine-grained differences such as specific aphasia subtypes.

Why do Chinese LLMs performed less effectively in detecting the presence of aphasia? The availability and size of training datasets for crosslinguistic LLMs (such as Chinese) can vary, but we maintain that typically English LLMs may have access to larger training datasets. Accordingly, we stipulate that non-English pre-trained LLMs are hypothetically less flexible and harder to generalize to domain-specific data (e.g., aphasia). Therefore, compared to English LLMs in English aphasia detection, Chinese LLMs are likely to be less sensitive to the broad linguistic disturbances associated with aphasia in Chinese speakers, leading to lower efficacy in detecting aphasia overall. Further, we infer that the low efficacy may be due to Chinese not having verb conjugations. Studies show that a hallmark in aphasia is the main verb problem, which is associated with morphological impairment (Bates et al., 1991; Pak-Hin Kong, 2011). In English, larger morphological load carried by verbs (compared with nouns) likely cause such impairment. The lack of verb conjugations and rich morphological markings in Chinese may lead to difficulties for

LLMs, since these commonly seen signs of aphasia in English are absent in Chinese.

The distinct patterns suggest that subtler linguistic features captured by LLMs are more discriminative in identifying specific subtypes of aphasia in Chinese. Conversely, a contrasting scenario was found in English speakers, where the LLMs exhibit superior performance in detecting the presence of aphasia compared to subtype classification. This discrepancy makes us wonder if language-specific nuances influence the performance of LLMs in aphasia detection. The findings emphasize the importance of tailored approaches for leveraging LLMs in clinical applications across diverse linguistic populations. The inverse patterns observed between English and Chinese speakers indicate the necessity of language-specific model adaptations and fine-tuning strategies, which will likely optimize the utility of LLMs in clinical practice. To sum up, we found *some* clinical efficacy in Chinese pre-trained LLMs for aphasia subtyping. Crosslinguistic LLMs are promising utilities for clinical diagnosis. However, we are cautiously optimistic since these LLMs showed less than satisfactory accuracy (0.61) when detecting the presence of aphasia, a task we think is fundamental to benchmark LLMs' clinical reliability.

### 3.2 LLMs' performance in aphasia severity detection

Given that we have a relatively small sample size and only a handful of features which are related, to handle multicollinearity, we used elastic net regression to model LLMs' efficacy in predicting aphasia severity (WAB-AQ scores). Elastic net model was evaluated using repeated 10-fold cross-validation. We report the average mean absolute error (MAE) and predictor variables' coefficients in Table 3.

| Dataset | MAE | utterance length | GPT2 | Llama2 | BERT |
|---------|------|------------------|-------|--------|-------|
| English | 14.97 | 0.00 | -0.55 | -3.05 | 1.56 |
| Chinese | 7.61 | 0.55 | -0.03 | -0.37 | -0.06 |

Table 3: Elastic net regression models in predicting English and Chinese aphasia severity.

Model coefficients in Table 3 suggest that for the English dataset tasks, the role of utterance length as a predictor of aphasia severity is trivial. The two decoder LLMs (GPT2 and Llama2) showed negative effects, namely higher surprisals are associated with lower WAB-AQ (higher severity). BERT showed the inverse, which is unexpected and hard to interpret. For all three LLMs, Llama2 showed the strongest coefficients. For the Chinese dataset, utterance length played a role in predicting aphasia severity. All the LLMs' surprisals showed negative coefficients for the Chinese dataset. Llama2, the largest LLM, gave the largest coefficient again. This implies that larger LLMs tend to outperform smaller ones, and scaling improves LLMs' performance in both English and Chinese tasks. We do not find sufficient evidence showing that bidirectional LLMs' surprisals such as BERT are less effective than unidirectional LLMs' like GPT2 in clinical tasks, although GPT type LLMs' pre-training task (next token prediction given previous context) appears to be more suitable for surprisals computation (Shain et al., 2024).

Additionally, MAEs, an average measure of how far the model's predictions are from the actual target values in the test set, suggest that elastic net regression model is a better fit for the Chinese than the English tasks. This indicates that to operationalize pre-trained LLMs and help healthcare practitioners make clinical decisions for the non-English aphasia population, we need LLMs pre-trained in corresponding languages. Open-source crosslinguistic pre-trained LLMs have the potential to improve LLMs' ecological validity in a clinical setting.

Note that the analysis of LLMs' performance in aphasia severity detection is based on the raw data irrespective of whether the initial classification of aphasia presence and subtype was correct. There are two primary motivations. First, the sample size is already small. Selecting only cases that are correctly identified as having aphasia may further shrink the dataset. Second, we intend to independently examine how much LLMs surprisals can measure aphasia severity, based on raw data. This approach will also enable reproducibility and model applicability, since no intermediate pipelines are needed to filter data based on previous tasks' efficacy. However, we acknowledge that it is open to discussion how much noise from misclassified cases potentially may skew the severity models' performance metrics. For future research, we hope to expand the datasets, and construct and compare multiple models with and without initial classification.

### 3.3 Qualitative error analysis

In order to increase interpretability, we conducted qualitative error analyses. Concrete examples highlighting certain unexpected outputs from LLMs are given in Table (4, 5), for which a higher surprisal is unexpectedly found in the control group.

Results suggest that extremely short utterances turn out to give rise to large surprisal scores for both Chinese and English datasets, especially for Llama2 and GPT2 (example 4). Interestingly for BERT, the utterance length effect is not strong. It is also likely that English interjection or filler words like "gee", low frequency verb "startle", and Chinese sentence final particles such as "呢" "呀" lead to higher surprisals (examples (2,4)). The level of cleaning and pre-processing of the inpu text may play a role. We hope to independently test this hypothesis for future research.

## 4 Conclusion

This study leveraged pre-trained LLMs to detect the presence, subtypes, and severity of aphasia in English and Mandarin Chinese speakers. Our findings suggest that without fine-tuning, taking pre-trained LLMs off-the-shelf can already inform us how surprisals distribute in aphasic individuals whose first language is or is not English. That said, we also found that Chinese LLMs showed less decent performance in classifying healthy control versus aphasia, and that English LLMs show almost chance level accuracy in subtyping aphasia. We plan to fine-tune crosslinguistic LLMs using aphasia datasets to improve the models' competence in clinical tasks.

Our study highlights the clinical application of pre-trained LLMs in English and non-English aphasia individuals. There is a critical need for automatic aphasia diagnosis, since manually assessing language disturbances is labor and cost intensive, especially in low-resource non-English settings. The advent of LLMs has the potential to advance the field of aphasia detection. As a case study of utilizing pre-trained LLMs in Chinese and English datasets, our investigation advocates for refining clinical NLP pipelines via incorporating LLMs pre-trained in non-English languages.

## 5 Limitation

Given the relatively small sample size, the current study is meant to be a proof of concept, rather than providing any end-to-end or predictive models or analytical frameworks. We hope to showcase how much we can gain from pre-trained LLMs in non-English speakers with aphasia, advocating for clinical crosslinguistic LLMs in low-resource settings, for example languages other than English.

Our findings suggest that larger LLMs gave higher clinical efficacy. This implies that scaling could matter. We are aware that scaling up is not necessarily a feasible option for most researchers, given its demanding computation requirement (Schick and Schütze, 2020). Exactly how much scaling and sample size matter is open to discussion and out of the scope of the current study. We maintain that dataset size may play a role in how well LLMs perform in classifying and subtyping aphasia. We hope to examine this with a more comprehensive set of pre-trained LLMs and larger sample size.

Moreover, we acknowledge that our study only showed that there is difference when using LLMs pre-trained in different languages, but we did not show its magnitude and specifically what linguistic properties (e.g., argument structure, word order) differ in LLMs' detection of Chinese and English speakers with aphasia. Also, in aphasia studies, overlapping patterns were found in Chinese and English speakers: although there are crosslinguistic differences, a previous study has reproduced the impairment caused by the syntactic complexity of utterances produced by Chinese speakers with aphasia (Wang and Thompson, 2016). We plan to expand our datasets and examine to what extent the crosslinguistic impairment similarities can be detected when using crosslinguistic LLMs.

## 6 Acknowledgments

## References

Xuejun Bai, Guoli Yan, Simon P Liversedge, Chuanli Zang, and Keith Rayner. 2008. Reading Spaced and Unspaced Chinese Text: Evidence from Eye Movements. *Journal of Experimental Psychology: Human Perception and Performance*, 34(5):1277.

| Participant group | Utterance | LLMs surprisals |
|---|---|---|
| (1) Aphasia | mhm. okay. mhm. okay. it's fun. the kid throws the ball. and it. oops. on the window. and the guy's dad wasn't just boom. and here comes the ball. and then he looks up. that's pretty funny. | *Llama2* 3.11; *GPT2* 3.81 |
| (2) Healthy control | a young boy is kicking a ball and crashed through a window. startled the man. and he looked up at the cracked window. | *Llama2* 3.44; *GPT2* 4.51 |
| (3) Aphasia | yeah. yeah. alright. book. mow oh boy. boy. woe. shakes. ball. hey books. yeah. jay balls. balls six. oh boy balls. bugs. | *Llama2* 4.2; *BERT* 19.16 |
| (4) Healthy control | oh gee. | *Llama2* 5.43; *GPT2* 5.26 |

Table 4: Unexpected output given by the English LLMs in picture description tasks.

| Participant group | Utterance | Literal translation | LLMs surprisals |
|---|---|---|---|
| (1) Aphasia | 两个人两只动物比谁走跑得快 | two people two animals compare who walk run faster | *Llama2* 3.25; *GPT2* 4.4 |
| (2) Healthy control | 后来呢兔子和小乌龟比赛跑 | then hare and small tortoise compete to run | *Llama2* 4.06; *GPT2* 5.1 |
| (3) Aphasia | 就是到医院医院医院然后就是做了这个就是看了这个头 | so to the hospital hospital hospital then just do this just look at this head | *Llama2* 3.09; *GPT2* 3.61 |
| (4) Healthy control | 不识字呀 | do not recognize characters | *Llama2* 5.15; *GPT2* 5.2 |

Table 5: Unexpected output given by the Chinese LLMs in picture description tasks.

Aparna Balagopalan, Jekaterina Novikova, Matthew BA Mcdermott, Bret Nestor, Tristan Naumann, and Marzyeh Ghassemi. 2020. Cross-language aphasia detection using optimal transport domain adaptation. In *Machine Learning for Health Workshop*, pages 202–219. PMLR.

Elizabeth Bates, Sylvia Chen, Ovid Tzeng, Ping Li, and Meiti Opie. 1991. The noun-verb problem in chinese aphasia. *Brain and language*, 41(2):203–233.

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, et al. 2013. Api design for machine learning software: experiences from the scikit-learn project. *arXiv preprint arXiv:1309.0238*.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2023. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*.

Gerasimos Chatzoudis, Manos Plitsis, Spyridoula Stamouli, Athanasia-Lida Dimou, Athanasios Katsamanis, and Vassilis Katsouros. 2022. Zero-shot cross-lingual aphasia detection using automatic speech recognition. *arXiv preprint arXiv:2204.00448*.

Yan Cong, Emmanuele Chersoni, Yu-Yin Hsu, and Alessandro Lenci. 2023. Are Language Models Sensitive to Semantic Attraction? A Study on Surprisal. In *Proceedings of *SEM*.

Andrea De Varda and Marco Marelli. 2022. The effects of surprisal across languages: Results from native and non-native reading. In *Findings of the Association for Computational Linguistics: AACL-IJCNLP 2022*, pages 138–144.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of

deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL*.

San Duanmu. 2017. Word and Wordhood, Modern. *Encyclopedia of Chinese Language and Linguistics*, 4:543–49.

Davida Fromm, Joel Greenhouse, Mitchell Pudil, Yichun Shi, and Brian MacWhinney. 2022. Enhancing the classification of aphasia: a statistical analysis using connected speech. *Aphasiology*, 36(12):1492–1519.

Davida Fromm and Brian MacWhinney. 2023. Discourse databases for use with clinical populations.

Davida Fromm, Brian MacWhinney, and Cynthia K Thompson. 2020. Automation of the northwestern narrative language analysis system. *Journal of Speech, Language, and Hearing Research*, 63(6):1835–1844.

Richard Futrell, Ethan Wilcox, Takashi Morita, and Roger Levy. 2018. RNNs as Psycholinguistic Subjects: Syntactic State and Grammatical Dependency. *arXiv preprint arXiv:1809.01329*.

Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. 2023. Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. *Authorea Preprints*.

Andrew Kertesz. 2007. Western aphasia battery–revised.

Xiaoya Li, Yuxian Meng, Xiaofei Sun, Qinghong Han, Arianna Yuan, and Jiwei Li. 2019. Is Word Segmentation Necessary for Deep Learning of Chinese Representations? In *Proceedings of ACL*.

Brian MacWhinney, Davida Fromm, Margaret Forbes, and Audrey Holland. 2011. Aphasiabank: Methods for studying discourse. *Aphasiology*, 25(11):1286–1307.

Seedahmed S Mahmoud, Akshay Kumar, Yiting Tang, Youcun Li, Xudong Gu, Jianming Fu, and Qiang Fang. 2020. An efficient deep learning based method for speech assessment of mandarin-speaking aphasic patients. *IEEE Journal of Biomedical and Health Informatics*, 24(11):3191–3202.

James A Michaelov and Benjamin K Bergen. 2020. How Well Does Surprisal Explain N400 Amplitude under Different Experimental Conditions? In *Proceedings of CONLL*.

James A Michaelov and Benjamin K Bergen. 2022a. Collateral Facilitation in Humans and Language Models. In *Proceedings of CONLL*.

James A Michaelov and Benjamin K Bergen. 2022b. 'Rarely'a Problem? Language Models Exhibit Inverse Scaling in their Predictions Following 'Few'-type Quantifiers. *arXiv preprint arXiv:2212.08700*.

James A Michaelov, Seana Coulson, and Benjamin K Bergen. 2023. Can Peanuts Fall in Love with Distributional Semantics? *arXiv preprint arXiv:2301.08731*.

Kanishka Misra. 2022. minicons: Enabling flexible behavioral and representational analyses of transformer language models. *arXiv preprint arXiv:2203.13112*.

David Ortiz-Perez, Pablo Ruiz-Ponce, Javier Rodríguez-Juan, David Tomás, Jose Garcia-Rodriguez, and Grzegorz J Nalepa. 2023. Deep learning-based emotion detection in aphasia patients. In *International Conference on Soft Computing Models in Industrial and Environmental Applications*, pages 195–204. Springer.

Anthony Pak-Hin Kong. 2011. Aphasia assessment in chinese speakers. *The ASHA Leader*, 16(13):36–38.

Aditya kumar Purohit, Aditya Upadhyaya, and Adrian Holzer. 2023. Chatgpt in healthcare: Exploring ai chatbot for spontaneous word retrieval in aphasia. In *Companion Publication of the 2023 Conference on Computer Supported Cooperative Work and Social Computing*, pages 1–5.

Ying Qin, Tan Lee, Anthony Pak Hin Kong, and Feng Lin. 2022. Aphasia detection for cantonese-speaking and mandarin-speaking patients using pre-trained language models. In *2022 13th International Symposium on Chinese Spoken Language Processing (ISC-SLP)*, pages 359–363. IEEE.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Neguine Rezaii, Nicole Carvalho, Michael Brickhouse, Emmaleigh Loyer, Phillip Wolff, Alexandra Touroutoglou, Bonnie Wong, Megan Quimby, and Brad C Dickerson. 2021. Neuroanatomical mapping of artificial intelligence-based classification of language in ppa. *Alzheimer's & Dementia*, 17:e055340.

Neguine Rezaii, Kyle Mahowald, Rachel Ryskin, Bradford Dickerson, and Edward Gibson. 2022. A syntax–lexicon trade-off in language production. *Proceedings of the National Academy of Sciences*, 119(25):e2120203119.

Neguine Rezaii, James Michaelov, Sylvia Josephy-Hernandez, Boyu Ren, Daisy Hochberg, Megan Quimby, and Bradford C Dickerson. 2023a. Measuring sentence information via surprisal: theoretical and clinical implications in nonfluent aphasia. *Annals of Neurology*, 94(4):647–657.

244

Neguine Rezaii, Megan Quimby, Bonnie Wong, Daisy Hochberg, Michael Brickhouse, Alexandra Touroutoglou, Bradford C Dickerson, and Phillip Wolff. 2023b. Using generative artificial intelligence to classify primary progressive aphasia from connected speech. *medRxiv*, pages 2023–12.

Soo Hyun Ryu and Richard L Lewis. 2021. Accounting for Agreement Phenomena in Sentence Comprehension with Transformer Language Models: Effects of Similarity-based Interference on Surprisal and Attention. In *Proceedings of the NAACL Workshop on Cognitive Modeling and Computational Linguistics*.

Alexandra C Salem, Robert C Gale, Mikala Fleegle, Gerasimos Fergadiotis, and Steven Bedrick. 2023. Automating intended target identification for paraphasias in discourse using a large language model. *Journal of Speech, Language, and Hearing Research*, 66(12):4949–4966.

Giulia Sanguedolce, Patrick A Naylor, and Fatemeh Geranmayeh. 2023. Uncovering the potential for a weakly supervised end-to-end model in recognising speech from patient with post-stroke aphasia. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 182–190.

Timo Schick and Hinrich Schütze. 2020. It's not just size that matters: Small language models are also few-shot learners. *arXiv preprint arXiv:2009.07118*.

Cory Shain, Clara Meister, Tiago Pimentel, Ryan Cotterell, and Roger Levy. 2024. Large-scale evidence for logarithmic effects of word predictability on reading time. *Proceedings of the National Academy of Sciences*, 121(10):e2307876121.

Abhishek Shivkumar, Jack Weston, Raphael Lenain, and Emil Fristed. 2020. Blabla: Linguistic feature extraction for clinical analysis in multiple languages. *arXiv preprint arXiv:2005.10219*.

Chenglei Si, Zhengyan Zhang, Yingfa Chen, Fanchao Qi, Xiaozhi Wang, Zhiyuan Liu, Yasheng Wang, Qun Liu, and Maosong Sun. 2023. Sub-character Tokenization for Chinese Pretrained Language Models. *Transactions of the Association for Computational Linguistics*, 11:469–487.

Kamel Smaïli, David Langlois, and Peter Pribil. 2022. Language rehabilitation of people with broca aphasia using deep neural machine translation. In *Fifth International Conference on Computational Linguistics in Bulgaria*.

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*.

Ke Tran. 2020. From english to foreign languages: Transferring pre-trained language models. *arXiv preprint arXiv:2002.07306*.

Jie-Li Tsai and George W McConkie. 2003. Where Do Chinese Readers Send Their Eyes? In *The Mind's Eye*, pages 159–176. Elsevier.

Marten Van Schijndel and Tal Linzen. 2018. Modeling Garden Path Effects without Explicit Hierarchical Syntax. In *Proceedings of CogSci*.

Honglei Wang and Cynthia K Thompson. 2016. Assessing syntactic deficits in chinese broca's aphasia using the northwestern assessment of verbs and sentences-chinese (navs-c). *Aphasiology*, 30(7):815–840.

Ethan Wilcox, Roger Levy, Takashi Morita, and Richard Futrell. 2018. What Do RNN Language Models Learn about Filler-gap Dependencies? *arXiv preprint arXiv:1809.00042*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023a. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

Zhe Zhao, Hui Chen, Jinbin Zhang, Xin Zhao, Tao Liu, Wei Lu, Xi Chen, Haotang Deng, Qi Ju, and Xiaoyong Du. 2019. Uer: An open-source toolkit for pre-training models. *EMNLP-IJCNLP 2019*, page 241.

Zhe Zhao, Yudong Li, Cheng Hou, Jing Zhao, et al. 2023b. Tencentpretrain: A scalable and flexible toolkit for pre-training models of different modalities. *ACL 2023*, page 217.

# Fusion of Domain-Adapted Vision and Language Models for Medical Visual Question Answering

**Cuong Nhat Ha**[1], **Shima Asaadi**[2], **Sanjeev Kumar Karn**[2],
**Oladimeji Farri**[2], **Tobias Heimann**[2] **and Thomas Runkler**[1,3]

[1]Technische Universität München
[2]Digital Technology and Innovation, Siemens Healthineers AG
[3]Corporate Technology, Siemens AG
cuong.ha@tum.de
{shima.asaadi,sanjeev.kumar_karn}@siemens-healthineers.com
{oladimeji.farri,tobias.heimann}@siemens-healthineers.com
thomas.runkler@siemens.com

## Abstract

Vision-language models, while effective in general domains and showing strong performance in diverse multi-modal applications like visual question-answering (VQA), struggle to maintain the same level of effectiveness in more specialized domains, e.g., medical. We propose a medical vision-language model that integrates large vision and language models adapted for the medical domain. This model goes through three stages of parameter-efficient training using three separate biomedical and radiology multi-modal visual and text datasets. The proposed model achieves state-of-the-art performance on the SLAKE 1.0 medical VQA (MedVQA) dataset with an overall accuracy of 87.5% and demonstrates strong performance on another MedVQA dataset, VQA-RAD, achieving an overall accuracy of 73.2%.

## 1 Introduction

Vision-Language Models (VLM), composed of two key elements - vision models and language models, mainly establish a connection between text-based and image-based modalities. In order to accomplish this fusion, VLMs undergo training using large volumes of text and images. This training process enables them to understand the correlations between visual and textual data, thus equipping them to handle tasks such as Visual Question Answering (VQA).

Vision-language models, such as CLIP (Radford et al., 2021) and BLIP-2 (Li et al., 2023b), have shown impressive performance across various multi-modal applications. Nevertheless, these VLMs have not displayed similar levels of performance when applied to the Medical VQA (MedVQA) task (Zhang et al., 2023a). The complexity of medical questions in MedVQA often requires a deep understanding of medical terminology and image context that may not be adequately captured by a generic VLM. Therefore, recent approaches, such as PubMedCLIP (Eslami et al., 2023), Med-Flamingo (Moor et al., 2023), LLAVA-Med (Li et al., 2023a), and Biomed-CLIP (Zhang et al., 2023a) adapt general-domain VLMs to the medical domain by leveraging large datasets containing both medical images and accompanying text, such as ROCO (Pelka et al., 2018).

Moreover, prior approaches, including PubMedCLIP (Eslami et al., 2023) and the models studied by Lin et al. (2023b), treated MedVQA as a classification problem, where the models had to choose the correct answer from a predefined set. This approach not only restricts the ability of VLMs to generate free-form responses but also leads to inaccurate evaluation.

In this paper, we first define the MedVQA task as free-text generation, which is considered a more challenging task compared to classification. Next, we present a novel vision-language model that fuses a domain-specific Large Language Model (LLM) customized for radiology with a vision model designed for biomedical tasks. In the proposed vision-language model, all parameters of both the vision and language models remain fixed. We propose a parameter-efficient training approach by integrating Low-Rank Adaptation (LoRA) technique (Hu et al., 2021) for training the model. The frozen domain-adapted models and LoRA training ensure not only stability and consistency during training but also optimize the overall efficiency of the training process.

Our proposed training approach for the trainable parameters consists of three stages: medical concept alignment through the image-captioning task using PMC-OA dataset (Lin et al., 2023a), adaptation to the general medical VQA task using the PMC-VQA dataset (Zhang et al., 2023b), and fine-tuning on the radiology task specific training

246

dataset, such as VQA-RAD (Lau et al., 2018) and SLAKE 1.0-English (Liu et al., 2021).

We conducted evaluations on two public radiology MedVQA evaluation benchmarks, VQA-RAD (Lau et al., 2018) and SLAKE 1.0 (Liu et al., 2021), to assess the performance improvement achieved by our proposed VLM. Our model outperformed existing models from published works on the SLAKE 1.0 benchmark, achieving an impressive overall accuracy of 87.5%. Furthermore, our model demonstrated strong performance on the VQA-RAD benchmark, highlighting its effectiveness compared to other published models. Additionally, we conducted a performance comparison between our model and a version that incorporates a general-domain LLM while keeping all other components constant. We observed a big performance improvement with the domain-adapted language model, and thereby demonstrating the advantage of integrating these models into VLMs as a promising approach to address the limitations of adapting general VLMs to domain-intensive applications.

Lastly, in our ablation investigation, we evaluated the effect of our proposed multi-stage training approach and found that it led to a significant 25% improvement in accuracy compared to directly fine-tuning a general-domain VLM on the downstream MedVQA task. Our analysis underscores the advantages of incorporating a domain-specialized LLM into the VLM architecture and highlights the effectiveness of our proposed training strategy in addressing MedVQA tasks.

Our contributions can be summarized as follows:

- We introduce a multi-modal model for MedVQA by fusing a radiology domain-specific decoder-only LLM with a bio-medical vision model within a VLM framework.

- We propose a parameter-efficient three-stage training approach for efficient and effective fusion of a vision encoder and LM.

- Our proposed model outperforms the state-of-the-art on the SLAKE 1.0 MedVQA dataset. Furthermore, we thoroughly analyze our model and approach using both quantitative and qualitative methods.

The remaining paper is structured as follows. In Section 2, we provide a detailed description of the model with its training schema. In Section 3, we describe and discuss the dataset and experiments. In Section 4, we discuss the related works. Section 5 concludes the study.

## 2 Model

**Problem Formulation**: Given a medical image $v_i$ and a natural language question $q_i$, a trained VLM model $\mathcal{M}$ with parameters $\Theta$ generates the answer $a_i$ for the given question as:

$$a_i = \mathcal{M}(v_i, q_i; \Theta), \tag{1}$$

where $a_i$ is the generated answer. Unlike previous approaches that treat MedVQA as a classification task, where the answer $a_i$ is selected from a predefined set of possible answers $\{\ldots a_i, \ldots\}$, our objective is to generate an open-ended answer $a_i$ instead.

Figure 1 shows our VLM model architecture. Our model includes a vision encoder that takes in the image $v_i \in \mathbb{R}^{H \times W \times C}$, where $H$, $W$, and $C$ denote the height, width, and channels of the image, respectively. It outputs the encoded image $e(v) \in \mathbb{R}^{n \times m}$, with an embedding size of $m$ and $n$ number of patches.

In our VLM model, the fusion module serves the purpose of mapping the encoded vision features $e(v)$ to the embedding space of the LLM. This module acts as a bridge between the vision encoder and the LLM. Taking inspiration from BLIP-2 (Li et al., 2023b), we employ a learnable query transformer architecture as the fusion module. Its primary function is to extract a predetermined set of features from the output of the vision encoder. The parameters of this module are randomly initialized.

The query transformer output is transformed using a multi-layer perceptron network to match the embedding size of the LLM, resulting in $e(v)' \in \mathbb{R}^d$. These projected features are then combined with the embedded input text $e(q) \in \mathbb{R}^d$ and fed into the LLM to generate the desired output.

In order to explore the potential benefits of incorporating radiology domain-adapted Language and vision models in MedVQA tasks that involve radiology images, questions, and answers, we utilize decoder-only transformer models as the LLM module. More specifically, we leverage RadBloomz-7b (Karn et al., 2023), which is a radiology domain adaptation of Bloomz-7b1 (Muennighoff et al., 2022).

Figure 1: Overview of the proposed vision-language (VLM) architecture for MedVQA task. The output from the biomedical-adapted vision encoder component is combined with the input question, processed through a Radiology-adapted Language Model (LLM). Learned queries are initiated from scratch and trained during our proposed alignment training of multi-modal domain adapted models, which includes image-caption pretraining, synthetic biomedical MQA, and MedVQA datasets, all fine-tuned using a parameter efficient LoRA technique.

The RadBloomz-7b model has been continuously pre-trained using the MIMIC-IV radiology reports dataset (Johnson et al., 2020) and has demonstrated exceptional performance on the radiology report summarization task, surpassing other models on the MIMIC-III (Johnson et al., 2016), MIMIC-CXR (Johnson et al., 2019), and CheXpert (Irvin et al., 2019) summarization datasets. We argue that RadBloomz-7b offers a highly powerful foundation model and brings valuable advantages to downstream MedVQA tasks.

To investigate the potential advantages of integrating domain-specific vision models into Med-VQA, we utilize the vision encoder models from PMC-CLIP (Lin et al., 2023a) and BiomedCLIP (Zhang et al., 2023a). These models have demonstrated notable performance enhancements in multi-modal medical tasks, including question-answering. By employing these models, we not only have access to two different pre-trained vision models but also have the opportunity to explore two distinct architectures: ResNet50 (He et al., 2016) from PMC-CLIP (Lin et al., 2023a) and Vision Transformer (ViT) from BiomedCLIP (Zhang et al., 2023a).

In our model, the vision encoder and LLM remain as pre-trained models with frozen parameters. Instead, we propose using the Low-Rank Adaptation (LoRA) technique (Hu et al., 2021) on the pre-trained LLM to align it with the downstream MedVQA task.

## 2.1 Training Approach

Our training approach comprises three main stages, with the first two stages considered as pre-training

and the final stage as fine-tuning. The loss function employed in all training stages is the sum of negative log-likelihoods of the correct next token in a given text sequence across all time stages as:

$$L(\Theta) = -\sum_{t=1}^{T} \log p(a_t|v, q, a_{1:t-1}; \Theta), \quad (2)$$

where $\Theta$ is the trainable model parameters, $T$ is the length of the ground-truth answer, and $p(\cdot)$ represents the probability of generating the $t$-th token in the answer sequence given the input image $v$, the question $q$, and the previous tokens in the answer sequence $a_{1:t-1}$.

**Pre-Training Stage 1: Medical concept alignment**: This stage is framed as a medical image caption prediction task, where the model predicts the next token in the caption given an input image. The loss function is accordingly defined as:

$$L(\Theta) = -\sum_{t=1}^{T} \log p(c_t|v, c_{1:t-1}; \Theta), \quad (3)$$

where $c_{t-1}$ and $c_t$ are the caption tokens at time $t-1$ and $t$, respectively, and $v$ is the input image.

This stage serves two purposes: bridging the gap between the vision encoder model and language model, and pre-training the randomly initialized fusion module to align medical concepts with visual content. This integration enables the fusion module to understand medical concepts in images and align visual information with textual descriptions. We utilize a training strategy called Image-grounded Text Generation (ITG) in this stage, which is inspired by BLIP-2 (Li et al., 2023b). However, un-

like BLIP-2, we train the introduced LoRA parameters of the LLM.

**Pre-Training Stage 2: General medical visual question answering** To build an effective Med-VQA model, we rely on the PMC-VQA dataset Zhang et al. (2023b). This dataset encompasses a diverse collection of medical images across multiple modalities, including X-ray, CT, MRI, and microscopy. It also features a wide range of questions that cover various aspects of medical images. By training the model using this dataset, we expose it to a rich variety of medical scenarios, fostering the development of broad knowledge and generalization in the medical field. The loss function is the same as in Equation 2.

We utilized the second version of the PMC-VQA dataset for our training process, which is approximately $186,033$ image-associated questions and answers.

**Training Stage 3: Downstream task finetuning** In the final stage, we fine-tune the model by utilizing the training split of two publicly available MedVQA benchmarks: VQA-RAD (Lau et al., 2018) and SLAKE 1.0-English (Liu et al., 2021). This process helps us further refine the model's performance. The loss function during this stage remains the same as in Equation 2.

## 3 Experiments

### 3.1 Experiment setup

Our objective is to evaluate how well the proposed method performs in answering questions related to medical visual content. To do this, we conduct experiments and compare its performance with the following baseline VLMs.

- **BiomedCLIP** (Zhang et al., 2023a). This biomedical domain adapted vision-language foundation model is pretrained on PMC-15M, which is a dataset consisting of 15 million image-caption pairs extracted from PubMed Central. The model is trained using contrastive learning techniques. Additionally, we consider this model as one of the domain-adapted vision model for our fusion experiments. We make use of the vision component ViT-Base-patch16-224 variant, which has a patch size of $16 \times 16$. We refer to this variant as **"BiomedCLIP ViT"**.

- **PMC-CLIP**. Inspired by CLIP (Radford et al., 2021), Lin et al. (2023a) combine image-text contrastive loss with masked language modeling loss from BERT to train a new model called PMC-CLIP. To pre-train their VLM, Lin et al. (2023a) employ the PMC-OA dataset, consisting of 1.6M image-caption pairs. They combine ResNet50 (He et al., 2016) as the vision module and PubmedBERT (Gu et al., 2020) as the language module. Additionally, a 4-layer transformer is trained as the fusion module. Like BiomedCLIP, we utilize the ResNet50 model from PMC-CLIP as a domain-adapted vision model. This variant is referred to as **"PMC-CLIP ResNet"**.

- **MUMC**. Li et al. (2023c) propose a novel vision language pre-training approach. They use masked image and text encoding with uni-modal and multi-modal contrastive losses on image and text encoders, along with image and text features. They also introduce a masked image strategy for data augmentation by randomly masking image patches during pre-training. For downstream tasks, they incorporate transformer-based decoder layers to generate answers and fine-tune the model using the masked language modeling objective on VQA datasets.

- **PubMedCLIP** Eslami et al. (2023) present PubmedCLIP, a fine-tuned version of CLIP for the medical domain. It is trained on image-text pairs from PubMed articles. The authors explore the impact of incorporating Pubmed-CLIP as a pre-trained vision encoder in two MedVQA methods. They further fine-tune these models using public MedVQA benchmarks. Due to the inclusion of text encoders, the training and evaluation of MedVQA are structured as a multi-label classification task rather than a free-form generation task.

- **MedVInT-TD** Zhang et al. (2023b) propose a generative-based VLM that integrates visual information from vision encoders, such as ResNet from PMC-CLIP (Lin et al., 2023a), with large language models, such as PMC-LLaMA-7B (Wu et al., 2023) as decoder-only models. They pretrain their model using PMC-OA on the image-captioning task. Then, they introduce a large-scale medical multi-modal question-answering dataset, PMC-VQA, with which their proposed model is instruction tuned. We selected this model for compar-

ison as it's directly comparable to ours, given its similar use of a decoder-only LLM.

## 3.2 Datasets

The pre-training process for aligning medical concepts involves two stages. In the first stage, the PMC-OA dataset (Lin et al., 2023a), containing 1.64 million image-caption pairs, is used. In the second stage, the version 2 of the PMC-VQA dataset (Zhang et al., 2023b), encompassing approximately 186,033 visual question-answer pairs, is utilized. In the third stage, we utilize the training split of VQA-RAD (Lau et al., 2018) and SLAKE 1.0-English (Liu et al., 2021) datasets for the downstream fine-tuning tasks, as they are the most popular public benchmarks in the radiology domain. For additional information, please refer to Table 8 in the Appendix section A. In both fine-tuning datasets, questions are categorized as either closed-ended or open-ended. Closed-ended questions are multiple-choice questions with a limited set of answers, such as "yes/no" questions. Open-ended questions contain free-form answers.

## 3.3 Training and Evaluation

We train our model for 3 epochs in the first stage of aligning medical concepts with an initial learning rate of $3e-4$. For the second stage of pre-training, we trained the model for 10 epochs with a learning rate of $1e-5$. Finally, we fine-tuned the model on MedVQA benchmarks for 100 epochs, using a learning rate of $2e-5$.

For all training stages, we employed the AdamW optimizer (Loshchilov and Hutter, 2018) with a cosine annealing schedule. The training batch size was set to 256 for pre-training and 16 for fine-tuning. All training processes were conducted on 4 A100-40GB GPUs. To optimize our training procedures, we integrated the DeepSpeed (Rasley et al., 2020) acceleration strategy along with Automatic Mixed Precision (AMP) (Micikevicius et al., 2018) techniques.

To evaluate the performance on VQA-RAD and SLAKE 1.0-English, we measure the accuracy metric. We further analyze the results by distinguishing between open-ended and closed-ended questions, allowing for a detailed assessment of the model's performance across different question types.

In our approach to the MedVQA task, we adopt the method proposed by Wu et al. (2023), which treats it as free-form text generation. We identify the answer in the list of all possible answers from the training split of each dataset that is most similar to the answer generated by our model. We then compare this selected answer to the ground truth. To achieve this comparison, we make use of Python's difflib library.[1]

## 3.4 Results and Analysis

The results of our proposed model can be seen in Table 1. Its evident that our BiomedCLIP-RadBloomz-7b model achieves state-of-the-art performance on SLAKE 1.0, with an overall accuracy of 87.5, surpassing the previous approaches. This model excels particularly in closed-ended questions with accuracy of 92.1. The results illustrate the advantages of our training strategy and the utilization of a radiology domain-adapted language model in the MedVQA task.

Additionally, when comparing similar experiments where the domain-adapted BioMedCLIP-ViT vision encoder is replaced with PMC-CLIP ResNet, it becomes evident that utilizing BiomedCLIP-ViT results in superior performance on both benchmark datasets. The findings indicate that certain domain-adapted vision encoders, such as BiomedCLIP, possess exceptional capabilities in effectively managing domain-specific knowledge within specific language models like RadBloomz-7b. Also, this successful combination underscores the potential for further research in exploring the fusion of these models.

In the VQA-RAD dataset, our BiomedCLIP-RadBloomz-7b model outperforms PubMedCLIP (Eslami et al., 2023) and Biomed-CLIP (Zhang et al., 2023a) models on the overall accuracy. It also demonstrates competitive performance with existing approaches on closed-ended questions. However, it does not perform as well on open-ended questions, where it falls behind compared to the MedVInt-TD model. We argue that the lower performance on open-ended questions can be attributed to several factors. One key factor is our formulation of the problem as free-form answer generation for both question types, as opposed to the baseline Biomed-CLIP and PubMedCLIP models. This means that our model is not constrained by a predefined set of answers in the training data.

To evaluate the influence of domain adaptation in the VLM, we performed experiments using two LMs, Bloomz-7b1 and RadBloomz-7b. The comparison results in Table 2 demonstrate that

---

[1] https://docs.python.org/3/library/difflib.html

| Model | VE | LM | SLAKE 1.0 | | | VQA-RAD | | |
|---|---|---|---|---|---|---|---|---|
| | | | Overall | Closed | Open | Overall | Closed | Open |
| Ours | BiomedCLIP ViT | RadBloomz-7b | **87.5** | **92.1** | **84.5** | 73.2 | 83.5 | 57.5 |
| Ours | PMC-CLIP ResNet50 | RadBloomz-7b | 82.5 | 88.5 | 78.6 | 67.6 | 79.4 | 49.7 |
| MedVInT-TD (Zhang et al., 2023b) | | | 85.2 | 86.3 | 84.5 | **81.6** | **86.8** | **73.7** |
| Biomed-CLIP (Zhang et al., 2023a) | | | 86.1 | 88.9 | 84.3 | 72.7 | 76.5 | 67.0 |
| PubMedCLIP (Eslami et al., 2023) | | | 80.1 | 82.5 | 78.4 | 72.1 | 80.0 | 60.1 |
| MUMC (Li et al., 2023c) | | | 84.9 | - | - | 79.2 | 84.2 | 71.5 |
| PMC-CLIP (Lin et al., 2023a) | | | 84.3 | 88.0 | 81.9 | 77.6 | 84.0 | 67.0 |

Table 1: Accuracy (%) results of VLMs on SLAKE 1.0-English and VQA-RAD datasets. Performance on open-ended and closed-ended questions as well as overall performance are reported. VE represents vision encoder.

| VE | LM | SLAKE 1.0 | | | VQA-RAD | | |
|---|---|---|---|---|---|---|---|
| | | Overall | Closed | Open | Overall | Closed | Open |
| BiomedCLIP ViT | Bloomz-7b1 | 80.0 | 86.8 | 75.7 | 68.3 | 80.9 | 49.2 |
| | Radbloomz-7b | **87.5** | **92.1** | **84.5** | **73.2** | **83.5** | **57.5** |
| PMC-CLIP ResNet | Bloomz-7b1 | 80.5 | 87.5 | 76.0 | 65.2 | 77.9 | 45.8 |
| | Radbloomz-7b | **82.5** | **88.5** | **78.6** | **67.6** | **79.4** | **49.7** |

Table 2: The table compares the accuracy (%) between a VLM with a radiology-adapted RadBloomz-7b LM and a general-domain Bloomz-7b1 LM, using the SLAKE 1.0-English and VQA-RAD datasets. Results for open-ended, closed-ended, and overall performance are included, with experiments conducted separately using two pretrained vision encoders (VE).

BiomedCLIP-RadBloomz-7b outperforms its general domain language model counterpart, Bloomz-7b1, on both datasets. There is a noticeable enhancement in overall accuracy on Slake 1.0, with an improvement of 7.5%. Similarly, on VQA-RAD, there is a significant increase in overall accuracy, with an improvement of 4.9%. This highlights the significant benefit of employing a domain-adapted language model, specifically RadBloomz-7b, as the backend language model for domain-intensive tasks in VLMs. The model's effectiveness is particularly evident in its performance on open-ended questions, demonstrating an average improvement of 8.5% in accuracy.

To evaluate the impact of including training of existing parameters in the fusion model, we conducted experiments on VLMs that employed trainable vision encoders. In this regard, we trained the vision encoder parameters alongside other trainable parameters throughout all training stages. Table 3 shows the results obtained from the VLMs using trainable BiomedCLIP-ViT. The two LMs, Bloomz and RadBloomz, were utilized in the experiments. Notably, the VLM utilizing the specialized-domain RadBloomz-7b achieves better performance with a reduced number of parameters compared to the VLM with a larger set of trainable parameters. We

argue that through an optimal fusion of the domain-adapted vision encoder and LM, there is no longer a need to train the vision encoder in our VLM. This results in a lightweight adaptation of the VLM.

To assess the effect of three different training stages on model performance, we explore the following scenarios: 1) Direct Fine-tuning, where the model is exclusively trained on VQA-RAD or SLAKE 1.0 datasets without any prior training phases. 2) One-stage Pre-Training, which includes pre-training stage 1, followed by fine-tuning on downstream datasets. 3) Full Pre-Training, where the model undergoes all three training stages. This comparison offers valuable insights into the most effective training pathway for this model architecture in domain-intensive MedVQA tasks.

Table 4 shows the comparison results with BiomedCLIP-RadBloomz-7b. The findings reveal significant improvements in final accuracy, with an approximate 25% increase in full pre-training (Scenario 3) compared to direct fine-tuning (Scenario 1). These results underscore the effectiveness of Pre-training stage 1, which greatly enhances the model's medical knowledge. Furthermore, full pre-training not only preserves the knowledge gained during stage 1 but also integrates medical concept alignment with specialized MedVQA training.

| VE | LM | Overall | Closed-ended | Open-ended |
|---|---|---|---|---|
| Trained BiomedCLIP ViT | Bloomz-7b1 | 69.4 | 80.1 | 53.1 |
| Frozen BiomedCLIP ViT | Bloomz-7b1 | 68.3 | 80.9 | 49.2 |
| Trained BiomedCLIP ViT | RadBloomz-7b | 71.4 | 81.3 | 56.4 |
| Frozen BiomedCLIP ViT | RadBloomz-7b | 73.2 | 83.5 | 57.5 |

Table 3: The table provides a comparison of accuracy (%) between two scenarios on the VQA-RAD dataset: one scenario where the vision encoder of VLMs is trained alongside alignment training, and another where the vision encoder is frozen during training. The table displays performance for open-ended and closed-ended questions, as well as overall performance.

| Scenarios | Overall | Closed-ended | Open-ended |
|---|---|---|---|
| 1 | 48.3 | 59.9 | 30.7 |
| 2 | 59.0 | 70.6 | 41.3 |
| 3 | 73.2 | 83.5 | 57.5 |

Table 4: The table demonstrates the performance of our VLM (BiomedCLIP ViT+Radbloomz-7b) on VQA-RAD under different training scenarios: 1) direct fine-tuning on VQA-RAD; 2) stage 1 pretraining followed by fine-tuning on VQA-RAD; and 3) full pre-training and fine-tuning on VQA-RAD. The accuracy metric is used, and performance is reported for open-ended, closed-ended questions, along with overall accuracy.

| Category | #Q | Bloomz-7b1 | RadBloomz-7b |
|---|---|---|---|
| Abnormality | 56 | 64.3 | **69.6** |
| Attribute | 20 | 90.0 | 90.0 |
| Color | 4 | 100.0 | 100.0 |
| Count | 6 | 66.7 | **83.3** |
| Modality | 33 | 45.5 | **48.5** |
| Organ | 10 | 20.0 | **40.0** |
| Plane | 26 | 73.1 | **76.9** |
| Position | 61 | 72.1 | 70.5 |
| Presence | 171 | 74.9 | **82.5** |
| Size | 46 | 87.0 | 82.6 |
| Other | 26 | 30.8 | 26.9 |

Table 5: Models' overall accuracy (%) across different question categories on VQA-RAD. Performance of two VLMs with Radbloomz-7b and Bloomz-7b1 as LLM component is reported separately. The vision encoder of VLMs is BiomedCLIP ViT. #Q: number of questions in the given category.

We examine the overall accuracy of VLMs using BiomedCLIP-ViT as the vision encoder across different question categories in both datasets. The results can be found in Tables 5 and 6. Our VLM with medical-tailored Radbloomz-7b shows better performance in most categories. RadBloomz-7b particularly excels in interpreting spatially-oriented queries, as evident from its leading performance in modality, abnormality, presence of objects/attributes, organ, and plane categories. This suggests a strong capability of RadBloomz-7b in analyzing the spatial arrangement in radiology images. However, the model can be further improved in shape, size, and position categories. Additionally, the distribution of categories in the training data has an impact on the model's performance.

Finally, we conduct a qualitative analysis of the model's predictions to identify areas where improvements may be needed for both the model and evaluation measures. Table 7 shows examples of questions from the VQA-RAD test split where the model's predictions are evaluated as incorrect during the evaluation. Notably, despite the model's responses being evaluated as incorrect according to our evaluation measure, a closer examination reveals a different perspective. The model provided responses that consist of terms that are either synonyms or contextually relevant to the given labels.

For instance, in question 1, the model identifies the modality as 'chest x-ray', which is essentially correct in the context of this question (See Figure 2). Similarly, for question 2, the model's prediction 't2 weighted' captures the essence of the 't2 weighted mri' label or in question 4, 'both sides' is predicted whereas the label is 'both'.

Given that traditional accuracy metrics may not fully capture the nuances and utilization of synonyms in the medical domain, conducting a manual evaluation of the predictions can be valuable in determining the actual performance of the model. However, it is worth noting that we have identified instances where the model generated incorrect answers, such as in questions 6 and 7. We asked a licensed medical expert to meticulously compare the model's predictions with the ground truth values and identify cases similar to those mentioned earlier. Following this rigorous human evaluation, we achieved an accuracy of 64.2%, surpassing the performance obtained using our automatic evaluation metric, which yielded an accuracy of 57.5%.

Although BiomedCLIP-RadBloomz-7b VLM demonstrates remarkable overall improvement in

(a) question 1      (b) question 4      (c) question 5

Figure 2: Image examples from VQA-RAD corresponding to questions in Table 7.

| Category | #Q | Bloomz-7b1 | RadBloomz-7b |
|---|---|---|---|
| Organ | 253 | 88.9 | **93.6** |
| Abnormality | 150 | 73.3 | **84.6** |
| Size | 65 | 86.1 | **87.6** |
| Position | 186 | 67.2 | **87.6** |
| Plane | 58 | 96.5 | **100.0** |
| Modality | 108 | 100.0 | 100.0 |
| Knowledge Graph | 148 | 68.9 | **75.0** |
| Color | 34 | 88.2 | **91.1** |
| Quantity | 52 | 59.6 | 59.6 |
| Shape | 7 | 85.7 | 71.4 |

Table 6: Model's overall accuracy (%) across different question categories on SLAKE 1.0-English. Performance of two VLMs with Radbloomz-7b and Bloomz-7b1 as LLM component is reported separately. The vision encoder of VLMs is BiomedCLIP ViT. #Q: the number of questions in the given category.

MedVQA, additional investigation of the model is necessary. Specifically, since the task is formulated as free-form generation, training a model to adhere to a restricted set of terminologies presents challenges and warrants further attention.

## 4 Background and Related Work

Language models (LMs) designed for general domains often face difficulties when applied to highly specialized fields. Additionally, data scarcity is a prevalent challenge in domain adaptation of LMs. Various methods have been developed to adapt pre-trained LMs to specific domains. One method involves continuous pre-training of model parameters using data specific to the target domain (Karn et al., 2023). Alternatively, synthetic data can be effectively incorporated into the training process for fine-tuning models to better adapt to specific target domains (Karn et al., 2021). Another approach includes using parameter-efficient fine-tuning methods (Xu et al., 2023) with task-specific training data. Our training schema amalgamates several of these methods like image-caption pretraining, synthetic biomedical MQA, and task-specific Med-VQA datasets, all fine-tuned using a parameter-efficient technique.

Among parameter-efficient fine-tuning approaches, the Low-Rank Adaptation (LoRA) technique (Hu et al., 2021) has received considerable interest for adapting Large LMs (LLMs). In the biomedical domain, domain-specific LLMs have been proposed either by fine-tuning the model's parameters (Luo et al., 2022; Wu et al., 2023) or by utilizing LoRA techniques (Gema et al., 2023). However, it's important to note that biomedical domain-adapted LLMs might not perform as effectively in the radiology domain. This is due to the complexity of terminologies in clinical NLP (Karn et al., 2022; Ghosh et al., 2023). Thus, there have been recent proposals for radiology domain-adapted LLMs (Karn et al., 2023).

The application of domain adaptation is not limited to LLMs. It also finds utility in the adaptation of multi-modal models like vision-language models (VLMs). In line with this, there have been recent proposed biomedical VLMs such as (Zhang et al., 2023a; Lin et al., 2023a; Moor et al., 2023; Chen et al., 2023; Li et al., 2023a). These have been successful in achieving state-of-the-art performance in downstream biomedical tasks, such as medical question-answering. In this study, we concentrate on developing a more efficient domain adaptation technique for VLMs within the challenging domain of Radiology.

## 5 Conclusion

We introduce a new vision-language model for medical visual question-answering by integrating a radiology large language model, RadBloomz-7b (Karn et al., 2023) and a biomedical vision encoder, BiomedCLIP-ViT (Zhang et al., 2023a), in to the VLM. Our main objective is to investigate the impact of integrating specialised LMs and vision encoders into VLMs for domain-specific tasks in the medical domain.

For this purpose, we propose a parameter-efficient training approach by deploying low-rank adaptation technique (Hu et al., 2021) to the

| | Question | Label | Prediction |
|---|---|---|---|
| 1 | What kind of image is this? | x-ray | chest x-ray |
| 2 | What type of MRI sequence is displayed in this image? | t2 weighted mri | t2 weighted |
| 3 | What modality was used? | plain film | plain film xray |
| 4 | Are pleural opacities located on the left, right, or both sides of the lung? | both | both sides |
| 5 | Are there multiple or just 1 metastatic focus? | one | just one |
| 6 | Which lung is clearer? | left | right |
| 7 | Is the anatomy of the brain gyri affected? | no | yes |

Table 7: Examples of our model's generated answers (Prediction) on closed- and open-ended questions in VQA-RAD evaluated as incorrect answer.

decoder-only LLM component in the VLM, which significantly reduces the number of trainable parameters while maintaining the model performance. Moreover, the vision encoder is kept frozen in the training process. We then propose a two-stage pre-training approach aiming to align our VLM to medical concepts by pre-training the model on the image-captioning task and acquiring general knowledge for medical visual question answering by pre-training it on a general MedVQA dataset. We finally finetune the model on the downstream MedVQA tasks.

Our results demonstrate state-of-the-art performance on a MedVQA SLAKE 1.0 dataset and strong performance on the VQA-RAD dataset. Furthermore, compared to a VLM with a general-domain LLM, we show that our proposed VLM leads to a higher performance using parameter-efficient training, while a VLM with general-domain LM benefits slightly from training the vision encoder as well. Finally, our findings suggest that the proposed pre-training approach significantly improves model performance in downstream MedVQA tasks.

## 6 Limitations

In this paper, we explored the generation ability of our adapted vision-language model on learning to generate free-form answers. While we observed impressive performance, we realized that in a few test cases, such as wh-questions, the model generates *yes/no* answers. Therefore, more investigation on optimizing the training to capture the type of the question is required.

We proposed a multi-modal model tailored for radiology-domain visual-question answering tasks. Therefore, we are aware that our model is not easily generalizable to diverse medical domains and tasks,

such as pathology image analysis. As a result, we didn't compare our model to SoTA generalized multi-modal models in other medical domains and tasks. Furthermore, the LLM model architecture we studied is restricted to a decoder-only type, thus its performance may not be directly comparable to different model architectures.

## 7 Ethics Statement

All datasets in this paper are publicly available for clinical NLP research. Trained models for Med-VQA tasks in this paper must be assessed carefully before considering them for final applications.

## 8 Disclaimer

The concepts and information presented in this paper are based on research results that are not commercially available. Future commercial availability cannot be guaranteed.

## 9 Acknowledgement

## References

Qiuhui Chen, Xinyue Hu, Zirui Wang, and Yi Hong. 2023. Medblip: Bootstrapping language-image pre-training from 3d medical images and texts.

Sedigheh Eslami, Christoph Meinel, and Gerard de Melo. 2023. PubMedCLIP: How much does CLIP benefit visual question answering in the medical domain? In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1181–1193, Dubrovnik, Croatia. Association for Computational Linguistics.

Aryo Gema, Luke Daines, Pasquale Minervini, and Beatrice Alex. 2023. Parameter-efficient fine-tuning of llama for the clinical domain. *arXiv preprint arXiv:2307.03042*.

Rikhiya Ghosh, Oladimeji Farri, Sanjeev Kumar Karn, Manuela Danu, Ramya Vunikili, and Larisa Micu. 2023. RadLing: Towards efficient radiology report understanding. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 640–651, Toronto, Canada. Association for Computational Linguistics.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison.

Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2020. Mimic-iv. *PhysioNet. Available online at: https://physionet. org/content/mimiciv/1.0/(accessed August 23, 2021)*.

Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. 2019. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

Sanjeev Kumar Karn, Francine Chen, Yan-Ying Chen, Ulli Waltinger, and Hinrich Schütze. 2021. Few-shot learning of an interleaved text summarization model by pretraining with synthetic data. In *Proceedings of the Second Workshop on Domain Adaptation for NLP*, pages 245–254, Kyiv, Ukraine. Association for Computational Linguistics.

Sanjeev Kumar Karn, Rikhiya Ghosh, Kusuma P, and Oladimeji Farri. 2023. shs-nlp at RadSum23: Domain-adaptive pre-training of instruction-tuned LLMs for radiology report impression generation. In *The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, pages 550–556, Toronto, Canada. Association for Computational Linguistics.

Sanjeev Kumar Karn, Ning Liu, Hinrich Schuetze, and Oladimeji Farri. 2022. Differentiable multi-agent actor-critic for multi-step radiology report summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1542–1553, Dublin, Ireland. Association for Computational Linguistics.

Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. 2018. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023a. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*.

Pengfei Li, Gang Liu, Jinlong He, Zixu Zhao, and Shenjun Zhong. 2023c. Masked vision and language pre-training with unimodal and multimodal contrastive losses for medical visual question answering. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2023: 26th International Conference, Vancouver, BC, Canada, October 8–12, 2023, Proceedings, Part I*, page 374–383, Berlin, Heidelberg. Springer-Verlag.

Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023a. PMC-CLIP: contrastive language-image pre-training using biomedical documents. In *Medical Image Computing and Computer Assisted Intervention - MICCAI 2023 - 26th International Conference, Vancouver, BC, Canada, October 8-12, 2023, Proceedings, Part VIII*, volume 14227 of *Lecture Notes in Computer Science*, pages 525–536. Springer.

Zhihong Lin, Donghao Zhang, Qingyi Tao, Danli Shi, Gholamreza Haffari, Qi Wu, Mingguang He, and Zongyuan Ge. 2023b. Medical visual question answering: A survey. *Artificial Intelligence in Medicine*, page 102611.

Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. 2021. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1650–1654. IEEE.

Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. BioGPT: Generative Pre-trained Transformer for Biomedical Text Generation and Mining. *Briefings in Bioinformatics*, 23(6):bbac409.

Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. 2018. Mixed precision training. In *International Conference on Learning Representations*.

Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Cyril Zakka, Yash Dalmia, Eduardo Pontes Reis, Pranav Rajpurkar, and Jure Leskovec. 2023. Med-flamingo: A multimodal medical few-shot learner. *arXiv preprint arXiv:2307.15189*.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.

Obioma Pelka, Sven Koitka, Johannes Rückert, Felix Nensa, and Christoph M Friedrich. 2018. Radiology objects in context (roco): a multimodal image dataset. In *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis: 7th Joint International Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*, pages 180–189. Springer.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.

Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Further fine-tuning llama on medical papers. *arXiv preprint arXiv:2304.14454*.

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*.

Sheng Zhang, Yanbo Xu, Naoto Usuyama, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, et al. 2023a. Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv preprint arXiv:2303.00915*.

Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023b. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*.

## A MedVQA datasets statistics

Table 8: Downstream dataset statistics of VQA-RAD and SLAKE 1.0, includes number of images and question-answer pairs (QAs). Questions are categorized as close-ended and open-ended.

| Dataset | VQA-RAD | | | SLAKE 1.0-English | | | |
|---|---|---|---|---|---|---|---|
| | Total | Train | Test | Total | Train | Validation | Test |
| #Images | 315 | 314 | 203 | 642 | 586 | 174 | 96 |
| #QAs | 3515 | 3064 | 451 | 12995 | 9835 | 2099 | 1061 |
| #Close-ended QAs | 2093 | 1821 | 272 | 5141 | 3881 | 844 | 416 |
| #Open-ended QAs | 1420 | 1241 | 179 | 7754 | 5854 | 1255 | 645 |

# LLM-Based Section Identifiers Excel on Open Source but Stumble in Real World Applications

**Saranya Krishnamoorthy, Ayush Singh, Shabnam Tafreshi**

inQbator AI at eviCore Healthcare

Evernorth Health Services

`firstname.lastname@evicore.com`

## Abstract

Electronic health records (EHR) even though a boon for healthcare practitioners, are growing convoluted and longer every day. Sifting around these lengthy EHRs is taxing and becomes a cumbersome part of physician-patient interaction. Several approaches have been proposed to help alleviate this prevalent issue either via summarization or sectioning, however, only a few approaches have truly been helpful in the past. With the rise of automated methods, machine learning (ML) has shown promise in solving the task of identifying relevant sections in EHR. However, most ML methods rely on labeled data which is difficult to get in healthcare. Large language models (LLMs) on the other hand, have performed impressive feats in natural language processing (NLP), that too in a zero-shot manner, i.e. without any labeled data. To that end, we propose using LLMs to identify relevant section headers. We find that GPT-4 can effectively solve the task on both zero and few-shot settings as well as segment dramatically better than state-of-the-art methods. Additionally, we also annotate a much harder real world dataset and find that GPT-4 struggles to perform well, alluding to further research and harder benchmarks.

## 1 Introduction

Modern day healthcare systems are increasingly moving towards large scale adoption of maintaining electronic health records (EHR) of patients (Congress, 2009). EHRs help healthcare practitioners with relevant information about a patient such as history, medications, etc. However, in recent times this practice has led to very long and convoluted EHRs (Rule et al., 2021). Naturally, the need for better information retrieval tools emerged due to the progressively lengthy and unstructured doctor notes. One such need is the accurate identification of sections in an EHR, pertinent to a physician's inquiry. For instance, a question like "What



Figure 1: Sample real world obscure image of an outpatient paper-based patient encounter form comprising of numerous sections (Hersh and Hoyt, 2018).

treatments has the patient undergone in the past?" concerning prior treatments administered to a patient necessitates the swift extraction of information from the "treatments" and "past medical history" sections, while excluding sections related to "ancestral medical history". This swift extraction is vital for timely decision-making in patient care. Additionally, during critical procedures such as the evaluation of medical necessity for prior authorization requests, it is customary for experienced clinicians to locate vital data within specific sections. An illustrative case entails examining the "physical exam" section to identify particular findings, such as signs of neurological disorders or movement-associated pain, indicating the need for additional diagnostic tests. The timely identification of such information is of utmost importance in ensuring the provision of appropriate care and reducing the risk of potential complications.

In general, regions found in EHR would often

258

have a section heading preceding the body of the section, as can be seen in example Table 1. Even though these section types have limited cardinality, however, more often than not, physicians would fail to adhere to standards and use lexical variations generated on the fly. Moreover, practitioners not only will generate lexical variations of sections on the fly but also completely new sections altogether for valid reasons like imaging reports, etc. Apart from these variations, oftentimes there would be no headers at all, even though the information present could ideally be part of a pre-existing section in a document or a new section altogether. While studies like Gao et al. (2022) utilize the Subjective, Objective, Assessment and Plan heading (SOAP) framework, real-world clinical notes often contain sections beyond these categories. This limitation is further emphasized in Landes et al. (2022), warranting further investigation and analysis.

The aforementioned factors have consequently contributed to the establishment of Section Identification (SI) as a distinct and enduring problem within the academic discourse (McKnight and Srinivasan, 2003), making it an indispensable component of any clinical natural language processing (NLP) pipeline. A SI task entails finding regions of text that are semantically related to an aspect of a patient's medical profile. More importantly, it helps to improve pre-existing information retrieval systems by enabling them to be more targeted and specific. Lastly, in light of recent findings of the negative impact of note bloat within EHRs on even the most sophisticated systems (Liu et al., 2022), using SI to shorten or create from EHR, a sub-EHR specific to a given task would prove to be a worthwhile effort for humans and machines both.

Because finding sections and hence their corresponding headers involves inherent variability, machine learning (ML) methods have played an important role in this natural language processing (Pomares-Quimbaya et al., 2019). ML has increasingly been shown to be efficient in finding relevant sections within a document, however, a key drawback of traditional ML methods has been the dependence on labeled data (Tepper et al., 2012). Reliance on annotated data for training ML models to be able to predict the beginning and end of section headers has stalled the field from fully solving the task. The emergence of large language models (LLMs) in contemporary research presents a promising avenue to overcome the limitations inherent in traditional machine learning approaches,

thereby expanding the scope of their applications.

LLMs have emerged as the de-facto system for NLP in scenarios where data is scarce (OpenAI, 2023). The key distinction between traditional Machine Learning (ML) models and Large Language Models (LLMs) lies in their ability to understand tasks in natural language. While traditional ML models require labeled data for training, LLMs can leverage pre-training on vast amounts of unstructured text data, enabling them to perform tasks with minimal task-specific fine-tuning. This makes ML possible in an unsupervised manner (no need for labeled data) and therefore opens room for applications in domains where annotated data is hard to acquire like healthcare. While LLMs have been evaluated on a wide array of NLP tasks in healthcare (Nori et al., 2023), they are yet to be evaluated on their effectiveness in segmenting a document into semantically relevant sections.

In this work, we address this gap and evaluate the efficacy of our approach on a widely-known datasets in the clinical medical domain. Findings show that GPT-4 (OpenAI, 2023) almost solved the section identification problem on the benchmark open-sourced dataset, however, on a private dataset the performance lags. Our contributions are threefold, listed as follows:

1. We show that GPT-4 can generate zero-shot headings of records with very high accuracy.

2. Contrary to the above, we find that its performance drops on internal real-world datasets.

3. An ontology of numerous section headers seen in real world EHR systems is shared which has much higher coverage.

## 2  Related Work

Traditionally, SI task has been done using a pre-defined dictionary of plausible candidates. Pomares-Quimbaya et al. (2019) performed a comprehensive survey and found that rule-based methods still dominated the array of methods proposed while ML systems increasingly achieved better coverage when combined in a hybrid manner with rule-based methods. McKnight and Srinivasan (2003) later on extracted bag-of-words from MedLINE abstracts and used a support vector machine to train a classifier to categorize sentences into either Introduction, Method, Result, or Conclusion, demonstrating promising results. Similarly, Hirohata et al.

| Allergies | Allergies: Patient recorded as having No Known Allergies to Drugs... |
| --- | --- |
| **History of Present Illness** | HPI: 61M w/ incidental L renal mass found during W/U for brachytherapy for low-grade [**Last Name (STitle) **], now w/ gradually worsening gross hematuria for the past several days. |
| **Labs Imaging** | Pertinent Results: [**2160-4-10**] 07:30AM BLOOD WBC-12.6* RBC-3.20* Hgb-8.2* Hct-24.5* MCV-77* MCH-25.6* MCHC-33.4 RDW-17.1* Plt Ct-438. |
| **Hospital Course** | Brief Hospital Course: 61M w/ low-grade [**Month/Day/Year **] awaiting brachytherapy and locally-advanced L renal mass w/ collecting system invasion, renal vein thrombus, and likely metastases, presented w/gradually worsening gross hematuria. |

Table 1: This figure illustrates a sample data point from the MIMIC-III database, highlighting the sections annotated with MedSecID corpus.

(2008) achieved very high accuracy by using conditional random fields to label scientific abstracts into Objectives, Methods, Results, and Conclusions.

Over time and with the inclusion of ML, the field re-framed this problem as one of span-level entity identification i.e. the system would be tasked with predicting whether each token in a sequence belongs to one of the predefined section types using the Inside-Outside-Beginning (IOB) tagging system (Ramshaw and Marcus, 1999). Tepper et al. (2012) addresses the task of segmenting clinical records into distinct sections using a two-step approach. First, the section boundaries are identified. Then, the sections are passed to the second step, where a classifier is used to label each token as *Begin*, *In* or *Out* of the span of a section. Nair et al. (2021) proposes several transfer learning models based on clinical contextual embeddings for classifying clinical notes into the major SOAP sections (Podder et al., 2023). Zhou et al. (2023) investigates the effectiveness of continued pre-training in enhancing the transferability of clinical note section classification models. Both of the above papers resemble our work, however, they restrict them to SOAP sections and train specific models to do so. While the techniques devised so far have shown promise, to the best of our knowledge none of the previous works have tried in an unsupervised manner.

With the advent of LLMs (Devlin et al., 2018; OpenAI, 2023), several works have shown the efficacy of LLMs in doing unsupervised zero-shot information extraction. The primary method for interacting with generative LLMs is by the use of natural language prompts. Wei et al. (2022) found a significant performance boost by asking the model to explain its chain of thought before answering the query. Further, Brown et al. (2020) showed that additional performance can be gained by passing some examples as part of the prompt, they named it

Few-Shot prompting. Wang et al. (2023); Bian et al. (2023); Ashok and Lipton (2023) have shown the efficacy of prompting the LLM to extract biomedical named entities from scientific articles. More recently, Liu et al. (2023) used GPT-4 to de-identify documents in a zero-shot manner. This hints at the immense document understanding capabilities of LLMs and opens doors to its application to a wide array of previously unresolved tasks such as SI.

Apart from the advancements in the field of ML and SI, to evaluate how well SI systems perform, a standardization of tasks as well as datasets is required. To that end, Uzuner et al. (2011) first proposed a SI task as part of Informatics for Integrating Biology and the Bedside (i2b2) benchmarks. Recently, Landes et al. (2022) argued that the previous dataset did not fully cover the nuances in SI task and proposed a dataset an order of magnitude larger as well as more comprehensive than one by Uzuner et al. (2011). However, the dataset proposed by Landes et al. (2022) is based on a clean source Johnson et al. (2016), which oftentimes is not the case in real-world scenarios. To that end, we also annotated a real-world dataset to evaluate LLMs on it as well.

## 3 Datasets

### 3.1 i2b2 2010

In their study, Tepper et al. (2012) meticulously curated a corpus comprising 183 annotated clinical notes extracted from a selection of discharge summaries within the i2b2 2010 (Uzuner et al., 2011) dataset. This dataset was annotated by an expert and served as a valuable resource for their research. However, owing to constraints imposed by Institutional Review Boards (IRBs), our current access to the i2b2 2010 dataset is limited. As a result, we were only able to procure clinical notes for 96 out of the originally annotated 183 documents.

| Dataset | MedSedId | i2b2 2010 | Real World |
|---|---|---|---|
| Document count | 2002 | 96 | 100 |
| Average token length | 2307 | 1283 | 7841 |
| Std. dev. token length | 1732 | 726 | 8093 |
| Average sections per doc | 12 | 17 | 12 |
| Std. dev. sections per doc | 5.7 | 6.2 | 8 |

Table 2: Corpus Statistics

## 3.2 MedSecID

MedSecID (Landes et al., 2022) is a publicly available corpus of 2,002 fully annotated medical notes from the MIMIC-III (Johnson et al., 2016) clinical record database. Each note has been manually annotated with section boundaries and section labels (See Table 1 for an example of a typical clinical note consisting of well-defined sections). The section labels correspond to different types of information that are typically found in clinical notes, such as history of present illness, physical exam findings, and progress notes.

## 3.3 Real-world

In an increasingly digital world, one would be inclined to assume healthcare data also lives digitally. Surprisingly, that is not the case almost 75% of the healthcare dataset still lives in faxes (CCSI, 2022) (see figure 1 for a sample handwritten and faxed clinical notes). Whereas all preexisting SI datasets are digitally derived from clean EHR systems, which even though offer us some insight into the performance of state of art, however, fail to paint the full picture. Therefore, we use an internal dataset of prior authorization requests derived from faxed-in images being transcribed to text via an optical character recognition system (OCR). These requests contain EHR of patients in the form of doctors' notes, submitted in both PDF and image formats. These documents lack a standardized structure, with segments and titles that can vary significantly in length. Although it's possible to group these titles into clusters of similar meaning, the language and number of titles differ across documents. Additionally, OCR inaccuracies arise from unclear text, spelling errors, complex table structures, and handwritten content, resulting in highly noisy input for any SI system to process.

## 4 Annotation Methods

In this section, we describe the dataset and the annotation design in our study. As we described before we decided to choose section identification (SI), a method to identify sections and sub-sections in EHR documents to split them into smaller text chunks and create some structure in these unstructured data. We designed a manual annotation task to identify these sections and create categorical section types. Below we explain the annotation task design, the result, and the challenges.

### 4.1 Annotation Design

We randomly selected 100 records from a pool of one million records we have in our corpus. These records are in two forms, PDF or fax images which doctors submit to insurance companies, and hence, can arrive from any arbitrary format. We refer to these records as documents in the span of this manuscript. These documents have no standard structures and sometimes they contain multiple patients information at the same time. Six annotators with higher education and non-native speakers of English carry the annotation task. Each annotates an equal amount and random selection of these documents.

We used Label Studio[1], an open source data labeling platform. PDF or image file of each record is uploaded to label studio and the task was to mark the section and sub-section in each file and manually enter the corresponding text of these sections and sub-sections. To instruct the annotators, we provided written instructions as well as held a video discussion session and explained the task to the annotators.

### 4.2 Annotation Result

We aggregate the sections per document to form the final section and sub-section list. A total of 912 sections and subsections are identified which makes 14 sections and sub-sections on average per document. Then one annotator, different from the ones who have annotated the documents, categorized these sections and sub-sections into more gen-

---

[1]https://labelstud.io/

261

Figure 2: Section categories which are selected based on observation of top-header sections in the corpus and human judgment to associate section names to their topic or category of representations.

eral categories based on the Consolidated Clinical Document Architecture (C-CDA) implementation guide[2]. In other words, the diverse categories are mapped to a category to unify them. This allows us to calculate IAA and be able to use the text semantic similarity method to find these sections in the unannotated documents. A total of 464 categories are coded of which 394 of these categories have a frequency of 1 and 70 categories have a frequency of 2 or more. We provide a small sample of the most frequent categories in Table 3 and Figure 2.

24 documents have been randomly selected and on each of these documents, a second annotator annotated the document. Further, we calculated the Jaccard similarity to report Inter-Annotator Agreement (IAA), The Jaccard similarity is a measure of the similarity between two sets of data. We obtained a Jaccard distance of 0.40, which is a fair agreement and an indication that the annotation task is challenging. The most diverse section and sub-section lists that each normalized into one section name are shown in table 4. Notably, the diversity of these two general categories indicates the challenge involved in structuring and identifying these sections in these documents. In some cases, categories such as *Order Report* or *Medication Reconciliation* can be both a section and sub-section according to the annotation results. This characteristic does not enforce the decision to select the general category for these types.

---

[2]C-CDA contains a library of CDA templates, incorporating and harmonizing previous efforts from Health Level Seven (HL7), Integrating the Healthcare Enterprise (IHE), and Health Information Technology Standards Panel (HITSP). https://www.hl7.org/ccdasearch/

## 5 Experimental Setup

Our task here is to take as input a document and output all the section headers found in it. For our underlying use case, we carried out testing with various LLMs like GPT-4 8k (OpenAI, 2023), LLaMa-2 7B (Touvron et al., 2023), and more recent Mistral 7B (Jiang et al., 2023) prompting strategies[3] (as shown in figure 3) and contrasted them with a baseline experiment that used keyword search, regex, MedSpacy library (Eyre et al., 2021) and the best model reported by Landes et al. (2022). MedSpacy is a clinical NLP toolkit built on the foundation of SpaCy, specifically designed to address the unique challenges of processing and extracting information from clinical text. This enables healthcare professionals to efficiently process and derive valuable insights from unstructured medical narratives. We did not restrict the tokens and used the entire clinical note for MedSecId. We extracted the actual section header using the header span mentioned in the MedSecId annotation and used it as the ground truth for our task. Because of the longer length of real-world data, we used the 32k version of GPT-4 while keeping all the hyper-parameters to default such as the temperature, frequency penalty, and presence penalty to 0 and max tokens to 1000. Lastly, in this study, we utilized a privately hosted instance of GPT-4 to ensure the prevention of any potential data leakage. Prior to initiating the experiment, we implemented a thorough anonymization procedure to protect the dataset Protected health information (PHI). This involved substituting all

---

[3]CoT A5, One Shot A4 and Close Ended A6 prompting strategies are elaborated in appendix A.

| | |
|---|---|
| **Medications Section** | Information about the current and past Medications |
| **Order Info** | This section consists of additional items that are required to conclude the assessments. Examples of such items are Mammograms, x-rays, etc., or the information about the provider of such items. |
| **Results Section** | Usually contains of lab results |
| **Physical Exam Section** | Result of physical exams such as Integumentary, Chest and Lung Exam, Cardiovascular, Abdomen, etc. |

Table 3: A sample of sections and subsections with the highest frequency.

| | |
|---|---|
| **Medications Section** | Medications, Medication Changes, Medication List at End of Visit, Medication, Medication Reconciliation, Preventive Medicine, Medication List, Medication List at End of Visith, Medications (active prior today), Medications (Added, Consumed or Stopped today), Medications (Added, Continued or Stopped today), Medications Changes, Medications Discontinued During This Encounter, Medications Ordered This Encounter, Medications Places This Encounter, MEDICATIONS PRESCRIBED THIS VISIT, Medications Reviewed As Of This Encounter, Meds, Outpatient Medications, Patients Medication, Preventive Medication, Previous Medications, Previous medications |
| **Order Info** | Orders Placed, Order Questions, Order, Order Details, Order Information, Order Providers, Order Report, Ordering Provider, Order Name, Order name, Order Number, Order Plain X-ray/Interpretation, Order Requisition, Order Tracking, Order Transmittal Tracking, Order User/Provider Detail, Order-Level Documents, Ordering Provider Information, Orders, Orders Placed This Encounter, Orders Requiring a Screening Form |

Table 4: The list of sections and subsections that are normalized into one section name.

---

You are a clinician and you read the given clinical document and identify section headers from them. Find section headers only from the clinical text.
For each section header, return the answer as a JSON object by filling in the following dictionary.
{section_title: string representing the section header}
Here are some clinical notes of a patient from a doctor. ### {*context_text*} ###

---

Figure 3: Basic Prompt Template

personal identifiers, such as names, identification numbers, and ages, with fictitious entities.

Apart from the basic prompts, we also experiment with combining them with Few-Shot (Brown et al., 2020) and CoT Prompting (Wei et al., 2022) where we ask the LLM to think step-by-step along with providing an example of the clinical note and a list of headings. We keep the prompts same across all the datasets. Lastly, the evaluation metric used here is the exact match (EM) accuracy as well as precision (P), recall (R), and F1-score calculated by comparing GPT-4's output to that of ground truth in the Inside-Outside-Beginning (IOB) scheme (Ramshaw and Marcus, 1999) as used in work by Landes et al. (2022). Similar GPT-4 experiments were conducted on i2b2 2010 dataset but as the context length of i2b2 was smaller, in all the experiments we use GPT-4 8K. Lastly, because of cost constraints, we chose the best-performing model on above mentioned benchmarks to be eval-

uated against our internal real-world dataset.

## 6 Results

Even though GPT-4 was able to perform very well on open source benchmark datasets, it was unable to reach the same level of performance on our internal corpus due to its complexity as shown in table 7. Experiments showed that GPT-4 was able to achieve an accuracy of only 37% in contrast to that of 96% on MedSecId corpus. LLaMa-2 and MedSpacy performed equally well, in that, former achieved higher recall than latter. This can be attributed to the global knowledge encoded in the LLMs, which is not the case with MedSpacy, while on the other hand MedSpacy would be much faster to run with less overhead. Results in table 5 and 6 show that one-shot GPT-4 OpenAI (2023) performed the best and achieved a new state of the art on MedSecId outperforming previous models by a significant margin. This unsupervised methodology

| Method | Accuracy(%) | Precision(%) | Recall(%) | F1(%) | EM(%) |
|---|---|---|---|---|---|
| Keyword Based | 36.07 | 100 | 36.07 | 53.01 | 36.05 |
| Regex | 49.24 | 100 | 30.07 | 46.24 | 50.8 |
| MedSpacy | 56.63 | 100 | 38.29 | 55.38 | 62.63 |
| GPT-4 Close Ended Prompt | 73.23 | 100 | 73.23 | 84.55 | 73.2 |
| GPT-4 Chain-of-Thought (CoT) | 94.9 | 100 | 88.62 | 93.97 | 92.47 |
| GPT-4 Zero Shot Prompt | 94.41 | 100 | 87.61 | 93.40 | 92.05 |
| GPT-4 One Shot Prompt | **96.86** | 100 | **92.93** | **96.24** | **96.11** |
| LLaMa-2 Close Ended Prompt | 39.96 | 100 | 39.96 | 57.10 | 39.94 |
| LLaMa-2 Zero Shot Prompt | 52.29 | 94.61 | 32.92 | 48.82 | 62.25 |
| LLaMa-2 One Shot Prompt | 13.95 | 94.57 | 6.86 | 12.80 | 16.86 |
| LLaMa-2 Chain-of-Thought (CoT) | 38.21 | 93.95 | 21.11 | 34.48 | 46.95 |
| Mistral Close Ended Prompt | 5.24 | 100 | 5.24 | 9.96 | 5.24 |
| Mistral Zero Shot Prompt | 11.51 | 97.43 | 5.23 | 9.93 | 14.45 |
| Mistral One Shot Prompt | 8.41 | 98.61 | 4.07 | 7.82 | 10.48 |
| Mistral Chain-of-Thought (CoT) | 11.99 | 98.61 | 5.64 | 10.67 | 15.53 |
| BiLSTM-CRF (Landes et al., 2022) | 82.2 | **95** | 95 | 95 | - |

Table 5: Results on MedSecId Corpus

| Method | Accuracy(%) | Precision(%) | Recall(%) | F1(%) | EM(%) |
|---|---|---|---|---|---|
| Keyword Based | 10.98 | 100 | 8.78 | 16.14 | 69.5 |
| Regex | 66.26 | 100 | 48.27 | 65.11 | 56.8 |
| MedSpacy | 38.45 | 100 | 21.92 | 35.96 | 38.14 |
| GPT-4 Close Ended Prompt | 11.82 | 78.24 | 8.46 | 15.27 | 73.8 |
| GPT-4 Chain-of-Thought (CoT) | 86.26 | 99.85 | 74.65 | 85.43 | 84.33 |
| GPT-4 Zero Shot Prompt | 89.47 | 100 | 78.46 | 87.93 | 84.58 |
| GPT-4 One Shot Prompt | **93.03** | **100** | **85.36** | **92.10** | **89.45** |
| LLaMa-2 Close Ended Prompt | 88.79 | 100 | 83.57 | 91.05 | 86.54 |
| LLaMa-2 Zero Shot Prompt | 56.2 | 100 | 36.62 | 53.61 | 58.59 |
| LLaMa-2 One Shot Prompt | 30.54 | 100 | 16.75 | 28.69 | 21.2 |
| LLaMa-2 Chain-of-Thought (CoT) | 40.23 | 99.83 | 22.61 | 36.87 | 50.7 |
| Mistral Close Ended Prompt | 10.41 | 100 | 6.65 | 12.48 | 19.34 |
| Mistral Zero Shot Prompt | 35.30 | 100 | 18.98 | 31.90 | 36.17 |
| Mistral One Shot Prompt | 6.58 | 100 | 3.24 | 6.29 | 7.80 |
| Mistral Chain-of-Thought (CoT) | 32.13 | 99.80 | 17.03 | 29.09 | 33.66 |
| Maximum Entropy (Tepper et al., 2012) | - | 91.1 | 90.8 | 91 | - |

Table 6: Results on i2b2 Corpus. While GPT-4 has superior performance, LLaMa-2 is not far behind.

| Method | A | P | R | F1 | EM |
|---|---|---|---|---|---|
| Regex | **67.64** | 98.69 | **51.30** | **67.51** | **71.9** |
| MedSpacy | 5.92 | 100 | 4.13 | 7.93 | 15.72 |
| GPT-4 ZS | 37.53 | 100 | 24.18 | 38.95 | 37.29 |
| LLaMa-2 ZS | 13.33 | 100 | 7.81 | 14.49 | 19.75 |
| Mistral ZS | 3.67 | 100 | 1.83 | 3.60 | 5.24 |

Table 7: Results on Real-World Corpus. ZS stands for Zero-Shot prompting

beats all the supervised models on the MedSecId corpus (Landes et al., 2022). Similarly, one-shot also had a state-of-the-art performance on i2b2 2010 dataset. On the other hand, LLaMa-2 did not perform as well as GPT-4, but nevertheless had on par performance with regex. Additionally, LLaMa-2 Touvron et al. (2023) performance on i2b2 dataset came very close to that of GPT-4 itself. This disparity in performance of LLaMa-2 as well as its variation in results across the experi-

ments leads to inconclusive results. Lastly, Mistral (Jiang et al., 2023) performance was sub-optimal, exhibiting only a marginal improvement than a naive keyword based approach.

## 7 Discussion

We performed an in-depth error analysis on the subset of records that GPT-4 was unable to predict correction. Our analysis found errors in the Med-SecId dataset itself, which is one of the reasons GPT-4 did not get a 100% performance. Error analysis reveals on the rest of 2.8% missed sections of the GPT-4 finds that 18% of the above stated 2.8% belong to the "Findings" section label and 13% belong to the "Image-Type" category. Most of the documents did not have those section headers explicitly mentioned and were hidden as part of the text. Even though the precision was 100% in i2b2 2010 dataset, the granularity of the subsections, the

| Section Categories | Number of Sections in Category | Frequency | Frequency (%) |
|---|---|---|---|
| Assessment & Plan | 413 | 958 | 60.98 |
| physical exam | 66 | 152 | 9.67 |
| Personal Info | 54 | 73 | 4.64 |
| Medication | 19 | 55 | 3.50 |
| History of Present Illness | 3 | 44 | 2.80 |
| Family History | 5 | 40 | 2.54 |
| Allergies | 4 | 40 | 2.54 |
| Order Info | 17 | 38 | 2.41 |
| Clinical Info | 16 | 36 | 2.29 |
| UNKNOWN | 13 | 25 | 1.59 |
| Additional Info | 4 | 18 | 1.14 |
| Appointment Date | 6 | 15 | 0.95 |
| Progress Notes | 1 | 15 | 0.95 |
| Results | 7 | 12 | 0.76 |
| Mental Status | 6 | 10 | 0.65 |
| History | 3 | 10 | 0.64 |
| Lab Results | 5 | 6 | 0.38 |
| Alcohol Use | 2 | 5 | 0.31 |
| Abdomen | 2 | 5 | 0.31 |
| Referral | 3 | 3 | 0.19 |
| Active Medication | 3 | 3 | 0.19 |
| References | 2 | 3 | 0.19 |
| Miscellaneous | 2 | 2 | 0.12 |
| All Reviewer List | 2 | 2 | 0.12 |
| Return Visit | 1 | 1 | 0.06 |

Table 8: Each section name is categorised to either its top-header section or a category is selected by human to represent the topic of the section. This annotation is done manually by two annotators where one selected a course-grained categories and the other selected a fine-grained categories. The one we show in this table is the coarse-grained category list, along with the number of of sections in each category, frequency, and frequency percentage. When the annotator were not able to asses a category they mark the section as *UNKNOWN*

presence of ambiguous language, or the lack of clear markers for section boundaries could be the contributors to the slight dip in recall of the section headers. We leave fixing the issues in the dataset and advanced prompting for future work.

Surprisingly, we found that GPT-4 was even able to extract sub-sections that were missed in the human annotations in MedSecId. This raises the question of whether GPT-4's superior performance on these datasets can be attributed to its prior exposure to them? We found out that MedSecId is derived from MIMIC dataset which forbids being used for LLM training, therefore, it is highly unlikely it was used during model training.

Further analysis of our internal dataset revealed that high variation in the structure of the document is the root cause of such a wide gap between benchmark and our internal datasets. The original version of our data is in the form of images and PDF files. While GPT was resilient to most OCR errors it did contribute to some misspelled sections. We acknowledge the difference in GPT's and the gold standard's approach to section title extraction. While the gold standard highlights literal text, GPT summarizes the content, potentially providing a more concise and informative overview. Example GPT output *Patient Information and Visit Details* encompasses multiple headers like *Chief Complaint, History of Present Illness*, and *Patient Information*. GPT also extracted irrelevant titles as section headers *Provider Information and Signature, Page Footer*, etc. We aim to work on addressing these issues by incorporating context awareness into the title-generation process.

The major challenge in performance drop on internal dataset is due to the nature of our data itself. More specifically, there is neither standard structure nor format. The situation exacerbates with the document being an out of an OCR system which introduces numerous morphological errors. Consequently, GPT-4's responses on our dataset are more creative and semantically similar which is something an exact match evaluation is unable to measure. As zero-shot was performing extremely well on public corpus and the improvement with other prompting techniques gave only minor improvements, we conducted only zero shot on our internal datasets.

Apart from conducting experiments on the state of art LLMs like GPT-4 (OpenAI, 2023), we also

wanted to experiment with smaller open-source models that offer flexibility. We experimented with two of the best-performing models LLaMa-2 (Touvron et al., 2023) and Mistral (Jiang et al., 2023). However, in reality, both the open source models found it hard to follow the prompts and the outputs are not consistent. The challenges were further exacerbated when the models were required to generate results in a uniform format. Sometimes, both LLaMa-2 and Mistral would just output the summarization of the text. LLaMa-2 demonstrated a significantly superior performance than Mistral on both i2b2 and MedSecID.

Further, each section name is categorised to either its top-header section or a category is selected by human to represent the topic of the section. This annotation is done manually by two annotators where one selected a course-grained category list and other selected a fine-grained one. The one we show in table 8 is the coarse-grained category list, along with the number of sections in each category, frequency, and frequency percentage. 25 categories are created by the annotator to represent the coarse-grained categories. There are some section names that both annotators are unable to assess or select a category. These sections are categorized as *UNKNOWN*. If we consider that the top nodes in an ontology network, on average each node will have 26 child nodes in this ontology.

## 8 Conclusion

In this work, we evaluated LLMs capabilities in segmenting a clinical document into individual sections. More specifically, we show that an unsupervised GPT-4 can nearly solve the Section Identification task. Even though GPT-4 has a very high accuracy on the benchmark datasets, however, its performance on a real-world dataset has a significant lag. We further analyze the reasons for such a wide gap and find that the source dataset has cleanly defined section headers which is not the case with its real-world counterpart. To show how diverse the real-world dataset is, we further derived an ontology using another set of annotators that we share with the community at large.

To that end, we create a harder benchmark, one that is derived from real-world data generating process. Moreover, we conducted an annotation study with five annotators to create the final dataset and found high ambiguity in the identification of headers on the newly introduced benchmark. As a take-

away, we suggest that if the source dataset or EHR is clean, then there is no need anymore to train specific supervised models to detect sections as an unsupervised LLM can perform that task.

## 9 Future Work

After realizing the close-to-perfect performance and poor performance on the internal real world dataset of an unsupervised LLM in this study, we believe currently released datasets do not paint a clear picture of how the techniques proposed so far would perform in real world scenarios. Using our own internal dataset, we would like to fine-tune the LLM to see whether it can improve performance in a way that is comparable to open-source. Lastly, because sharing sensitive patient data is not possible, we plan to work on de-identifying and training an LLM to generate synthetic but realistic datasets which could lead to better real world benchmarks.

## 10 Limitations

One of the self-evident limitations of our approach is the reliance on GPT-4 to perform SI task. Using GPT-4 incurs both high overhead costs and significant data leakage risks if not set up properly. Therefore, the technique itself cannot be run in an isolated environment as it depends on an external API. Another drawback common with ML systems is if tomorrow new sections emerge and GPT-4 is not updated, the if will fail to capture the new section types.

## 11 Ethics

The datasets used in the study involved sensitive patient data. Therefore, we decided not to disclose the internal data. Additionally, even for the data based on MIMIC (Johnson et al., 2016), we used a privately hosted instance of GPT-4 that sits in a HIPAA compliant environment. Separately, the annotators were provided fully de-identified data, and the identification of the annotators themselves was anonymized during the annotation process. We have released the taxonomy at our github[4] and kindly request the community to report any further advancements to us via email.

---

[4] https://github.com/inQbator-eviCore/LLM_section_identifiers

# References

Dhananjay Ashok and Zachary C Lipton. 2023. Promptner: Prompting for named entity recognition. *arXiv preprint arXiv:2305.15444*.

Junyi Bian, Jiaxuan Zheng, Yuyi Zhang, and Shanfeng Zhu. 2023. Inspire the large language model by external knowledge on biomedical named entity recognition.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

CCSI. 2022. Six healthcare workflows primed for cloud faxing. https://healthitsecurity.com/news/six-healthcare-workflows-primed-for-cloud-faxing. Accessed: 2023-12-15.

US Congress. 2009. Hr 1: American recovery and reinvestment act of 2009. *Washington, DC (February 2009)*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

H. Eyre, A. B. Chapman, K. S. Peterson, J. Shi, P. R. Alba, M. M. Jones, T. L. Box, S. L. DuVall, and O. V. Patterson. 2021. Launching into clinical space with medspaCy: a new clinical text processing toolkit in Python. *AMIA Annu Symp Proc*, 2021:438–447.

Yanjun Gao, Dmitriy Dligach, Timothy Miller, Samuel Tesch, Ryan Laffin, Matthew M Churpek, and Majid Afshar. 2022. Hierarchical annotation for building a suite of clinical natural language processing tasks: Progress note understanding. In *LREC... International Conference on Language Resources & Evaluation:[proceedings]. International Conference on Language Resources & Evaluation*, volume 2022, page 5484. NIH Public Access.

William R Hersh and Robert E Hoyt. 2018. *Health Informatics: Practical Guide Seventh Edition*. Lulu. com.

Kenji Hirohata, Naoaki Okazaki, Sophia Ananiadou, and Mitsuru Ishizuka. 2008. Identifying sections in scientific abstracts using conditional random fields. In *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-I*.

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

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

Paul Landes, Kunal Patel, Sean S Huang, Adam Webb, Barbara Di Eugenio, and Cornelia Caragea. 2022. A new public corpus for clinical section identification: Medsecid. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3709–3721.

Jinghui Liu, Daniel Capurro, Anthony Nguyen, and Karin Verspoor. 2022. "note bloat" impacts deep learning-based nlp models for clinical prediction tasks. *Journal of Biomedical Informatics*, 133:104149.

Zhengliang Liu, Xiaowei Yu, Lu Zhang, Zihao Wu, Chao Cao, Haixing Dai, Lin Zhao, Wei Liu, Dinggang Shen, Quanzheng Li, Tianming Liu, Dajiang Zhu, and Xiang Li. 2023. Deid-gpt: Zero-shot medical text de-identification by gpt-4.

Larry McKnight and Padmini Srinivasan. 2003. Categorization of sentence types in medical abstracts. In *AMIA 2003, American Medical Informatics Association Annual Symposium, Washington, DC, USA, November 8-12, 2003*. AMIA.

Namrata Nair, Sankaran Narayanan, Pradeep Achan, and KP Soman. 2021. Clinical note section identification using transfer learning. In *Proceedings of Sixth International Congress on Information and Communication Technology: ICICT 2021, London, Volume 1*, pages 533–542. Springer.

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.

Vivek Podder, Valerie Lew, and Ghassemzadeh Sassan. 2023. *SOAP Notes*. StatPearls Publishing.

Alexandra Pomares-Quimbaya, Markus Kreuzthaler, and Stefan Schulz. 2019. Current approaches to identify sections within clinical narratives from electronic health records: a systematic review. *BMC Medical Research Methodology*, 19.

Lance A Ramshaw and Mitchell P Marcus. 1999. Text chunking using transformation-based learning. In *Natural language processing using very large corpora*, pages 157–176. Springer.

Adam Rule, Steven Bedrick, Michael F. Chiang, and Michelle R. Hribar. 2021. Length and Redundancy of Outpatient Progress Notes Across a Decade at an Academic Medical Center. *JAMA Network Open*, 4(7):e2115334–e2115334.

Michael Tepper, Daniel Capurro, Fei Xia, Lucy Vanderwende, and Meliha Yetisgen-Yildiz. 2012. Statistical section segmentation in free-text clinical records. In *Lrec*, pages 2001–2008.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Özlem Uzuner, Brett R South, Shuying Shen, and Scott L DuVall. 2011. 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text. *Journal of the American Medical Informatics Association*, 18(5):552–556.

Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*.

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.

Weipeng Zhou, Meliha Yetisgen, Majid Afshar, Yanjun Gao, Guergana Savova, and Timothy A Miller. 2023. Improving model transferability for clinical note section classification models using continued pretraining. *medRxiv*.

# A  Appendix

Figure 4 illustrates an example of "One Shot" prompt method. It contains the segmentations and the seed list of heading found in MedSecId. In the end we present the entire patient notes received from the doctors. Figure 5 shows an example of "CoT" prompt. We observe that in this method the prompt should instruct the LLMs to think rationally and ask them to extract the section headers from the patient notes. Lastly, figure 6 shows an example

of "Close Ended" prompt method. This method restricts the responses to be one of the 50 class labels that is obtained from MedSecId annotation.

Table 9 demonstrates the top 50 populated section names that we observed in our corpus. The numbers are extracted from the aggregated annotation results. We observe that "Allergies", "Family History", and "Social History" are top 3 populated sections in the corpus. The full list is published in our GitHub which is provided in section 11.

Figure 2 shows the sections categories. The annotation is done by two annotators. One annotator chooses course-grained categories and the other chooses more fine-grained categories. These categories are selected based on observation of top-header sections in the corpus and human judgment to associate these section names to their topic or category of representations. Our findings show that "Assessment & Plan" is the most populated category with 958 sections and "Return Visit" us the least populated one with only 1 section. The sections are extracted from the aggregated annotation result of our study. Statistics such as number of sections per category, frequency, and frequency percentage is shown in Table 8.

You are a clinician and you read the given clinical document and identify section headers from them.
Find section headers only from the clinical text.
Example clinical text: {sample_text}
Answer { List of section headers from the corpus. }
For each section header return the answer as a JSON object by filling in the following dictionary.
{section_title: string representing the section header}
Here are some clinical notes of a patient from a doctor. ### {*context_text*} ###

Figure 4: One Shot Prompt: provide examples of segmentation as well as provide a seed list of headings found in MedSecId.

You are a clinician and you read the given clinical document and identify section headers from them.
Find section headers only from the clinical text.
For each section header, return the answer as a JSON object by filling in the following dictionary.
{section_title: string representing the section header
CoT: string describing thinking step by step }
Here are some clinical notes of a patient from a doctor. ### {*context_text*} ###

Figure 5: CoT Prompt: make the LLM think rationally and try to extract all possible section headers in the clinical notes

You are a clinician and you read the given clinical document and identify section headers from them.
Classify the section headers into one of the following section type labels.
section types: {List of section types from the MedSecId training corpus.}
If the section headers do not belong to any of the above section type labels, classify them as Ńone.
Only print the section types identified in a list. Here are some clinical notes of a patient from a doctor.
### {*context_text*} ###

Figure 6: Close Ended Prompt: restrict the responses to one of the 50 class labels obtained from the MedSecId annotation.

| Section Names | Frequency | Percentage (%) |
| --- | --- | --- |
| Allergies | 36 | 2.3% |
| Family History | 36 | 2.3% |
| Social History | 34 | 2.2% |
| Past Medical History | 29 | 1.9% |
| Physical Exam | 28 | 1.8% |
| Subjective | 25 | 1.6% |
| Objective | 24 | 1.5% |
| Plan | 24 | 1.5% |
| Surgical History | 24 | 1.5% |
| HPI | 23 | 1.5% |
| Assessment | 21 | 1.3% |
| Chief Complaint | 20 | 1.3% |
| History of Present Illness | 20 | 1.3% |
| Review of Systems | 19 | 1.2% |
| Impression | 17 | 1.1% |
| Medications | 16 | 1.0% |
| Vital signs | 16 | 1.0% |
| Additional Documentation | 15 | 1.0% |
| Progress Notes | 15 | 1.0% |
| ROS | 14 | 0.9% |
| Medication Changes | 13 | 0.8% |
| Orders Placed | 13 | 0.8% |
| Visit Diagnoses | 13 | 0.8% |
| Assessment/Plan | 12 | 0.8% |
| Current Medications | 11 | 0.7% |
| Past Surgical History | 11 | 0.7% |
| Vitals | 11 | 0.7% |
| Assessments | 10 | 0.6% |
| Examination | 10 | 0.6% |
| Musculoskeletal | 10 | 0.6% |
| Problems | 10 | 0.6% |
| Technique | 10 | 0.6% |
| Communications | 9 | 0.6% |
| Comparison | 9 | 0.6% |
| Exam | 9 | 0.6% |
| Findings | 9 | 0.6% |
| Reason for Appointment | 9 | 0.6% |
| Diagnosis | 8 | 0.5% |
| Medical History | 8 | 0.5% |
| Medication List at End of Visit | 8 | 0.5% |
| Screening | 8 | 0.5% |
| Skin | 8 | 0.5% |
| Cardiovascular | 7 | 0.4% |
| General | 7 | 0.4% |
| History | 7 | 0.4% |
| Tobacco Use | 7 | 0.4% |
| Treatment | 7 | 0.4% |
| Eyes | 6 | 0.4% |
| Instructions | 6 | 0.4% |
| Patient Information | 6 | 0.4% |

Table 9: Top 50 Sections Names quantified by their frequencies and percentages in the entire corpus. We observe that "Allergies", "Family History", and "Social History" are top 3 most populated sections in the corpus.

# Adapting Abstract Meaning Representation Parsing to the Clinical Narrative – the SPRING THYME parser

**Jon Z. Cai[1], Kristin Wright-Bettner[1]**
**Martha Palmer[1], Guergana K. Savova[2], James H. Martin[1]**
[1]University of Colorado Boulder
[2]Boston Children's Hospital and Harvard Medical School

## Abstract

This paper is dedicated to the design and evaluation of the first AMR parser tailored for clinical notes. Our objective was to facilitate the precise transformation of the clinical notes into structured AMR expressions, thereby enhancing the interpretability and usability of clinical text data at scale. Leveraging the colon cancer dataset from the Temporal Histories of Your Medical Events (THYME) corpus, we adapted a state-of-the-art AMR parser utilizing continuous training. Our approach incorporates data augmentation techniques to enhance the accuracy of AMR structure predictions. Notably, through this learning strategy, our parser achieved an impressive F1 score of 88% on the THYME corpus's colon cancer dataset. Moreover, our research delved into the efficacy of data required for domain adaptation within the realm of clinical notes, presenting domain adaptation data requirements for AMR parsing. This exploration not only underscores the parser's robust performance but also highlights its potential in facilitating a deeper understanding of clinical narratives through structured semantic representations.

## 1 Introduction

Abstract Meaning Representation (Banarescu et al., 2013)(AMR)is a highly adaptable and expressive framework designed to capture the semantics of natural language expressions. Automatic AMR parsing is a natural language processing (NLP) method that translates natural language inputs into formal AMR expressions – representations which have proven to be useful across a wide range of downstream applications (Kapanipathi et al., 2021; Liu et al., 2015; Liao et al., 2018; Li and Flanigan, 2022; Bonial et al., 2020; Bai et al., 2021) including those in the biomedical domain (Garg et al., 2016; Rao et al., 2017).

Formally, AMR expressions take the form of labeled, rooted, directed, and acyclic graphs, $g = $ $(V, E)$, where $V$ represents the set of AMR nodes, which can be of type predicate, abstract concept and attributes; $E$ represents the possible semantic relations between nodes such as prototypical agent and patient denoted by `arg0` and `arg1`. The AMR graph structure underpinned by Neo-Davidsonian semantics can then effectively encapsulate the abstract concepts, relationships, and entities present in individual sentences or utterances.

From a practical standpoint, AMR expressions encompass the semantic content typically addressed by individual representation schemes such as semantic role labeling (Palmer et al., 2005), named entities (Wang et al., 2022), and coreference chains (Joshi et al., 2020), thereby unifying these diverse aspects of meaning into a single comprehensive representation. Figure 1 illustrates an AMR expression selected from the clinical domain.

As Figure 1 demonstrates, concepts including events, entities and properties are captured as nodes in the graph, while the relations among the concepts are captured by labeled edges connecting the nodes. Events are represented using PropBank frames (Palmer et al., 2005), and the semantic relations of both entities and events to these predicates are specified either by a frame's numbered argument or one of the relations from AMR's role inventory. For example, the see-09 predicate represents the event of "visit/consultation by a medical professional." In this case, the agent of the seeing event is "Dr. Chandler Bing", represented by see-09's `ARG0` relation, and the semantic role of patient for the event is "she" indicated by the `ARG1` semantic relation. AMR graphs also specify the temporal information in a formal way. In the above example, the time of the seeing event is specified by two temporal modifier subgraphs. It is a conjunction of "after now" and "within this week" which makes "later this week" a concrete time range.

AMR parsers based on pretrained large lan-

271

Figure 1: the AMR graph of sentence "We will have her see Dr. Chandler Bing in surgical consultation later this week following her testing."

guage models and sequence-to-sequence (encoder-decoder) architectures have demonstrated impressive accuracy when trained and evaluated on standard datasets. The use of AMR parsers has contributed to improved performance across a range of NLP tasks including question answering (Fu et al., 2021), information retrieval (Liao et al., 2018), knowledge-graph construction (Ribeiro et al., 2022), and text generation (Bai et al., 2022).

These successes have sparked growing interest in employing AMR in domains that diverge from the existing training data, such as human-robot interaction tasks, educational applications involving classroom discourse analysis, and diverse biomedical use cases. Unfortunately, as language form and meaning deviate from the general language captured in generic training data, parsing performance shows a rapid decline. This decline stems from disparities in vocabulary, syntax, and overall discourse structure. Addressing these challenges necessitates dedicated human expert annotation efforts to create domain-specific AMR resources. However, such endeavors can be costly and time-consuming. Hence, the preference lies in maximizing the utilization of existing data and parsers and adapting them to new domains, rather than building entirely new systems from scratch.

The contributions of this paper include:

- We adapted the high-performance SPRING parser (Bevilacqua et al., 2021) to the clinical domain, specifically leveraging the Temporal Histories of Your Medical Events (THYME) corpus (Wright-Bettner et al., 2020), and achieved state-of-the-art performance in AMR parsing within this context..

- We demonstrated that by tailoring an existing general domain English neural AMR parser with a relatively modest amount of gold-standard in-domain data, we could attain significantly high accuracy.

- We showcased data augmentation techniques that effectively enhance the parser's robustness across different domains.

## 2 Data

Supervised training data for AMR parsers consists of pairs of linguistic expressions along with their associated human annotated gold-standard AMR expressions. The current standard dataset for AMR development is AMR 3.0 (Knight et al., 2020) available from the Linguistic Data Consortium as LDC2020T02. This general domain dataset is the basis for our baseline efforts prior to domain adaptation. AMR 3.0 consists of over 59k English expressions from a variety of broadcast conversations, newswire, weblogs, web discussion forums, fiction and web text. To facilitate evaluation and model comparison, AMR 3.0 is divided into standard training, development and test splits consisting of 55,635, 1,722, and 1,898 expressions respectively.

To adapt AMR to the clinical narrative, we developed 8,327 in-domain AMRs (separate paper with detailed description under review) on a subset of the THYME colon cancer corpus (Styler et al., 2014; Wright-Bettner et al., 2020). The colon cancer part of the THYME corpus consists of 594 de-identified physicians' notes for 198 patients with colon cancer. Each patient is represented by one pathology note and two clinical notes. The corpus has undergone several prior annotation efforts, including temporal and coreference annotation (Styler et al., 2014; Wright-Bettner et al., 2019, 2020) and entity tagging as defined by the Unified Medical Language System (UMLS (Bodenreider, 2004)). As part of our AMR annotation process, we adopted seven clinical-domain named entity

(NE) types (anatomical-site, clinical-attribute, devices, disease-disorder, medications-drugs, sign-symptom) from the UMLS project and relied heavily on the UMLS in classifying many AMR concepts.

Like other genre-specific AMR tasks (Bonial et al., 2019; Bonn et al., 2020), we found it necessary to modify the standard AMR annotation approach to support meaningful annotation of domain-unique linguistic phenomena. Two phenomena are pervasive in the clinical narrative. First, physician notes frequently drop eventive mentions when they are inferable by human readers. For example, "Declines tetanus" does not mean the patient declined having tetanus; they declined a tetanus immunization. We expanded AMR's guidelines to permit explicit rendering of certain implicit concepts like the immunization:

```
(d / decline-02
     :ARG1 (s / shot-13 :implicit +
          :ARG3 (d2 / disease-disorder :
   name (n / name :op1 "tetanus"))))
```

Second, like other specialized domains, clinical texts are rife with semantically dense noun phrases (NPs) (Grön et al., 2018). In AMR, NPs must be treated in one of two ways: Either all components are extracted and related (white marble = marble that is white), or they are analyzed as single units of meaning, i.e., NEs (White House). However, semantic compositionality exists on a spectrum (Nakov, 2013), and many specialized NPs in particular strain the adequacy of a binary approach. This can be seen even in simple clinical NPs: One annotator might decide "blood pressure" is a single, cohesive unit of meaning and annotate it as an NE, while another might decide "pressure" is an extractable property of "blood". To address this, we implemented a two-pass strategy: In the first pass, for NPs that fell under one of the clinical NE types mentioned above, an experienced annotator made these compositionality judgments and added each unique phrase to a searchable, phrasal NE Dictionary along with an AMR fragment that "defined" the compositionality for each phrase. Annotators then referenced the Dictionary when building the AMR graphs in the second pass. This approach supported consistency and speed of annotation.

Finally, the THYME corpus contains frequent repetition of many other multiword expressions and phrases. For extremely formulaic phrases, such as those found in Vital Signs sections (Height = 167.60 cm, e.g.), we implemented a template-filling

script that deterministically produced the AMRs, again saving significant manual annotation time. Of the 8,327 AMRs, 1,640 were produced by this script; the rest were created manually. The final 8,327 THYME-AMR data are split into training, development and test sets randomly with 4,955, 1,641 and 1,731 sentence-AMR pairs, respectively. All of the model training is conducted on the training set of the AMR 3.0 and THYME AMR corpora. We show the Inter Annotator Agreement between three annotators on 107 THYME-AMRs in Table 1

| Comparison | P | R | F1 |
|---|---|---|---|
| gold vs annotator 1 | 0.93 | 0.93 | 0.93 |
| gold vs annotator 2 | 0.93 | 0.93 | 0.93 |
| annotator 1 vs annotator 2 | 0.91 | 0.90 | 0.90 |

Table 1: Smatch scores on 107 manuall THYME AMRs, representing three clinical notes

## 3 Methods

We treat the AMR parsing task as a supervised machine learning problem and train a parameterized model to map natural language expressions to their corresponding AMR graphs. Various model architectures and training methods and paradigms have been employed over the years (Flanigan et al., 2014; Foland and Martin, 2017; Lyu and Titov, 2018; Cai and Lam, 2019; Zhang et al., 2019; Wang et al., 2015; Ballesteros and Al-Onaizan, 2017; Fernandez Astudillo et al., 2020; Hoang et al., 2021), resulting in a continuous improvement in the state of the art on the general domain AMR dataset(i.e. AMR 2.0 and 3.0 corpus (LDC2020T2)). However, these improvements are highly dependent on the availability of significant amounts of annotated training data hampering the development of parsers for specific genres and languages other than English. Our approach here is to leverage an existing high-performance parser and adapt it to the clinical domain using the modest amount of domain-specific training data described in the last section.

Meanwhile, the great advances of the pre-trained foundational models has introduced a new modeling paradigm in the field of NLP as well as to structure-prediction problems such as AMR parsing. In particular, the sequence-to-sequence modeling, originally developed for machine translation, has proven a highly effective approach for AMR parsing (Bevilacqua et al., 2021; Konstas et al., 2017; Xu et al., 2020). In this approach, two neu-

Figure 2: AMR graph to PENMAN linearization pipeline. The transformation map between the AMR graphical representation and its linearized representation is one-to-one-and-onto.



Figure 3: The SPRING parser modeling diagram. A transformer-based self-attention mechanism is used to produce embeddings for the input expression. The decoder then uses cross attention to drive autoregressive generation of a sequence of AMR output tokens.

ral network components are involved: an encoder, which takes the natural language sentence as input and maps it to a continuous manifold as a sequence of high-dimensional vectors, and a decoder, which takes the embedded sentence representation vectors and maps them to the output embedding space, corresponding to the target sequence tokens.

Here we make use of the SPRING parser (Bevilacqua et al., 2021), one of the state-of-the-art AMR parsers on AMR 3.0 evaluation. The underlying pre-trained language model is BART-large (Lewis et al., 2020), a transformer-based language model that has been trained using a set of denoising pre-training objectives, such as a masked language modeling objective and a document reconstruction objective, on general domain unlabeled English text. The neural network architecture relies on the self-attention and cross-attention mechanism to learn patterns from natural language texts. This pre-trained model is then fine-tuned on the AMR 3.0 training data to map English inputs to linearized AMR graphs, which consist of a sequence of AMR tokens. We show the linearization correspondence of an AMR graph to its sequence of AMR tokens in Figure 2.

A critical aspect of using sequence-to-sequence models for structured prediction tasks, like parsing, is transforming the task itself. In AMR parsing, the AMR graph is converted into a sequence of tokens through a linearization algorithm. Note that the vocabulary of the decoder differs from that of the encoder model, as the target sequence consists of AMR-specific tokens such as the relations `arg0` and `arg1`, and predicates like `test-01`. During fine-tuning, we utilize the vocabulary derived from the AMR 3.0 corpus, which ensures consistency and accuracy in the parsing process. The parsing problem is then to convert an input text sequence into a valid sequence of AMR tokens that can be deterministically transformed into a directed AMR graph. The overall SPRING approach is depicted in Figure 3. Given a high-performing SPRING model, we adapt it to the THYME domain by fine-tuning on the THYME-AMR training set (4,955 expressions). Here, fine-tuning involves continuous gradient-based updates to the original model parameters with a small learning rate ($5 \times 10^{-6}$) with batch size to be 20, we keep the maximum sequence length to be 1024..

## 3.1 Evaluation

The standard metric to evaluate AMR parsing performance is SMATCH, which decomposes an AMR graph into triples that capture the edge list representation of a graph structure. For instance, the AMR for the sentence "He had never undergone a screening colonoscopy." can be decomposed into its edge list representation as AMR1 and edge list 1 as follows:

```
AMR1:
(c / colonoscopy-01 :polarity -
      :arg1 (h / he)
      :arg2 (s2 / screen-01
            :arg1 h))

AMR2:
(c1 / colonoscopy-01 :polarity -
      :arg1 (s / she)
      :arg2 (s2 / screen-01
            :arg1 s))
```

```
Decomposed edge list1:
      instance(c, colonoscopy-01)
      instance(h, he)
      instance(s2, screen-01)
      polarity(c, -)
      arg1(c, h)
      arg2(c, s2)
      arg1(s2, h)

Decomposed edge list2:
      instance(c1, colonoscopy-01)
      instance(s, she)
      instance(s2, screen-01)
      polarity(c1, -)
      arg1(c1, s)
      arg2(c1, s2)
      arg1(s2, s)
```

We conjured another slightly altered AMR2 with the `he` node replaced with a `she` node, indicating a potential mistake in the parser generated AMR. In the above decomposition of AMR graphs, `instance()` represents the nodes in the graph while the rest are the edges. Given the edge lists for a hypothetical parse and its corresponding gold-standard parse, the SMATCH metric produces precision (p), recall (r), and F1-measure scores as follows:

$$p = \frac{N_{correct}}{N_{predicted}}, r = \frac{N_{correct}}{N_{reference}}, F_1 = \frac{2pr}{p+r}$$

A complication in computing these scores is that we need to know which of the proposed AMR nodes in the parse are supposed to correspond to which ones in the correct set. In other words, the graphs need to be matched before they can be scored. This issue originates from the encoding of AMR nodes with variables, through which different instantiations of a concept can be encoded. The standard SMATCH scorer (Cai and Knight, 2013) employs a greedy heuristic method to provide the required alignment to avoid computing a computationally expensive optimal alignment.

Finally, AMR representations are an amalgamation of semantic representations including predicate-argument relations, named entities, and coreference components. The SMATCH score represents an average over these component categories, obscuring the model performance over the various categories of information in AMR expressions, thus making it difficult to assess the usability of the results in downstream applications. To address this, a more fine-grained analysis tool[1] provides precision, recall and F1 measures across the various component AMR tasks. We will discuss the fine-grained categories in section 4.3.

## 4 Experiments

We present the domain adaption training experiments in this section to show the characteristics of the text from THYME corpus when it comes to AMR parser developement.

### 4.1 Domain Adaptation

Table 2 provides the results of our primary domain adaptation experiments. The first column presents the evaluation results of the off-the-shelf SPRING AMR parser trained solely with the AMR 3.0 training data. The 83.0 SMATCH score for the SPRING parser reaches near state-of-the-art performance on the AMR 3.0 test set, whereas, the performance on the THYME-AMR test set is significantly lower at 51.7 SMATCH. The second column shows the results of the same parser fine-tuned using the THYME-AMR training data. Here, we see that the fine-tuned parser achieves excellent results on the THYME-AMR corpus test set with a 35.3 point absolute improvement over the original model.

| Test \ Train | AMR 3.0 | THYME-AMR | AMR 3.0 + THYME-AMR |
|---|---|---|---|
| AMR 3.0 | 83.0 | 77.0 | 80.0 |
| THYME-AMR | 51.7 | 87.0 | 88.0 |

Table 2: SPRING THYME-AMR parser performance with different training sources. All scores are Smatch F1

---

[1] https://github.com/mdtux89/amr-evaluation

## 4.2 Avoiding Forgetting

Catastrophic forgetting is a frequently observed problem when fine-tuning large pre-trained models on domain specific data (Li and Hoiem, 2018; Riemer et al., 2019; Scialom et al., 2022). While fitting the model's parameters to the new domain, there is often a significant loss in terms of the model's performance on its original domain. To assess the robustness and potential forgetting of general domain AMR knowledge, we evaluated the THYME-AMR fine-tuned parser on the AMR 3.0. The results showed a decrease in performance from 83.8 to 77.0, indicating significant forgetting of the general domain AMR.

Based on this observation, we deployed a joint training approach to mitigate this forgetting phenomenon. In this experiment, we fine-tuned the parser on a mixture sampled from both the AMR 3.0 and THYME-AMR data. Considering the differing sizes of the two corpora, we sampled them in a 12-to-1 ratio between THYME-AMR and AMR 3.0 sources. As can be seen from Table 2, this modest infusion of general domain data allowed the parser to attain high performance on the THYME-AMR test set while also largely maintaining its performance on the AMR 3.0 test set. This observation underscores the effectiveness of domain-specific annotation in improving semantic parsing in a joint fashion. This means that the understanding of semantics improves collectively rather than independently, thanks to domain-specific data. As more representative data are collected, we expect further improvements in the parser's performance, making it even more adept at comprehending the semantics in the given domain.

## 4.3 Fine-Grained Performance

Table 3 presents detailed results of our best-performing parser across the semantic components that comprise AMR graphs. AMR representations are an amalgamation of semantic representations including predicate-argument relations, named entities, and coreference components. The SMATCH score represents an average over these sub-categories. To leverage the in-depth analytical power of these linguistic sub categories, a more fine-grained analysis tool[2] provides precision, recall and F1 measures across the various component AMR tasks. We list the fine-grained performance metric category definitions briefly as follows:

- *Unlabeled* category assesses the parsing performance on the AMR graph, disregarding the edge labels.

- *No WSD* category evaluates the parsing performance while ignoring the Propbank word sense labels (e.g., `see-09` becomes just `see`).

- *Concepts* category considers only the abstract concept node matches.

- *Named Entity* category focuses on the matches of named entity subgraphs.

- *Negation* category concerns the matches of the negation attribute nodes(e.g. the `:polarity` edges).

- *Reentrancy* category examines only the concept re-entrancy subgraphs(usually a back reference node).

- *Semantic Role Label (SRL)* category pertains to the performance of each predicate argument structure generation.

We observe that the mixed data augmentation technique significantly improves performance across the board, impacting almost every sub-category of evaluation. Notably, the off-the-shelf parser faced significant challenges in understanding the semantics in the new domain. The performance drop due to domain shifting was not uniform across different sub-categories. The most significant drop in performance was seen in *Named Entity* Recognition, which is expected due to the abundance of medical-related terminology. On the other hand, the data-augmented parser excelled in *Concept* predication and *Named Entity* recognition aspects of AMR parsing, while the performance in the *Negation* and *Reentrancy* category was relatively less impressive compared to the other categories.

## 4.4 Data Requirements for Successful Adaptation

Manual annotation of AMR data is time consuming and expensive. At the current time, the standard AMR 3.0 still consists of only 60k sentences, nearly 10 years after the initial data release. The results shown in Table 2 raise the question of how

---

[2] https://github.com/mdtux89/amr-evaluation

| Sub-category | Training Set | Precision | Recall | F1 |
|---|---|---|---|---|
| SMATCH | THYME-AMR + AMR 3.0 | 0.89 | 0.88 | 0.88 |
| | THYME-AMR | 0.88 | 0.87 | 0.87 |
| | AMR 3.0 | 0.53 | 0.45 | 0.49 |
| Unlabeled | THYME-AMR + AMR 3.0 | 0.90 | 0.90 | 0.90 |
| | THYME-AMR | 0.90 | 0.88 | 0.89 |
| | AMR 3.0 | 0.60 | 0.51 | 0.55 |
| No WSD | THYME-AMR + AMR 3.0 | 0.89 | 0.88 | 0.88 |
| | THYME-AMR | 0.88 | 0.87 | 0.87 |
| | AMR 3.0 | 0.55 | 0.46 | 0.50 |
| Concepts | THYME-AMR + AMR 3.0 | 0.93 | 0.92 | 0.93 |
| | THYME-AMR | 0.93 | 0.91 | 0.92 |
| | AMR 3.0 | 0.52 | 0.46 | 0.49 |
| Named Ent. | THYME-AMR + AMR 3.0 | 0.94 | 0.93 | 0.93 |
| | THYME-AMR | 0.93 | 0.92 | 0.92 |
| | AMR 3.0 | 0.18 | 0.05 | 0.08 |
| Negation | THYME-AMR + AMR 3.0 | 0.86 | 0.85 | 0.85 |
| | THYME-AMR | 0.84 | 0.86 | 0.85 |
| | AMR 3.0 | 0.45 | 0.42 | 0.44 |
| Reentrancies | THYME-AMR + AMR 3.0 | 0.78 | 0.79 | 0.78 |
| | THYME-AMR | 0.78 | 0.76 | 0.77 |
| | AMR 3.0 | 0.48 | 0.37 | 0.41 |
| SRL | THYME-AMR + AMR 3.0 | 0.88 | 0.87 | 0.87 |
| | THYME-AMR | 0.87 | 0.85 | 0.86 |
| | AMR 3.0 | 0.55 | 0.47 | 0.51 |

Table 3: SPRING parser performance analytical breakdowns comparison among three models trained on different combination of the fine-tuning data source. The evaluation is on the THYME-AMR test set.

much data is actually required to attain high levels of parser accuracy through adaptation. To address this question, we conducted a series of experiments training models with progressively larger snapshots of the available training data. Specifically, we gradually augmented the training set size for each model by random sampling without replacement from the training data (resulting in training sets of size 500, 1,000, 2,000, 3,000, 4,000 and 4,955). The results in Figure 4 illustrate the parser's performance across these training sets.

As can be seen, performance rapidly rises from the non-adapted baseline to 80 SMATCH with 1,000 training examples; the model trained on only 2,000 samples achieves 90% of the performance of our best parser trained on all available training data. This rapid improvement with domain specific data is a positive indication of the effectiveness of continued training from a generic model and its ability to rapidly generalize from the domain-specific data.



Figure 4: The performance curve with different sample sizes of the THYME-AMR training set. The x axis is the sample size of the training data; the y axis represents the SMATCH F1 performance score(with the unit of percentage) of the parsers evaluated on the same withheld test set (THYME-AMR test set)

## 5   Discussion

Our results have highlighted the advantages of employing data augmentation techniques for domain adaptation fine-tuning. This opens up the possibility for additional follow-up studies, including the incorporation of data from domain-specific Propbank roleset development. For instance, in the case of THYME, leveraging example sentences for newly added named-entity types like "anatomical-site" could prove beneficial. Initializing the word embedding vectors with such domain-specific concepts would enable a better fit with the pre-trained foundational models. Future investigations involving more sophisticated foundational models and data augmentation approaches hold great promise for enhancing AMR parsing in the medical domain and other specialized domains. By harnessing the capabilities of cutting-edge language models and innovative data augmentation strategies, we can expect significant advancements in semantic parsing tasks and domain adaptation techniques.

With these advances, AMR parses have wide applicability to core information extraction tasks from the clinical narrative such as entity recognition, negation detection, uncertainty detection, coreference, temporality and relation extraction.

## 6   Conclusion

In our investigation, we have presented substantial evidence highlighting the critical role of domain-specific AMR annotations in the context of domain adaptation. Our findings illuminate how variances in the distribution between original and target domains can precipitate a marked decline in the performance of AMR parsing. This phenomenon underscores the challenge of catastrophic forgetting, a significant hurdle in the training of neural network models where new learning can disrupt previously acquired knowledge.

To counteract this issue, we demonstrated the critical role of data augmentation techniques. Specifically, by integrating domain-specific examples into the training dataset, we significantly bolstered the model's capability to acclimate to the nuances of the new domain while preserving its proficiency in the original domain. This strategic approach of coupling domain-specific annotation with thoughtful data augmentation has emerged as a formidable solution, ensuring both the robustness and accuracy of AMR parsing across different domain adaptation scenarios.

Our study reaffirms the indispensability of domain-specific annotation in achieving effective domain adaptation and also supports data augmentation as an essential tool in maintaining a delicate balance between learning new domain characteristics and retaining essential knowledge from the original domain. This balanced approach provides a promising avenue for future research and development in the field of AMR parsing, potentially paving the way for more nuanced and adaptable AI systems capable of navigating other domains with limited data yet maintain robustness.

## 7   Limitations and Future Work

Our study faced constraints primarily due to computational limitations, which necessitated a focus on a specific subset of model and data augmentation strategies. A reasonable extension of this research could involve the exploration of more advanced foundational models, including GPT-3.5, GPT-4, and their publicly accessible counterparts such as LLAMA. These platforms present opportunities for experimenting with zero- or few-shot learning techniques. Importantly, our use of clinical data mandates adherence to stringent privacy standards; thus, it is imperative that any models employed can be locally installed and operated within a secure, firewall-protected environment. This requirement currently excludes the use of proprietary models like those within the GPT family, which are tailored for commercial applications and do not meet the privacy criteria essential for our research objectives.

## 8   Acknowledgements

Danielle Bitterman, Piet de Groen, and Dmitriy Dligach.

## Ethics Statement

In our exploration of clinical notes analysis and the design of automation systems, we navigate through a terrain rich with sensitive personal data and entwined with ethical complexities. Our work is fundamentally rooted in a profound respect for the dignity, rights, and welfare of the individuals whose lives and experiences are documented in these notes. Guided by a set of core ethical principles, our research endeavors to uphold the highest standards of integrity and respect.

Foremost, we prioritize the privacy and confidentiality of patient data. In this paper, all examples have been rigorously de-identified to ensure no personal information can lead back to individuals. Moreover, recognizing the critical importance of obtaining informed consent, we actively collaborate with institutional review boards (IRB) to ethically justify and secure consent approvals for utilizing all data involved in our research.

We are acutely aware of the potential biases in our analysis and interpretation of clinical narratives. This awareness extends to biases that might emerge from the data collection process, the selection of narratives for analysis, and our own preconceptions. We are committed to making concerted efforts to ensure that our analysis encompasses diverse perspectives, thereby avoiding the perpetuation of stereotypes or inequalities.

We urge downstream users of our parser to conscientiously consider the potential impact of their findings on the individuals depicted in the clinical narratives, as well as on wider patient populations. This involves thoughtful reflection on how the research could affect public perceptions, clinical practice, and policy making. A crucial aspect of our approach is to balance the dissemination of research findings with the imperative to prevent harm or distress.

Lastly, our pursuit of transparency in our methodology and findings is relentless. We advocate for the use of Abstract Meaning Representation (AMR) as a superior tool compared to opaque, "black-box" models. AMR offers a fully transparent and verifiable representation of the semantics in clinical narratives, which aligns with our commitment to fostering trust and accountability.

Our approach is a testament to our dedication to ethical research practices, emphasizing the protection of privacy, the mitigation of bias, the thoughtful consideration of impacts, and the advancement of transparency and accountability. These principles are the bedrock of our efforts to contribute meaningful and ethically sound advancements in the field of clinical notes analysis and automation system design.

## References

Xuefeng Bai, Yulong Chen, Linfeng Song, and Yue Zhang. 2021. Semantic representation for dialogue modeling. 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 4430–4445, Online. 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.

Miguel Ballesteros and Yaser Al-Onaizan. 2017. AMR parsing using stack-LSTMs. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1269–1275, Copenhagen, Denmark. Association for Computational Linguistics.

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.

Michele Bevilacqua, Rexhina Blloshmi, and Roberto Navigli. 2021. One SPRING to rule them both: Symmetric AMR semantic parsing and generation without a complex pipeline. In *Proceedings of AAAI*.

Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic Acids Res.*, 32(Database-Issue):267–270.

Claire Bonial, Lucia Donatelli, Mitchell Abrams, Stephanie M. Lukin, Stephen Tratz, Matthew Marge, Ron Artstein, David Traum, and Clare Voss. 2020. Dialogue-AMR: Abstract Meaning Representation for dialogue. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 684–695, Marseille, France. European Language Resources Association.

Claire N. Bonial, Lucia Donatelli, Jessica Ervin, and Clare R. Voss. 2019. Abstract Meaning Representation for human-robot dialogue. In *Proceedings of the Society for Computation in Linguistics (SCiL) 2019*, pages 236–246.

Julia Bonn, Martha Palmer, Zheng Cai, and Kristin Wright-Bettner. 2020. Spatial AMR: Expanded spatial annotation in the context of a grounded Minecraft corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4883–4892, Marseille, France. European Language Resources Association.

Deng Cai and Wai Lam. 2019. Core semantic first: A top-down approach for AMR parsing. 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 3799–3809, Hong Kong, China. Association for Computational Linguistics.

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.

Ramón Fernandez Astudillo, Miguel Ballesteros, Tahira Naseem, Austin Blodgett, and Radu Florian. 2020. Transition-based parsing with stack-transformers. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1001–1007, Online. Association for Computational Linguistics.

Jeffrey Flanigan, Sam Thomson, Jaime Carbonell, Chris Dyer, and Noah A. Smith. 2014. A discriminative graph-based parser for the Abstract Meaning Representation. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1426–1436, Baltimore, Maryland. Association for Computational Linguistics.

William Foland and James H. Martin. 2017. Abstract Meaning Representation parsing using LSTM recurrent neural networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 463–472, Vancouver, Canada. Association for Computational Linguistics.

Ruiliu Fu, Han Wang, Xuejun Zhang, Jun Zhou, and Yonghong Yan. 2021. Decomposing complex questions makes multi-hop QA easier and more interpretable. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 169–180, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Sahil Garg, Aram Galstyan, Ulf Hermjakob, and Daniel Marcu. 2016. Extracting biomolecular interactions using semantic parsing of biomedical text. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, page 2718–2726. AAAI Press.

Leonie Grön, Ann Bertels, and Kris Heylen. 2018. The interplay of form and meaning in complex medical terms: Evidence from a clinical corpus. In *Proceedings of the Joint Workshop on Linguistic Annotation, Multiword Expressions and Constructions (LAW-MWE-CxG-2018)*, pages 18–29, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Thanh Lam Hoang, Gabriele Picco, Yufang Hou, Young-Suk Lee, Lam M. Nguyen, Dzung T. Phan, Vanessa López, and Ramon Fernandez Astudillo. 2021. Ensembling graph predictions for AMR parsing. In *Advances in Neural Information Processing Systems*.

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. SpanBERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8:64–77.

Pavan Kapanipathi, Ibrahim Abdelaziz, Srinivas Ravishankar, Salim Roukos, Alexander Gray, Ramón Fernandez Astudillo, Maria Chang, Cristina Cornelio, Saswati Dana, Achille Fokoue, Dinesh Garg, Alfio Gliozzo, Sairam Gurajada, Hima Karanam, Naweed Khan, Dinesh Khandelwal, Young-Suk Lee, Yunyao Li, Francois Luus, Ndivhuwo Makondo, Nandana Mihindukulasooriya, Tahira Naseem, Sumit Neelam, Lucian Popa, Revanth Gangi Reddy, Ryan Riegel, Gaetano Rossiello, Udit Sharma, G P Shrivatsa Bhargav, and Mo Yu. 2021. Leveraging Abstract Meaning Representation for knowledge base question answering. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3884–3894, Online. Association for Computational Linguistics.

Kate Keahey, Jason Anderson, Zhuo Zhen, Pierre Riteau, Paul Ruth, Dan Stanzione, Mert Cevik, Jacob Colleran, Haryadi S. Gunawi, Cody Hammock, Joe Mambretti, Alexander Barnes, François Halbach, Alex Rocha, and Joe Stubbs. 2020. Lessons learned from the chameleon testbed. In *Proceedings of the 2020 USENIX Annual Technical Conference (USENIX ATC '20)*. USENIX Association.

Kevin Knight, Bianca Badarau, Laura Baranescu, Claire Bonial, Madalina Bardocz, Kira Griffitt, Ulf Hermjakob, Daniel Marcu, Martha Palmer, Tim O'Gorman, and Nathan Schneider. 2020. Abstract meaning representation (amr) annotation release 3.0. Web Download. LDC Catalog No.: LDC2020T02, DOI: https://doi.org/10.35111/44cy-bp51.

Ioannis Konstas, Srinivasan Iyer, Mark Yatskar, Yejin Choi, and Luke Zettlemoyer. 2017. Neural AMR: Sequence-to-sequence models for parsing and generation. In *Proceedings of the 55th Annual Meeting of*

the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 146–157, Vancouver, Canada. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Changmao Li and Jeffrey Flanigan. 2022. Improving neural machine translation with the Abstract Meaning Representation by combining graph and sequence transformers. In *Proceedings of the 2nd Workshop on Deep Learning on Graphs for Natural Language Processing (DLG4NLP 2022)*, pages 12–21, Seattle, Washington. Association for Computational Linguistics.

Zhizhong Li and Derek Hoiem. 2018. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12):2935–2947.

Kexin Liao, Logan Lebanoff, and Fei Liu. 2018. Abstract Meaning Representation for multi-document summarization. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1178–1190, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Fei Liu, Jeffrey Flanigan, Sam Thomson, Norman Sadeh, and Noah A. Smith. 2015. Toward abstractive summarization using semantic representations. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1077–1086, Denver, Colorado. Association for Computational Linguistics.

Chunchuan Lyu and Ivan Titov. 2018. AMR parsing as graph prediction with latent alignment. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 397–407, Melbourne, Australia. Association for Computational Linguistics.

Preslav Nakov. 2013. On the interpretation of noun compounds: Syntax, semantics, and entailment. *Natural Language Engineering*, 19:291 – 330.

Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The Proposition Bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106.

Sudha Rao, Daniel Marcu, Kevin Knight, and Hal Daumé III. 2017. Biomedical event extraction using Abstract Meaning Representation. In *BioNLP 2017*, pages 126–135, Vancouver, Canada,. Association for Computational Linguistics.

Leonardo F. R. Ribeiro, Mengwen Liu, Iryna Gurevych, Markus Dreyer, and Mohit Bansal. 2022. FactGraph: Evaluating factuality in summarization with semantic graph representations. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3238–3253, Seattle, United States. Association for Computational Linguistics.

Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, , and Gerald Tesauro. 2019. Learning to learn without forgetting by maximizing transfer and minimizing interference. In *International Conference on Learning Representations*.

Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. 2022. Fine-tuned language models are continual learners. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6107–6122, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

IV Styler, William F., Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, and James Pustejovsky. 2014. Temporal Annotation in the Clinical Domain. *Transactions of the Association for Computational Linguistics*, 2:143–154.

Chuan Wang, Nianwen Xue, and Sameer Pradhan. 2015. A transition-based algorithm for AMR parsing. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 366–375, Denver, Colorado. Association for Computational Linguistics.

Yu Wang, Hanghang Tong, Ziye Zhu, and Yun Li. 2022. Nested named entity recognition: A survey. *ACM Trans. Knowl. Discov. Data*, 16(6).

Kristin Wright-Bettner, Chen Lin, Timothy Miller, Steven Bethard, Dmitriy Dligach, Martha Palmer, James H. Martin, and Guergana Savova. 2020. Defining and learning refined temporal relations in the clinical narrative. In *Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis*, pages 104–114, Online. Association for Computational Linguistics.

Kristin Wright-Bettner, Martha Palmer, Guergana Savova, Piet de Groen, and Timothy Miller. 2019. Cross-document coreference: An approach to capturing coreference without context. In *Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019)*, pages 1–10, Hong Kong. Association for Computational Linguistics.

Dongqin Xu, Junhui Li, Muhua Zhu, Min Zhang, and Guodong Zhou. 2020. Improving AMR parsing with sequence-to-sequence pre-training. In *Proceedings*

*of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2501–2511, Online. Association for Computational Linguistics.

Sheng Zhang, Xutai Ma, Kevin Duh, and Benjamin Van Durme. 2019. AMR parsing as sequence-to-graph transduction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 80–94, Florence, Italy. Association for Computational Linguistics.

# SERPENT-VLM 🐍 : Self-Refining Radiology Report Generation Using Vision Language Models

**Manav Nitin Kapadnis**** **Sohan Patnaik*** **Abhilash Nandy** **Sourjyadip Ray**
**Pawan Goyal** **Debdoot Sheet**
iammanavk@gmail.com     sohanpatnaik106@gmail.com
Indian Institute of Technology Kharagpur
India

## Abstract

*Radiology Report Generation* (R2Gen) demonstrates how Multi-modal Large Language Models (MLLMs) can automate the creation of accurate and coherent radiological reports. Existing methods often *hallucinate* details in text-based reports that don't accurately reflect the image content. To mitigate this, we introduce a novel strategy, **SERPENT-VLM** (**SE**lf **R**efining Radiology Re**P**ort G**EN**era**T**ion using **V**ision **L**anguage **M**odels), which improves the R2Gen task by integrating a self-refining mechanism into the MLLM framework. We employ a unique *self-supervised loss* that leverages similarity between pooled image representations and the contextual representations of the generated radiological text, alongside the standard Causal Language Modeling objective, to refine image-text representations. This allows the model to scrutinize and align the generated text through dynamic interaction between a given image and the generated text, therefore reducing hallucination and continuously enhancing nuanced report generation. SERPENT-VLM outperforms existing baselines such as LlaVA-Med, BiomedGPT, etc., achieving SoTA performance on the IU X-ray and Radiology Objects in COntext (ROCO) datasets, and also proves to be robust against noisy images. A qualitative case study emphasizes the significant advancements towards more sophisticated MLLM frameworks for R2Gen, opening paths for further research into self-supervised refinement in the medical imaging domain.

## 1 Introduction

*Radiology Report Generation* (R2Gen) serves as a crucial link between medical imaging and natural language processing, to automate the interpretation of radiological images into comprehensive text reports. This task requires models to learn long-range dependencies effectively while generating

the report, a challenge that remains largely unmet in current systems. The primary goal of R2Gen is to generate accurate and comprehensive medical reports from radiological imagery, an essential step toward enhancing diagnostic accuracy and efficiency. Prevailing methods (Vinyals et al., 2015; Xu et al., 2015; Tang et al., 2023; You et al., 2016; Tang et al., 2021) in R2Gen often rely on (1) large datasets for pre-training to impart domain-specific knowledge, and (2) typically utilizing compute-intensive encoder-decoder architectures for fine-tuning. These approaches are fraught with drawbacks, such as omission of minor yet clinically significant details (Wang et al., 2022b; You et al., 2021; Wang et al., 2021) and the persistent issue of *hallucination* as seen in Fig. 1, where generated reports from LlaVA-Med and BiomedGPT wrongly include details not present in the images. Minimizing hallucinations in radiology report generation is crucial since these inaccuracies can lead to misdiagnoses, directly impacting patient treatment plans and outcomes. Moreover, reducing hallucinations ensures the reliability and trustworthiness of automated reports, which is vital for maintaining clinical credibility and facilitating effective patient care. Therefore, the limitations pertaining to existing approaches underscore the necessity for a more refined approach for accurate medical diagnosis, addressing the critical gaps in R2Gen.

In this paper, we introduce a streamlined pipeline, SERPENT-VLM, which begins by processing a given X-ray image by passing it through a visual encoder and mapping it to a vector representation in a high-dimensional space. This process facilitates a nuanced understanding of the medical imagery. The encoded image, alongside a report generation prompt, is then passed as inputs to a Large Language Model (LLM) for text generation. We employ a cross-entropy loss for the **causal language modeling objective** and **introduce a novel self-refining objective** that leverages the pooled

---

| | Ground Truth | LlaVA-Med | BiomedGPT | SERPENT-VLM |
|---|---|---|---|---|
| | The lungs are hyperexpanded . cardiomediastinal silhouette is within normal limits . no pleural effusion focal airspace opacities or pneumothorax . no free subdiaphragmatic air . | clear and well-defined lung field. ribs are well-defined, no obvious fractures or dislocations. The costophrenic angles are sharp, and there are no pleural effusions. However, there appears to be a **small, faint, and irregular shadow in the lower left lung field, which could be indicative of a small nodule or mass.** | Lung fields are clear, with well-defined contours. Ribs show no fractures or dislocations, appearing intact. Costophrenic angles maintain sharpness, indicating absence of pleural effusions. A **slight haziness is noted in the right upper lung field**, potentially indicating **early signs of infiltrative process.** | heart size and mediastinal contours are within normal limits . pulmonary vascularity is within normal limits . no focal consolidation pleural effusion or pneumothorax identified . visualized osseous structures of the thorax are without acute abnormality . |
| | The cardiomediastinal silhouette and pulmonary vasculature are within normal limits in size. The lungs are clear of focal airspace disease, pneumothorax, or pleural effusion. There are no acute bony findings. | Clear, well-defined diaphragm. **Slight irregularity, lower right lung lobe, potential small nodule, early pulmonary issue.** Normal heart shadow, no cardiomegaly. Normal-sized aorta, no pneumothorax, no pleural effusion. Intact ribs, vertebral column, no fractures, dislocations. Centrally located mediastinum, no masses, lymphadenopathy. Normal lung vascular markings, no atelectasis, pneumonia. | clear, well-defined diaphragm and lung fields, except for **faint linearity in the middle right lung,** suggesting **early fibrotic changes.** The heart shadow is normal, excluding cardiomegaly, while the aorta is of normal size with no evidence of pneumothorax or pleural effusion. Bones, including ribs and vertebral column, show no fractures or dislocations, despite **minimal displacement in the right clavicle.** | the heart size and mediastinal contours are within normal limits . the lungs are clear without focal airspace opacity pleural effusion or pneumothorax . there is no visible pneumothorax . there is no evidence of acute bony abnormality |

Figure 1: Generated report samples on IU-Xray dataset. We qualitatively analyze reports generated by medical pre-trained LLMs LlaVA-Med and BioMedGPT with SERPENT-VLM. Hallucinated information in the reports is highlighted using yellow.

image representation and the generated report's contextual representation. This allows for tuning the network without compromising inference latency, while significantly improving performance evaluated using metrics such as $Bleu$, $Rouge_L$, $BertScore$.

The contributions of our work are summarized as follows:

1. Our approach **does not compromise on inference latency**, adopting a refining strategy through a novel loss function used only for fine-tuning

2. The introduction of a self-refining loss ensures the generation of nuanced, **hallucination-free** radiology reports

3. Our system not only matches but surpasses the performance of leading generalistic pre-trained medical LLMs.

4. Our approach demonstrates **robustness against noisy image** inputs, maintaining the generation of comprehensive reports.

This marks a substantial advancement in the field of R2Gen, setting new benchmarks for accuracy, efficiency, and robustness.

The remainder of the paper is organized as follows: We begin by delving into the literature review in Section 2, focusing on current and past state-of-the-art (SoTA) methodologies in the domain of radiological report generation. Section 3

discusses the proposed strategy for the self-refining fine-tuning our approach. The datasets, baselines, experimental setups, and ablation studies are detailed in Section 4. Finally, we conclude with a summary of our findings in Section 5.

## 2 Related Work

**Medical Report Generation (MRG)**: Medical Report Generation has been extensively studied through ML models. (Jing et al., 2018) proposed a co-attention network that aligns visual and textual information to generate comprehensive radiology reports. Further enhancing the capabilities, a memory-driven transformer (Chen et al., 2020) integrates memory modules for encoding and decoding processes, allowing for more sophisticated report generation (Chen et al., 2020, 2021). Cross-modal learning (Wang et al., 2022a) utilizes prototype matrices and contrastive losses to refine the learning of visual-textual correlations, complemented by a self-boosting framework to align image features with report text (Wang et al., 2021). (Liu et al., 2021) addressed the problem of mitigating inherent biases through a data-driven method, introducing a prior-posterior knowledge-based report generation. (Nooralahzadeh et al., 2021) leveraged curriculum learning to extract global concepts to create a bridge between images and text. Task-specific architecture with sentence-level attention mechanism across visual features (Yuan et al., 2019) allows the model to capture key medical concepts from

images. A weakly supervised paradigm to amplify hard negative samples (Yan et al., 2021) addresses the medical data scarcity challenge.

**Large Language Models and Vision language Models**: The advent of Large Language Models (LLMs) such as GPT-4, Claude, BARD showcase excellent zero-shot language understanding (bro, 2020; Li et al., 2021; Liu et al., 2021; Irvin et al., 2019); image understanding and visual question answering (Team et al., 2023) capabilities. Open-source LLMs, like LLaMA and BLOOM, and Multi-modal LLMs such as LlaVA (Liu et al., 2024), Open Flamingo (Awadalla et al., 2023) have also democratized access to cutting-edge generative technology (Ouyang et al., 2022; Pan et al., 2020). Furthermore, domain-specific models LlaVA-Med (Li et al., 2023) and BiomedGPT (Zhang et al., 2024) have shown promising results in pathology and radiology-related tasks. However, knowledge grounding for medical reports (Hyland et al., 2023), thereby reducing hallucination produced by these models remains a challenge.

**Source & Representation of Feedback**: Iterative refinement in MRG has traditionally relied on human feedback to achieve high-quality outputs (Tandon et al., 2022). Scalar reward functions and domain-specific feedback tools, such as compilers, were proposed as cost-effective alternatives to human feedback (Le et al., 2022; Yasunaga and Liang, 2020). Recent developments show that Large Language Models (LLMs) can self-evaluate their responses. However, applying this to Multi-modal Large Language Models remains largely unexplored in terms of generating grounded and hallucination-free responses.

We now discuss the proposed methodology in the subsequent section.

## 3 Methodology

### 3.1 Overview of SERPENT-VLM

We summarize the pipeline of SERPENT-VLM in Figure 2. It consists of two branches to establish the learning optimization criterion. **1) Causal Language Modeling Objective** enforces standard cross-entropy loss (step 4 in Fig. 2) for supervised radiology report generation. Our approach consists of a visual encoder that extracts information from chest X-ray images (step 1 in Fig. 2), a visual mapper that projects low dimensional image features onto high dimensional feature space (step 2 in Fig. 2) and a Large Language Model that au-

toregressively generates the diagnostic radiological report (step 3 in Fig. 2). To further reduce hallucination, we construct a pooled representation of the given X-ray image, a contextual representation leveraging the attention weights and last hidden states of the generated report and enforce **2) Self Refining Objective** that tries to maximise the similarity between pooled image representation and the contextual representation of the generated report through a self-supervised loss criterion (step 5 in Fig. 2). We train the network through a weighted combination of both the losses (step 6 in Fig. 2), thereby enabling SERPENT-VLM to continuously refine itself by aligning generated text with the input image. We now discuss the details of each component.

### 3.2 SERPENT-VLM Framework

The architecture of SERPENT-VLM can be partitioned into three different modules - a visual encoder, a visual mapper and a large language model (LLM). Formally, consider a chest X-ray image $I_v \in \mathbb{R}^{C \times H \times W}$, where $C$ is the number of input channels, $H$, $W$ being the height and width of the image respectively. $I_v = [I_{v_1}, I_{v_2}, \cdots I_{v_k}]$ comprises of a sequence of $k$ patches with $I_{v_i} \in \mathbb{R}^{C \times P \times P}$ being the $i^{th}$ patch, and $P$ is the patch size. We leverage a transformer-based visual encoder $V_{enc}$ to encode and obtain contextual representation $\tilde{e}_{v_i} \in \mathbb{R}^{d_v}$ denoted by Eq. 1 and aggregate each encoded patch to obtain a global image representation $\tilde{e}_v$ depicted by Eq. 2.

$$\tilde{e}_{v_1}, \tilde{e}_{v_2}, \cdots \tilde{e}_{v_k} = V_{enc}(I_{v_1}, I_{v_2}, \cdots I_{v_k}) \quad (1)$$

$$\tilde{e}_v = V_{pooler}(\tilde{e}_{v_1}, \tilde{e}_{v_2}, \cdots \tilde{e}_{v_k}) \quad (2)$$

The encoded image features inherently reside in a visual feature space, which is distinct and not directly compatible with the textual feature space, and hence need to be aligned with the word embedding space of the LLM. To ensure this, we use a learnable visual mapper $V_{map}$ to project the patch embeddings $\tilde{e}_{v_i}$ onto the word embedding space. Formally, $e_{v_i} = V_{map}(\tilde{e}_{v_i})$. We construct a seed prompt $T$ instructing the LLM to generate a report conditioned on the image $I_v$, and obtain the corresponding tokens $\mathcal{T}_{tokens} = [t_1, t_2, \cdots, t_{|\mathcal{T}_{tokens}|}]$ which is given as input to the $Embedding$ module of the LLM to construct the token embeddings (refer Eq. 3),

$$e_{t_1}, e_{t_2}, \cdots, e_{t_{|\mathcal{T}_{tokens}|}} = Embedding(t_1, t_2, \cdots, t_{|\mathcal{T}_{tokens}|})$$
$$(3)$$

Figure 2: Overview of the SERPENT-VLM pipeline. The X-ray image is processed using a visual encoder (step 1) and projected onto a high-dimensional space using a visual mapper (step 2). The encoded image with the report generation prompt is fed into the LLM (step 3). Cross-entropy loss is employed (step 4) for the causal language modeling objective. The pooled image representation and the Contextual representation of the generated report are used to compute the self-refining loss (step 5). A weighted combination of both objectives is used to train the network (step 6).

We concatenate the sequence of projected image patch embeddings $e_{v_i}$ with the seed prompt text embeddings $e_{t_j}$ to obtain a sequence of input embeddings $e_{\mathcal{I}} = [e_v; e_t]$ which are given as input to the decoder-only LLM denoted by $TD$ for generating the logits of the response tokens in autoregressive fashion. $V_{enc}$, $V_{pooler}$, $V_{map}$ and $TD$ are trained through cross-entropy loss $\mathcal{L}_{report}$ enforced between the generated logits and the actual responses. To further guide the report generation process by aligning the generated response with the input image, we enforce a self-supervised *refining loss*.

### 3.3 Self-refining Strategy

We construct an aggregated representation of the generated text by utilizing the attention weights of the last layer of $TD$. Consider the logit distribution for each generated token as $l_i \in \mathbb{R}^d$, where $d$ is the vocabulary size of $TD$. To encode the representation of each generated token, which is further used to compute the *self-refining* loss in a differentiable fashion, we leverage Gumbel-Softmax on the logit distribution to obtain $\hat{l}_i$ for each predicted token. We construct the aggregated representation $\hat{e}_i^p = \sum_{j=1}^d e_j \hat{l}_{ij}$ of each predicted token by

taking a weighted sum of the embedding matrix $E = e_1, e_2, \cdots, e_d$ with $\hat{l}_i$ being the corresponding weights. Formally,

$$\hat{l}_{ij} = \frac{e^{(log(l_{ij})+g_{ij})/\tau}}{\sum_{j=1}^d e^{(log(l_{ij})+g_{ij})/\tau}} \quad (4)$$

Since, the gumbel-softmax operator makes the logit distribution peaky, taking a weighted sum effectively yields the predicted token embeddings. Further, we construct an aggregated representation $h_t \in \mathbb{R}^{d_t}$ of the predicted token embeddings by leveraging the attention weights from the last layer of $TD$. We hypothesize that aligning the aggregated representation of the generated report with the pooled input image representation would reduce hallucination and ground the report generation task. For this, we enforce a *self-refining loss* between $h_t$ and $e_v$ depicted by Eq. 3.3

$$\mathcal{L}_{refine} = \frac{1}{b} \sum_i^b e^{-h_t^T e_v}, \quad (5)$$

where $b$ is the batch size.

Minimizing the negative exponential of the similarity between the image and generated text representation pushes the representation closer, thus further grounding the report generation process. We

286

optimize our network with a weighted combination of both the causal language modeling objective and the self-refining objective. The total loss is denoted by Eq. 6

$$\mathcal{L}_{total} = \lambda_{report} \, \mathcal{L}_{report} + \lambda_{refine} \, \mathcal{L}_{refine} \quad (6)$$

$\mathcal{L}_{report}$ depicts the standard causal language modeling objective that ensures the conditional generation of radiological report text based on the input image, whereas $\mathcal{L}_{refine}$ ensures that the generated report is grounded in context of the input image, thereby establishing a robust pipeline for radiology report generation.

## 4 Experiments and Evaluation

We now discuss the details corresponding to the experiments and ablation studies carried out and enumerate the observations.

### 4.1 Implementation Details

We discuss the technical details and hyper-parameter settings for all the experiments. For the visual encoder $V_{enc}$, we employed the base version of Swin-Transformer-V2[1] and a feed-forward neural network for $V_{map}$. We leverage LLaMA2-7B[2] as our primary LLM. Further, the hidden dimension of $d_v$ of $V_{enc}$ and $d_t$ of $TD$ are 768 and 1024 respectively. We freeze the weights of $V_{enc}$, however keep $V_{map}$ trainable. We employ LoRA with a rank and $\alpha$-scaling factor of 16 each to fine-tune the underlying LLM $TD$. We train SERPENT-VLM for 15 epochs on IU-Xray dataset and 20 epochs on the ROCO dataset with mixed precision on an effective batch size (BS) of 6 using one NVIDIA A40 48GB GPU using a learning rate of $1 \times 10^{-4}$ with linear rate scheduler through AdamW optimizer. For inference, we leverage beam search decoding with beam size configured to 3.

### 4.2 Datasets and Evaluation Metrics:

We evaluate SERPENT-VLM on two commonly used datasets diverse modality -

1. **IU X-Ray** which is a widely used publicly available dataset for medical report generation tasks containing 3,955 fully de-identified radiology reports with sections such as Impression, Findings, Indication, etc., each associated with frontal and/or lateral chest X-rays, totaling 7,470 images;

2. **ROCO** which has 'radiology' and 'out-of-class' subsets (synthetic radiology images, clinical photos, portraits, compound radiology images, and digital art) of roughly 65,460 and 8,182 'radiology', and 4,902 and 613 'out-of-class' images in the train and test set respectively.

Since the reports are verbose and need to be accurately measured with word-level precision, we compute overlap-based metrics like BLEU and Rouge-L, and a semantic similarity-based metric BertScore for evaluating the efficacy of our approach.

| Dataset | Train | Val | Test | Image Views |
|---------|-------|------|------|-------------|
| IU X-Ray | 2769 | 791 | 395 | Frontal and Lateral |
| ROCO | 65460 | 8183 | 8182 | Frontal |

Table 1: Statistics of Evaluation Datasets

### 4.3 Performance of SERPENT-VLM on Radiology Report Generation

Table 2 illustrates the comprehensive comparison of SERPENT-VLM against various state-of-the-art baselines across the IU-Xray and ROCO datasets. In comparison with traditional non-LLM approaches such as Show-Tell (Vinyals et al., 2015), Att2in (Xu et al., 2015), and R2Gen (Chen et al., 2020), SERPENT-VLM exhibits significant improvements. For instance, on the IU-Xray dataset, SERPENT-VLM achieves a $Bleu_4$ score of 0.190, surpassing Show-Tell's 0.078 and R2Gen's 0.165, and even outperforming the more advanced R2GenCMN, which scores 0.170. This indicates not only an improvement in capturing long-range dependencies but also a notable reduction in detail hallucination, a common issue in earlier models. Furthermore, when compared to Medical LLMs and generalistic Vision-Language Models such as LlaVA-Med (Li et al., 2023), BiomedGPT (Zhang et al., 2024), and MiniGPT4 (Zhu et al., 2023), SERPENT-VLM demonstrates superior performance, marking a significant leap in R2Gen. For example, against LlaVA-Med, which records a $Bleu_4$ of 0.186 on IU-Xray, SERPENT-VLM shows a marked improvement with a score of 0.190. Similarly, in the context of $BertScore$, SERPENT-VLM achieves an impressive 0.935 compared to LlaVA-Med's 0.845 and BiomedGPT's 0.793, underscoring its enhanced textual coherence.

| | **IU-Xray** | | | | | | **ROCO** | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Methods** | $Bleu_1$ | $Bleu_2$ | $Bleu_3$ | $Bleu_4$ | $Rouge_L$ | $BertScore$ | $Bleu_1$ | $Bleu_2$ | $Bleu_3$ | $Bleu_4$ | $Rouge_L$ | $BertScore$ |
| Show-Tell | 0.243 | 0.13 | 0.108 | 0.078 | 0.307 | 0.378 | 0.104 | 0.076 | 0.051 | 0.027 | 0.089 | 0.34 |
| Att2in | 0.248 | 0.134 | 0.116 | 0.091 | 0.309 | 0.386 | 0.106 | 0.077 | 0.052 | 0.027 | 0.091 | 0.347 |
| AdaAtt | 0.284 | 0.207 | 0.15 | 0.126 | 0.311 | 0.442 | 0.122 | 0.089 | 0.060 | 0.031 | 0.104 | 0.397 |
| Transformer | 0.372 | 0.251 | 0.147 | 0.136 | 0.317 | 0.579 | 0.159 | 0.116 | 0.079 | 0.041 | 0.137 | 0.521 |
| M2transformer | 0.402 | 0.284 | 0.168 | 0.143 | 0.328 | 0.626 | 0.172 | 0.125 | 0.085 | 0.044 | 0.148 | 0.563 |
| R2Gen | 0.47 | 0.304 | 0.219 | 0.165 | 0.371 | 0.732 | 0.201 | 0.147 | 0.099 | 0.052 | 0.173 | 0.658 |
| R2GenCMN | 0.475 | 0.309 | 0.222 | 0.17 | 0.375 | 0.74 | 0.169 | 0.148 | 0.100 | 0.052 | 0.175 | 0.665 |
| MSAT | 0.481 | 0.316 | 0.226 | 0.171 | 0.372 | 0.749 | 0.212 | 0.150 | 0.102 | 0.053 | 0.177 | 0.673 |
| METransformer | 0.483 | 0.322 | 0.228 | 0.172 | 0.38 | 0.752 | 0.211 | 0.151 | 0.102 | 0.053 | 0.178 | 0.676 |
| R2GenGPT (Deep) | 0.480 | 0.316 | 0.216 | 0.169 | 0.377 | 0.748 | 0.213 | 0.150 | 0.101 | 0.053 | 0.177 | 0.672 |
| MiniGPT4 | 0.494 | 0.329 | 0.220 | 0.179 | 0.390 | 0.767 | 0.219 | 0.156 | 0.103 | 0.056 | 0.183 | 0.689 |
| BiomedGPT | 0.516 | 0.343 | 0.233 | 0.183 | 0.403 | 0.793 | 0.229 | 0.163 | 0.109 | 0.058 | 0.189 | 0.712 |
| LlaVA-Med | 0.528 | 0.346 | 0.237 | 0.186 | 0.422 | 0.845 | 0.234 | 0.164 | **0.111** | **0.061** | 0.198 | 0.759 |
| SERPENT-VLM | **0.547** | **0.356** | **0.242** | **0.190** | **0.452** | **0.935** | **0.243** | **0.169** | 0.108 | 0.057 | **0.212** | **0.84** |

Table 2: Results of SERPENT-VLM on Benchmark datasets

## 4.4 Discussion on the Impact of different Design Choices for SERPENT-VLM

We carry experiments pertaining to two different design choices for SERPENT-VLM and establish the efficacy of the proposed architecture through the comparative analysis across experiments.

1. **Effect of relative importance of two losses:** We vary the relative importance self-refining loss ($\lambda_{refine}$) and report-generation loss ($\lambda_{report}$) in Eq. 6. Table 3 shows that combining the two losses yields much better performance for IU X-ray and ROCO compared to just using the report generation loss (row 5 vs. row 2). This highlights that self-refining loss complements the report generation loss by grounding the generated report on the input image, thereby reducing hallucination. Further, it is observed that using only self-refining loss (row 1) leads to a degradation in performance because SERPENT-VLM is trained only through a self-supervised paradigm without any kind of supervision. As observed, this equilibrium is not merely about avoiding hallucinations but also about fostering a synergistic effect where each loss component reinforces the other, thereby elevating the overall quality and reliability of the automated radiology reports. The findings from our experiments provide compelling evidence for the critical role of balanced loss parameters in achieving the desired outcomes, advocating for a nuanced approach in their application within the framework of SERPENT-VLM.

2. **Effect of contextual representation design strategy:** We explore different aggregation

strategies for obtaining the contextual representation of the generated report. As depicted in Table 4, attention-based aggregation outperforms other aggregation strategies by a significant margin by obtaining a BertScore of 0.935 and 0.840; BLEU$_1$ score of 0.547 and 0.243 on IU X-ray and ROCO respectively. Average pooling (average of token representations), Max pooling (token representation with maximum L2-norm) and Top-k average pooling (average top $k = 5$ token representations based on attention-weights) give suboptimal performance on both IU X-ray and ROCO benchmark, thereby establishing the critical importance of sophisticated feature integration methods in enhancing the model's capability to synthesize coherent and contextually relevant radiology reports. Exploration into different aggregation strategies reveals that the sophistication and adaptability of the aggregation mechanism play a pivotal role in the efficacy of medical report generation models.

## 4.5 How robust is SERPENT-VLM to noisy images?

We assess the robustness of SoTA methods LlaVA-Med and BiomedGPT, with our method SERPENT-VLM, by introducing Gaussian noise to radiological images. Fig. 3 demonstrate that SERPENT-VLM significantly outperforms the current SoTA models, LlaVA-Med and BiomedGPT, across all Gaussian Noise scales, maintaining higher BLEU$_1$ ( 5-6% higher) and BertScore ( 9-10% higher) metrics, thus showcasing superior robustness in report generation under noisy and corrupted images. This also highlights SERPENT-VLM's ability to focus

| Dataset | $\lambda_{Report}$ | $\lambda_{Refine}$ | $Bleu_1$ | $Bleu_2$ | $Bleu_3$ | $Bleu_4$ | $Rouge_L$ | $BertScore$ |
|---|---|---|---|---|---|---|---|---|
| | 0 | 1.0 | 0.416 | 0.270 | 0.184 | 0.144 | 0.344 | 0.711 |
| | **0.3** | **0.7** | **0.547** | **0.356** | **0.242** | **0.190** | **0.452** | **0.935** |
| IU-Xray | 0.5 | 0.5 | 0.492 | 0.320 | 0.218 | 0.171 | 0.407 | 0.842 |
| | 0.7 | 0.3 | 0.479 | 0.311 | 0.212 | 0.166 | 0.396 | 0.818 |
| | 1 | 0.0 | 0.451 | 0.311 | 0.200 | 0.157 | 0.373 | 0.771 |
| | 0 | 1 | 0.187 | 0.130 | 0.083 | 0.044 | 0.163 | 0.647 |
| | **0.3** | **0.7** | **0.243** | **0.169** | **0.108** | **0.057** | **0.212** | **0.840** |
| ROCO | 0.5 | 0.5 | 0.214 | 0.149 | 0.095 | 0.050 | 0.187 | 0.739 |
| | 0.7 | 0.3 | 0.207 | 0.144 | 0.092 | 0.048 | 0.180 | 0.714 |
| | 1 | 0 | 0.194 | 0.135 | 0.086 | 0.046 | 0.170 | 0.672 |

Table 3: Impact of combining self-refining loss (weight $\lambda_{refine}$) with report-generation loss (weight $\lambda_{report}$). Fusing both the loss components gives optimal performance.

| Dataset | Design Strategy | $Bleu_1$ | $Bleu_2$ | $Bleu_3$ | $Bleu_4$ | $Rouge_L$ | $BertScore$ |
|---|---|---|---|---|---|---|---|
| | **Attention based aggregation** | **0.547** | **0.356** | **0.242** | **0.190** | **0.452** | **0.935** |
| IU-Xray | Average pooling | 0.410 | 0.267 | 0.182 | 0.143 | 0.339 | 0.701 |
| | Top k average pooling | 0.465 | 0.303 | 0.206 | 0.162 | 0.384 | 0.795 |
| | Max pooling | 0.383 | 0.249 | 0.169 | 0.133 | 0.316 | 0.655 |
| | **Attention based aggregation** | **0.243** | **0.169** | **0.108** | **0.057** | **0.212** | **0.840** |
| ROCO | Average pooling | 0.190 | 0.132 | 0.084 | 0.044 | 0.165 | 0.655 |
| | Top k average pooling | 0.199 | 0.139 | 0.089 | 0.047 | 0.174 | 0.689 |
| | Max pooling | 0.170 | 0.118 | 0.076 | 0.040 | 0.148 | 0.588 |

Table 4: Performance comparison of different design strategies for contextual representation. Attention weights-based aggregation displays superior performance.

on relevant parts of the image, thereby mitigating the effects of added noise and grounding the generated report - an indication of reduction in hallucination phenomena. The integration of SERPENT-VLM could markedly enhance diagnostic accuracy, aiding radiologists in delivering faster and more accurate patient care.

## 5 Summary and Conclusion

In this paper, we propose SERPENT-VLM, an innovative method for producing detailed and accurate radiology reports from Chest X-rays without hallucinations. The process utilizes a frozen visual encoder to transform X-ray images into a high-dimensional space, which a Large Language Model (LLM) then uses to generate initial reports. These reports undergo further refinement through a novel combination of self-refining loss and Causal Language Modeling Loss, significantly surpassing existing methods as detailed in Section 4. Our experiments in Section 4 and supplementary materials, confirm the effectiveness of our self-refining approach, even with distorted noisy images. Our future works involve the extension of our method to other medical imaging types, such as MRIs and CT

scans, and to incorporate diagnostic RADreports to enhance report accuracy further.

## Limitations

The SERPENT-VLM has shown significant advancements in creating radiology reports from chest X-rays, reducing inaccuracies, and better matching the content of the images compared to earlier models. However, this research has its limitations. The testing of the model's performance and adaptability has been limited to particular datasets (IU X-Ray and ROCO), which do not encompass the broad spectrum of radiological images or health conditions. It remains unclear how well this would work in actual medical situations. Furthermore, although the model's ability to handle low-quality images is emphasized, the wide range of image quality in real-life scenarios could pose challenges that have yet to be evaluated.

## Ethics Statement

The deployment of SERPENT-VLM in clinical settings involves significant ethical considerations. The model's potential to generate erroneous interpretations from radiological images, despite re-

(a) Performance metrics for ROCO dataset with varying levels of Gaussian noise added to input radiological images.



(b) Performance metrics for IU-Xray dataset with varying levels of Gaussian noise added to input radiological images.

Figure 3: Comparative performance metrics for ROCO and IU-Xray datasets.

duced hallucinations, necessitates cautious application, especially since incorrect reports could lead to misdiagnoses or inappropriate treatments. The use of large datasets for training also raises privacy concerns, requiring stringent data handling and patient consent protocols.

# References

2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.

Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. 2021. Cross-modal memory networks for radiology report 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 5904–5914, Online. Association for Computational Linguistics.

Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. 2020. Generating radiology reports via memory-driven transformer. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1439–1449, Online. Association for Computational Linguistics.

Stephanie L Hyland, Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Mercy Ranjit, Anton Schwaighofer, Fernando Pérez-García, Valentina Salvatelli, Shaury Srivastav, Anja Thieme, et al. 2023. Maira-1: A specialised large multimodal model for radiology report generation. *arXiv preprint arXiv:2311.13668*.

Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI'19/IAAI'19/EAAI'19. AAAI Press.

Baoyu Jing, Pengtao Xie, and Eric Xing. 2018. On the automatic generation of medical imaging reports. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics.

Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Hoi. 2022. CodeRL: Mastering code generation through pretrained models and deep reinforcement learning. In *Advances in Neural Information Processing Systems*.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. Llava-med: Training a large language-and-vision assistant for biomedicine in one day.

Yehao Li, Yingwei Pan, Jingwen Chen, Ting Yao, and Tao Mei. 2021. X-modaler: A versatile and high-performance codebase for cross-modal analytics. In *Proceedings of the 29th ACM International Conference on Multimedia*, MM '21, page 3799–3802, New York, NY, USA. Association for Computing Machinery.

Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. 2021. Exploring and distilling posterior and prior knowledge for radiology report generation. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.

Farhad Nooralahzadeh, Nicolas Perez Gonzalez, Thomas Frauenfelder, Koji Fujimoto, and Michael

Krauthammer. 2021. Progressive transformer-based generation of radiology reports. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2824–2832, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.

Yingwei Pan, Ting Yao, Yehao Li, and Tao Mei. 2020. X-linear attention networks for image captioning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.

Niket Tandon, Aman Madaan, Peter Clark, and Yiming Yang. 2022. Learning to repair: Repairing model output errors after deployment using a dynamic memory of feedback. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 339–352, Seattle, United States. Association for Computational Linguistics.

Mingkang Tang, Zhanyu Wang, Zhenhua LIU, Fengyun Rao, Dian Li, and Xiu Li. 2021. Clip4caption: Clip for video caption. In *Proceedings of the 29th ACM International Conference on Multimedia*, MM '21, page 4858–4862, New York, NY, USA. Association for Computing Machinery.

Mingkang Tang, Zhanyu Wang, Zhaoyang Zeng, Xiu Li, and Luping Zhou. 2023. Stay in grid: Improving video captioning via fully grid-level representation. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(7):3319–3332.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. 2015. Show and tell: A neural image caption generator. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3156–3164, Los Alamitos, CA, USA. IEEE Computer Society.

Jun Wang, Abhir Bhalerao, and Yulan He. 2022a. Cross-modal prototype driven network for radiology report generation. In *Computer Vision – ECCV 2022*, pages 563–579, Cham. Springer Nature Switzerland.

Zhanyu Wang, Hongwei Han, Lei Wang, Xiu Li, and Luping Zhou. 2022b. Automated radiographic report generation purely on transformer: A multicriteria supervised approach. *IEEE Transactions on Medical Imaging*, 41(10):2803–2813.

Zhanyu Wang, Luping Zhou, Lei Wang, and Xiu Li. 2021. A self-boosting framework for automated radiographic report generation. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2433–2442.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2048–2057, Lille, France. PMLR.

An Yan, Zexue He, Xing Lu, Jiang Du, Eric Chang, Amilcare Gentili, Julian McAuley, and Chun-Nan Hsu. 2021. Weakly supervised contrastive learning for chest X-ray report generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4009–4015, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Michihiro Yasunaga and Percy Liang. 2020. Graph-based, self-supervised program repair from diagnostic feedback. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org.

Di You, Fenglin Liu, Shen Ge, Xiaoxia Xie, Jing Zhang, and Xian Wu. 2021. AlignTransformer: Hierarchical alignment of visual regions and disease tags for medical report generation. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2021*, Lecture notes in computer science, pages 72–82. Springer International Publishing, Cham.

Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4651–4659.

Jianbo Yuan, Haofu Liao, Rui Luo, and Jiebo Luo. 2019. Automatic radiology report generation based on multi-view image fusion and medical concept enrichment. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, pages 721–729, Cham. Springer International Publishing.

Kai Zhang, Jun Yu, Eashan Adhikarla, Rong Zhou, Zhiling Yan, Yixin Liu, Zhengliang Liu, Lifang He, Brian Davison, Xiang Li, Hui Ren, Sunyang Fu, James Zou, Wei Liu, Jing Huang, Chen Chen, Yuyin Zhou, Tianming Liu, Xun Chen, Yong Chen, Quanzheng Li, Hongfang Liu, and Lichao Sun. 2024. Biomedgpt: A unified and generalist biomedical generative pre-trained transformer for vision, language, and multimodal tasks.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models.

# ERD: A Framework for Improving LLM Reasoning
# for Cognitive Distortion Classification

**Sehee Lim**[*]
Yonsei University
sehee0706@yonsei.ac.kr

**Yejin Kim**[*]
Yonsei University
yjkim.stat@yonsei.ac.kr

**Chi-Hyun Choi**[*]
EverEx
leo@everex.co.kr

**Jy-yong Sohn**[†]
Yonsei University
EverEx
jysohn1108@yonsei.ac.kr

**Byung-Hoon Kim**[†]
Yonsei University
EverEx
egyptdj@yonsei.ac.kr

## Abstract

Improving the accessibility of psychotherapy with the aid of Large Language Models (LLMs) is garnering a significant attention in recent years. Recognizing cognitive distortions from the interviewee's utterances can be an essential part of psychotherapy, especially for cognitive behavioral therapy. In this paper, we propose ERD, which improves LLM-based cognitive distortion classification performance with the aid of additional modules of (1) extracting the parts related to cognitive distortion, and (2) debating the reasoning steps by multiple agents. Our experimental results on a public dataset show that ERD improves the multi-class F1 score as well as binary specificity score. Regarding the latter score, it turns out that our method is effective in debiasing the baseline method which has high false positive rate, especially when the summary of multi-agent debate is provided to LLMs.

## 1 Introduction

Large Language Models (LLMs) are dominating the research areas in machine learning and artificial intelligence, broadening its usage in various applications (Radford et al., 2018, 2019; Brown et al., 2020; OpenAI, 2023; Ouyang et al., 2022). Especially in the medical domain, PaLM (Chowdhery et al., 2022) and its variants, such as Med-PaLM (Singhal et al., 2022), are equipped with medical data and instructions to answer the questions from clinical field (Chowdhery et al., 2022; Singhal et al., 2023). In addition, conversational AI assistant chatbots are devised to support patients with mental health issues (Rathje et al., 2023; Vaidyam et al., 2019; Saha et al., 2022; Stock et al., 2023; Liu et al., 2023; Welivita et al., 2021; Sharma et al., 2020).

Recognizing the fact that individuals with mental disorders hesitate to seek in-person medical con-

sultations (Steinberg et al., 1980), previous studies (Yang et al., 2023; Lee et al., 2023; Chen et al., 2023b) attempt to enhance the accessibility and quality of psychotherapy through the use of LLMs with Chain-of-Thought (CoT) reasoning (Wei et al., 2022). These models aim to detect the user's personality and interpret their mental state in order to generate more empathetic responses.

For example, Diagnosis-of-Thought (DoT) uses LLMs to classify cognitive distortions from utterances, which is a crucial part of Cognitive Behavior Therapy (CBT) (Chen et al., 2023b).

While the DoT method holds promise, one key challenge that remains an open issue is the tendency of the model to overdiagnose cognitive distortions, incorrectly inferring irrational thought patterns even when the user's statements are benign. In addition, the distortion classification performance of DoT in multi-class setup is close to that of random guessing, which limits its usage in practice.

In this paper, we tackle these issues by proposing a new framework for classifying cognitive distortions from the user utterances, by introducing modules for debiasing the overdiagnosing tendency of existing methods and for improving the performance on classifying distortion types inferred from the utterances.

Our main contributions can be summarized as below:

- We introduce ERD, a new framework for classifying cognitive distortions in the user utterances using three steps: Extraction, Reasoning, and Debate, each of which uses LLMs. The first step lets LLM extract a part of the utterances that is related with the distortion, the second step uses LLM to generate the thought process of estimating cognitive distortions from the extracted part, and the third step uses multi-agent LLMs to discuss the thought process described in the second

---

[*]Equal contribution
[†]Corresponding authors

Figure 1: The pipeline of Extraction-Reasoning-Debate (ERD), which detects and classify the cognitive distortion from the input user speech. It begins with the identification and extraction of potential cognitive distortions from the user speech. These extracted elements are then utilized to construct an intermediate reasoning step. Subsequently, a debate is conducted, wherein multiple LLM agents deliberate to assess the presence and type of cognitive distortion. Finally, a judge integrates the entire debate process to get the final answer on the distortion classification problem.

step and make the final decision.

- Compared with existing baselines, ERD improves the multi-class F1 score for distortion classification task by more than 9% and improves the distortion assessment specificity score by more than 25%, when tested on the cognitive distortion detection dataset with 2530 samples in Kaggle.

- We provide factor analysis on ERD, showing that (1) multiple rounds of debate in ERD is beneficial for improving the classification score, and (2) the summarization and the validity evaluation processes during the debate step enhance the debiasing effect.

## 2   ERD

We propose Extraction-Reasoning-Debate (ERD), a framework for classifying distortions in a given user speech, as shown in Fig. 1. The prompts we used can be found in Figure 3 in Appendix. Below we elaborate each step in our framework.

| Input | Distortion Classification |
|---|---|
| User Speech | $15.28_{0.65}$ |
| Distorted Part of User Speech | $\mathbf{27.08_{0.27}}$ |

Table 1:  Multi-class F1 score of DoT (Chen et al., 2023b) for the cognitive distortion classification problem, when two different inputs are given. The first option uses the user speech as the input, as done in (Chen et al., 2023b). The second option is considered by us, which only puts the ground-truth distorted part within the user speech. Putting only the distorted part significantly improves the classification performance, which motivates the Extraction step in ERD framework.

### 2.1   Extraction

To provide the motivation for the Extraction step proposed in our method, we first share our empirical results showing that extracting the distorted parts of user speech is beneficial for distortion classification. Table 1 shows the multi-class F1 score of Diagnosis-of-Thought (DoT) method for distortion classification problem (predicting out of 10 classes), tested on a cognitive distortion detection dataset with 2530 samples in Kaggle[1]. We test on two different options: (1) putting the user speech as it is, and (2) putting the ground-truth part (provided in the 'distorted part' column of the dataset) within the speech, that indicates the distortion. Table 1 shows that the multi-class F1 score increases more than 10% when the ground-truth distorted part is extracted before running DoT.

Motivated by this result, prior to the Reasoning step (e.g., DoT) which outputs the thought process for assessing/classifying the distortion, we add an Extraction step which instructs LLMs isolate the segments from the user's utterance that may potentially exhibit cognitive distortions. This process of extraction is done *without* paraphrasing or summarizing, thereby preserving the original context and nuances for the subsequent thought process. In summary, Extraction process ensures that the LLMs' responses hinge on the most informative facets of the utterance, which in turn enhance the quality of the distortion classification performance.

### 2.2   Reasoning

Our target task (cognitive distortion classification from the user speech) is naturally considered as a

---

[1] https://www.kaggle.com/
datasets/sagarikashreevastava/
cognitive-distortion-detetction-dataset

| Method | Distortion Assessment (True/False) | | | Distortion Classification (out of 10 types) |
| | Sensitivity | Specificity | F1 Score | Weighted F1 Score |
| --- | --- | --- | --- | --- |
| Reasoning | $99.29_{0.19}$ | $6.79_{0.34}$ | $\mathbf{78.26_{0.16}}$ | $15.28_{0.65}$ |
| +Extraction | $\mathbf{99.83_{0.03}}$ | $0.93_{0.22}$ | $\underline{77.48_{0.04}}$ | $\mathbf{24.40_{0.69}}$ |
| +Debate | $73.10_{0.26}$ | $\mathbf{33.05_{0.58}}$ | $68.89_{0.24}$ | $22.18_{0.99}$ |
| +Extraction+Debate | $74.89_{2.31}$ | $\underline{30.74_{3.92}}$ | $69.49_{0.62}$ | $\underline{24.27_{1.14}}$ |

Table 2: Cognitive distortion assessment/classification results of ERD when various modules (Extraction and Debate) are added. Here, we test on cognitive distortion detection dataset in Kaggle, and use DoT (Chen et al., 2023b) method for the Reasoning step. Upon the above results, Extraction improves the distortion classification performance and Debate increases the distortion assessment specificity significantly. Combining both Extraction and Debate takes the sweet spot, simultaneously enhancing both performances.



Figure 2: Confusion matrices of ERD when tested on 2530 samples: (Left) only Reasoning is used, (Right) Extraction, Reasoning and Debate steps are used. Including Extraction and Debate modules increases the number of true negatives from 61 to 322, thus correctly identifying the samples with 'no distortion'.

task that requires logical thinking, if we imagine how doctors classify the patients. In recent years, various methods propose letting LLMs mimic the logical thought process or reasoning steps. For example, chain-of-thought (CoT) prompting and its variants (Wei et al., 2022; Kojima et al., 2022; Yao et al., 2023; Besta et al., 2023; Chen et al., 2023a; Yang et al., 2023; Lee et al., 2023) provide a significant performance improvement in various reasoning tasks including common sense reasoning and mathematical reasoning.

Our Reasoning step chooses any existing methods which let LLMs output the thought process for performing the target task. By default, we use diagnosis-of-thought (DoT) (Chen et al., 2023b) comprised of three critical stages (subjectivity assessment, contrastive reasoning, and schema analysis) that construct rationales for the detection of cognitive distortions. At the *subjectivity assessment* stage, the input utterances are differentiated between the objective facts and the subjective thoughts. This is followed by the *contrastive reasoning* stage, where the process elicits both supportive and contradictory perspectives to the speaker's viewpoint. The final stage, *schema analysis*, involves delving into the underlying thought schema, which refers to the subconscious cognitive patterns or frameworks that shape and influence a person's specific thought process and behavior.

## 2.3 Debate

Several recent works on using LLMs for reasoning tasks show that multiple LLM agents debating their thought processes significantly improve the performance (Liang et al., 2023; Zheng et al., 2023; Xiong et al., 2023; Chan et al., 2023; Du et al., 2023). Motivated by this observation, we add multi-agent debate (or Debate) step following the Reasoning step. In Figure 1, ERD employs three LLM agents, each designated with the role of "physician" to simulate a professional medical debate. The discussion between first two agents (two debators) is overseen by the third agent (called the judge agent), bearing the role of "head doctor" who monitors the entire debate to ensure a fair evaluation. The third agent is introduced, motivated by recent result showing that LLMs can behave as a good judge (Zheng et al., 2023). The first debater presents arguments for the presence or absence of cognitive distortion in the user speech, based on the LLM outputs obtained in the Extraction and Reasoning steps. Subsequently, the second debator counters the initial assertions, presenting a contradicting viewpoint. The first debater then responds to this counterargument, followed by a second round of rebuttal from the second debater, resulting in two rounds of argumentation. One can consider repeating this iterative exchange of thoughts for multiple rounds. After this iterative process, the judge agent integrates the entire discourse, employing two proposed methodologies to reach a final decision.

We consider two different options for controlling the behavior of the judge agent to get better performances. The first option involves a straightforward summarization of the total debate process. The second option involves summarizing the debate and evaluating which side's arguments are more valid. By adding such summarization process, we expect that the final answer of ERD is based on a comprehensive consideration of all presented viewpoints.

| Distortion Type | Count |
|---|---|
| All-or-nothing thinking | 100 |
| Emotional Reasoning | 134 |
| Fortune-telling | 143 |
| Labeling | 165 |
| Magnification | 195 |
| Mental filter | 122 |
| Mind Reading | 239 |
| Overgeneralization | 239 |
| Personalization | 153 |
| Should statements | 107 |
| No Distortion | 933 |
| In Total | 2530 |

Table 3: Details of dataset used in this paper; we report the number of samples for each class including 10 types of distortions and "no distortion" type whose utterance does not contain distortion.

## 3 Experiments

**Settings** We use a cognitive distortion detection dataset (Shreevastava and Foltz, 2021) composed of speeches that correspond to 10 types of "cognitive distortions" and neutral speeches categorized as "no distortion" type. This dataset, sourced from Kaggle[1], contains 2530 annotated examples by experts and the least number of examples for each type of distortion is 100, which can be found in Table 3. This dataset is designed to facilitate two tasks: distortion assessment and distortion classification.

In the distortion assessment task, the model determines whether cognitive distortion is present in the patient's utterance. In the distortion classification task, the model identifies the specific type of cognitive distortion. We report the Sensitivity, Specificity and F1 score for the distortion assessment task, and the weighted F1 score for the distortion classification task. We run 3 random trials and report the mean and standard deviation values. We employ the model gpt-3.5-turbo with the temperature as 0.1. Every result reported in this paper is based on the zero-shot prompting.

**Experimental results** Table 2 shows the performances of ERD, when different modules are plugged in. Compared with naive method using Reasoning module only, adding Extraction module improves the distortion classification score by more than 9%, and adding Debate module not only improves the distortion classification score by around 7%, but also improves the distortion assessment specificity by more than 25%. Fig. 2 shows the confusion matrix of the ERD for two cases: (1) when only the Reasoning module is used, and (2) when Extraction, Reasoning and

Debate are used. This result shows that adding Extraction and Debate modules promotes the correct estimation of utterances with no distortion. This qualitative result can be supported by our qualitative results (in Fig. 4 and Fig.5 in Appendix) showing the effectiveness of Debate step for improving the estimation performance. For a given speech (that does not have cognitive distortion), Fig. 4 and Fig. 5 show the responses of LLM, when Debate step is in-activated and activated, respectively. While LLM without Debate incorrectly estimates that the speech contains cognitive distortion (of type "Labeling"), LLM with Debate correctly estimates that the speech does not contain cognitive distortions.

Recall that in Debate step of ERD, we consider different prompting techniques to control the behavior of the judge agent when making the final decision. Table 4 shows the effect of such prompts for three variants:

(1) "ERD without summarization" does not instruct judge to summarize the claims of debate and just directly make decision, (2) "ERD with summarization" instructs judge to summarize the claims before making the decision, and "ERD with summarization and validity evaluation" instructs judge to summarize and evaluate the claims of debate before making the decision. Note that the specificity is keep improved as we provide more detailed instructions to the judge agent.

Table 5 shows how the performance improves as we increase $r$, the number of Debate rounds used in ERD. The results show that increasing the number of Debate rounds led to enhancements in both the binary F1 score and the multi-class F1 score. The performance saturates after $r = 2$, thus better to use two rounds of debate considering the token efficiency. This finding aligns with the results presented in a related work on multi-agent debate of LLMs, demonstrating a similar pattern in the impact of the number of debate rounds on the model performance (Du et al., 2023).

## 4 Conclusion

We introduce ERD, a framework using LLMs to estimate the cognitive distortion contained in the user utterances through three steps: Extracting distorted parts within the utterances, Reasoning the estimation of the corresponding distortion classes, and Debating the initial estimation using multiple agents. Compared with existing baselines only having the reasoning step, including the extraction

| | Distortion Assessment | | | Distortion Classification |
|---|---|---|---|---|
| | Sensitivity | Specificity | F1 Score | Weighted F1 Score |
| ERD without Summarization | $\mathbf{92.13}_{0.38}$ | $11.01_{0.66}$ | $\mathbf{75.48}_{0.23}$ | $\mathbf{25.28}_{0.46}$ |
| +Summarization | $\underline{86.10}_{0.58}$ | $\underline{19.58}_{1.36}$ | $\underline{73.88}_{0.36}$ | $23.96_{1.05}$ |
| +Summarization+Validity Evaluation | $74.89_{2.31}$ | $\mathbf{30.74}_{3.92}$ | $69.49_{0.62}$ | $\underline{24.27}_{1.14}$ |

Table 4: Comparison of ERD with three different prompting options that control the behavior of the judge. For all three options, the `Extraction` and `Reasoning` modules are active in all cases, with differences applied exclusively to the `Debate` module. For the first option, judge predicts the cognitive distortion type only based on the debate process log, without any summarization step. For the second option, judge first *summarizes* the debate and then predicts the cognitive distortion type. In the final option, judge *summarizes* the debate, *evaluates the validity* of the claims in the debate, and then predicts the cognitive distortion type. Both summarization and validity evaluation steps improve the performance in terms of specificity. Note that the number of `Debate` rounds is set to $r = 2$.

| Metric | Round 1 | Round 2 | Round 3 |
|---|---|---|---|
| Binary F1 | $52.13_{1.25}$ | $69.49_{0.62}$ | $\mathbf{70.74}_{0.44}$ |
| Multi-class F1 | $22.79_{1.62}$ | $24.27_{1.14}$ | $\mathbf{24.83}_{0.81}$ |

Table 5: F1 scores for different $r$, the number of `Debate` rounds. The performances improve as $r$ increases.

and `debating` steps improve the distortion classification performance by 9% and improve the distortion assessment specificity by over 25%. Such improvements is crucial to cognitive behavior therapy since ERD is more adept at correctly identifying cases without distortions, avoiding the pitfall of over-diagnosing cognitive distortions. Furthermore, experimental results reveal that we can control the behavior of ERD with various prompting options.

# References

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023a. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks.

Zhiyu Chen, Yujie Lu, and William Yang Wang. 2023b. Empowering psychotherapy with large language models: Cognitive distortion detection through diagnosis of thought prompting. *arXiv preprint arXiv:2310.07146*.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate.

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.

Yoon Kyung Lee, Inju Lee, Minjung Shin, Seoyeon Bae, and Sowon Hahn. 2023. Chain of empathy: Enhancing empathetic response of large language models based on psychotherapy models. *arXiv preprint arXiv:2311.04915*.

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. 2023. Encouraging divergent thinking in large language models through multi-agent debate.

June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi Liao, and Jiamin Wu. 2023. Chatcounselor: A large language models for mental health support. *arXiv preprint arXiv:2309.15461*.

R OpenAI. 2023. Gpt-4 technical report. *arXiv*, pages 2303–08774.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Steve Rathje, Dan-Mircea Mirea, Ilia Sucholutsky, Raja Marjieh, Claire Robertson, and Jay J Van Bavel. 2023. Gpt is an effective tool for multilingual psychological text analysis.

Tulika Saha, Saichethan Reddy, Anindya Das, Sriparna Saha, and Pushpak Bhattacharyya. 2022. A shoulder to cry on: Towards a motivational virtual assistant for assuaging mental agony. In *Proceedings of the 2022 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 2436–2449.

Ashish Sharma, Adam S. Miner, David C. Atkins, and Tim Althoff. 2020. A computational approach to understanding empathy expressed in text-based mental health support.

Sagarika Shreevastava and Peter Foltz. 2021. Detecting cognitive distortions from patient-therapist interactions. In *Proceedings of the Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access*, pages 151–158.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2022. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*.

Herbert Steinberg, Moshe Torem, and Stephen M Saravay. 1980. An analysis of physician resistance to psychiatric consultations. *Archives of General Psychiatry*, 37(9):1007–1012.

Anna Stock, Stephan Schlögl, and Aleksander Groth. 2023. Tell me, what are you most afraid of? exploring the effects of agent representation on information disclosure in human-chatbot interaction. In *International Conference on Human-Computer Interaction*, pages 179–191. Springer.

Aditya Nrusimha Vaidyam, Hannah Wisniewski, John David Halamka, Matcheri S Kashavan, and John Blake Torous. 2019. Chatbots and conversational agents in mental health: a review of the psychiatric landscape. *The Canadian Journal of Psychiatry*, 64(7):456–464.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.

Anuradha Welivita, Yubo Xie, and Pearl Pu. 2021. A large-scale dataset for empathetic response generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1251–1264.

Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. 2023. Examining inter-consistency of large language models collaboration: An in-depth analysis via debate.

Tao Yang, Tianyuan Shi, Fanqi Wan, Xiaojun Quan, Qifan Wang, Bingzhe Wu, and Jiaxiang Wu. 2023. Psycot: Psychological questionnaire as powerful chain-of-thought for personality detection. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

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. *arXiv preprint arXiv:2305.10601*.

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.

You are a physician participating in a debate discussing the presence or absence of cognitive distortions in a given speech.
Your task is to 1) finish a few test of thoughts questions to analyze the thought patterns of the patient.
Then based on the test of thoughts analysis, 2) make an initial claim and 3) defend your position against another debater's refutation.

**Debater 2's Instruction**

You are a physician participating in a debate discussing the presence or absence of cognitive distortions in a given speech. Your task is to refute another debater's claim.
The other debater analyzed the person's speech and determined whether cognitive distortions were present, based on the analysis. They also identified specific types of cognitive distortions, if any were found.
Now it's your turn to refute the previous claim.

**Judge's Instruction**

You are a head doctor and your task is to monitor the entire debate and make a final decision.
In the debate, the two physicians discussed whether there was a cognitive distortion present and, if so, what type it was.

**Debater 1) Reasoning**

Given a speech of a person: "__{person's speech}__", finish the following diagnosis of thought questions:
0. Extract the sentence that most clearly shows distortion from the patient's utterance.
1. Based on the extracted sentence, answer the following question: what is the situation? Find out the facts that are objective; What is the person thinking or imagining? Find out the thoughts or opinions that are subjective.
2. What makes the person think the thought is true or is not true? Find out the reasoning processes that support and do not support these thoughts.
3. Why does the person come up with such reasoning process supporting the thought? What's the underlying cognitive mode of it?
Each stage should be constructed based on the findings of the preceding one.

**Debater 1) Initial Claim**

Then, based on the above analysis, Respond to the task.
1. Identify if there is cognitive distortion in the speech.
2. Recognize the specific types of the cognitive distortion.

**Debater 2) Refutation**

The opposing team's argument: "__{debater1's initial claim}__"
Refute the opposing team's arguments and present your new claim.

**Debater 1) Defense**

This is the opposing team's refutation of your opinion: "__{debater2's refutation}__"
Reaffirm and defend your original argument.

**Debater 2) Second Refutation**

This is the opposing physician's defense against your refute: "__{debater1's defense}__"
Counter this defense.

**Judge) Summary and Reasoning**

The discussion on whether there is cognitive distortion in this speech is as follows : "__{entire debate process}__"
Summarize the entire debate and evaluate which side presents a more valid argument.

**Judge) Final Decision**

Make a final decision.

Figure 3: Prompts of ERD proposed in this paper.

298

**Speech**

My mother and I have had a rocky history, mostly as a result of a childhood divorce and her alcohol addiction. My mother cheated on my father, who I always been extremely close to, and re-married immediately to a younger man. Subsequently over the years, she received treatment and changed her life, and so we re-established a relationship and have grown close.

**Reasoning and Answering**

1) There is a presence of cognitive distortion in the speech.
2) The specific cognitive distortion is "Labeling," as the patient reduces their mother to a single characteristic of betrayal and abandonment based on her actions. This oversimplification of the mother's behavior contributes to a distorted view of their relationship dynamics.

Figure 4: The output of ERD for a given speech, when the Debate step is not activated. LLM estimates that the speech contains cognitive distortion (type "Labeling") even though it does not.

**Speech**

My mother and I have had a rocky history, mostly as a result of a childhood divorce and her alcohol addiction. My mother cheated on my father, who I always been extremely close to, and re-married immediately to a younger man. Subsequently over the years, she received treatment and changed her life, and so we re-established a relationship and have grown close.

**Debater1's Initial Claim**

1) There is a presence of cognitive distortion in the speech.
2) The specific cognitive distortion is "Labeling," as the patient reduces their mother to a single characteristic of betrayal and abandonment based on her actions. This oversimplification of the mother's behavior contributes to a distorted view of their relationship dynamics.

**Debater2's Refutation**

I respectfully disagree with the opposing team's assessment. The patient's speech does not exhibit the cognitive distortion of labeling. The patient accurately describes specific actions of betrayal and abandonment by the mother, which are factual events, not an oversimplification or reduction of the mother to a single characteristic. The patient acknowledges the complexity of the relationship dynamics by also mentioning the mother's alcohol addiction and subsequent treatment, indicating a nuanced understanding of the situation. Therefore, there is no evidence of cognitive distortion in the form of labeling in the speech.

**Debater1's Defense**

I maintain that the patient's speech demonstrates the cognitive distortion of labeling. By focusing solely on the mother's actions of betrayal and abandonment, the patient reduces her to a single characteristic, overlooking the complexity of human behavior. While the events described may be factual, the patient's emphasis on these aspects distorts the overall perception of the relationship. The patient's strong emotional response to these specific actions indicates a biased interpretation, aligning with the cognitive distortion of labeling. Therefore, the presence of labeling in the speech is evident despite the acknowledgment of other factors.

**Debater2's Second Refutation**

I respectfully maintain my position that the patient's speech does not exhibit the cognitive distortion of labeling. The patient's detailed account of specific events does not equate to reducing the mother to a single characteristic. Acknowledging the complexity of the relationship dynamics by mentioning the mother's alcohol addiction and treatment shows a nuanced understanding, not oversimplification. The patient's emotional response to betrayal and abandonment does not automatically indicate labeling; it reflects a genuine reaction to significant events. Therefore, the presence of labeling in the speech remains unsubstantiated.

**Judge's Summary and Reasoning**

The debate revolves around whether the patient's speech exhibits the cognitive distortion of labeling. Debater 1 argues that the patient reduces the mother to a single characteristic of betrayal and abandonment, indicating labeling. On the other hand, Debater 2 contends that the patient's detailed account of events does not equate to labeling, as the patient acknowledges the complexity of the situation. While both sides make valid points, Debater 2's argument that the patient's nuanced understanding and emotional response do not necessarily indicate labeling appears stronger, as it considers the broader context and complexity of human behavior.

**Judge's Final Decision**

Presence of Cognitive Distortion: False
Classification of Cognitive Distortion Types: N/A

Figure 5: The output of ERD for a given speech (same as that shown in Fig. 4), when the `Debate` step is activated. LLM correctly estimates that the speech does not have cognitive distortions. This shows the debiasing effect of adding `Debate` step in ERD.

# Leveraging Prompt-Learning for Structured Information Extraction from Crohn's Disease Radiology Reports in a Low-Resource Language

**Liam Hazan**[1], **Gili Focht**[5], **Naama Gavrielov**[2], **Roi Reichart**[1], **Talar Hagopian**[3],
**Mary-Louise C. Greer**[4], **Ruth Cytter Kuint**[3], **Dan Turner**[5], and **Moti Freiman**[2]

[1]Faculty of Data and Decision Sciences, Technion - Israel Institute of Technology, Haifa, Israel
[2]Faculty of Biomedical Engineering, Technion - Israel Institute of Technology, Haifa, Israel
[3]Department of Radiology, Shaare Zedek Medical Center, Jerusalem, Israel
[4]Department of Radiology, Hospital for Sick Children, Toronto, Canada
[5]The Juliet Keidan Institute of Pediatric Gastroenterology,
Shaare Zedek Medical Center, Jerusalem, Israel

## Abstract

Automatic conversion of free-text radiology reports into structured data using Natural Language Processing (NLP) techniques is crucial for analyzing diseases on a large scale. While effective for tasks in widely spoken languages like English, generative large language models (LLMs) typically underperform with less common languages and can pose potential risks to patient privacy. Fine-tuning local NLP models is hindered by the skewed nature of real-world medical datasets, where rare findings represent a significant data imbalance. We introduce SMP-BERT, a novel prompt learning method that leverages the structured nature of reports to overcome these challenges. In our studies involving a substantial collection of Crohn's disease radiology reports in Hebrew (over 8,000 patients and 10,000 reports), SMP-BERT greatly surpassed traditional fine-tuning methods in performance, notably in detecting infrequent conditions (AUC: 0.99 vs 0.94, F1: 0.84 vs 0.34). SMP-BERT empowers more accurate AI diagnostics available for low-resource languages.

## 1 Introduction

Medical imaging, particularly Computed Tomography (CT) and Magnetic Resonance Imaging (MRI), emerges as a key element in the management of complex conditions such as Crohn's Disease (CD) (Minordi et al., 2022) serving as a cornerstone for diagnosis, monitoring, and guiding treatment decisions (Bruining et al., 2018). Large-scale analyses of imaging data in CD hold promise for advancing research on the inflammatory burden in the bowel and developing predictive models of disease progression (Gu et al., 2024). The critical clinical information extracted from these images is typically embedded in free-text radiology reports, presenting a significant challenge for large-scale analysis.



Figure 1: Comparison of the median AUC and F1-score of three models (Standard Fine-tuning, SMP-BERT Zero-Shot, and SMP-BERT + tuning) over all phenotypes with 10+ positives. Error bars represent the Interquartile Range (IQR).

Manually extracting phenotypes and other pertinent information from radiology reports is labor-intensive and requires domain-specific expertise in radiology. Furthermore, CD exhibits high heterogeneity in the disease course, necessitating manual evaluation of a wide range of potential conditions (Torres et al., 2017). This task's time-consuming nature and impracticality for large-scale applications pose significant challenges in achieving efficient and accurate data extraction.

Recent attempts to automate this extraction process have utilized generative Large Language Models (LLMs) such as GPT-4, which leverage free-text instructions instead of requiring annotated data for training (Liu et al., 2023b). While these models hold promise, concerns regarding low-resource languages and data privacy remain a challenge.

Other approaches have involved directly fine-tuning open-source language models on a manually labeled subset of the data (Smit et al., 2020; Yan

Figure 2: Example of SMP-BERT Input and Output. A medical radiology report section relevant to a patient's CD diagnosis. The section labeled "Findings" serves as the input for the SMP-BERT model, similar to its pre-training phase.

et al., 2022). However, fine-tuning performance suffers from significant data imbalance, a common challenge in medical datasets and particularly in the case of CD, which features some rare conditions.

To address these limitations, we propose SMP-BERT, a novel prompt learning method built upon the "pre-train, prompt, and predict" framework (Liu et al., 2023a), specifically tailored for the structured nature of radiology reports. SMP-BERT leverages a new pre-training task called Section Matching Prediction (SMP). This task leverages the structured format of radiology reports, where key findings reside in some "Impression" section. By pre-training on this task, SMP-BERT can infer in a zero-shot setting and also further fine-tune using a relatively small amount of annotated data. This approach not only mitigates the challenge of data imbalance but also eliminates the need for massive training corpora during pre-training. This advantage makes SMP-BERT readily applicable to low-resource languages, paving the way for a more inclusive and efficient method of extracting information from radiology reports.

## 2 Related Work

### 2.1 Radiology Reports Information Extraction

Various natural language processing approaches have been used in the past to extract information and identify findings on radiology reports, from rule-based methods to deep learning–based language models (Smit et al., 2020; Mozayan et al., 2021; Tejani et al., 2022; Fink et al., 2022). While deep learning models like ClinicalBERT (Huang et al., 2019), and RadBERT (Yan et al., 2022) exploited the use of pre-training on clinical notes and radiology reports, they still require human annotation and a somewhat balanced dataset for fine-tuning.

Generative LLMs, such as GPT-4 and Cluade, may have clear advantages: They don't require extra training and can be easily instructed in natural language to do the task with high performance (Liu et al., 2023b). Unfortunately, radiology reports are usually confidential and can't be sent as a query through the Internet. Although open-source LLMs might be the solution (Mukherjee et al., 2023) they are still focused on English and struggle when it comes to low-resource languages. Moreover, even GPT4 gets comparable results to those of fine-tuned BERT in German (Adams et al., 2023) and an open-source model Vicuna-13B also gets comparable results to BERT-based model (Mukherjee et al., 2023).

### 2.2 Prompt Learning

Prompt learning (Liu et al., 2023a) is a recent advancement in Natural Language Processing (NLP) that offers a powerful alternative to traditional supervised learning methods which rely on extensive datasets for training a model $P(y|x; \theta)$. Utilizing pre-trained language models (LMs), this approach employs specific input prompts to extend the models' capabilities to tasks beyond their original training. It capitalizes on the input text's probability $P(x; \theta)$, enabling effective use of the comprehensive knowledge amassed by LMs during pre-training. Prompt learning's benefits include its efficient use of data, versatility across different tasks, and reduced need for additional extensive training.

Most prompt learning techniques are based on token-level pre-training tasks such as Left-to-Right Language Modeling (Radford et al., 2019; Brown et al., 2020) or Masked Language Modeling (Schick and Schütze, 2021a,b). However, a handful of approaches operate at the sentence level, such as (Wang et al., 2021), which reformulates the classification task into an entailment task between two sentences.

302

NSP-BERT (Sun et al., 2022) is another technique that employs sentence-level pre-training through the Next Sentence Prediction (NSP) task. It uses a structured input format beginning with a [CLS] token, followed by two sentences, A and B, separated by a [SEP] token. The training model balances instances where B genuinely follows A (IsNext) with cases where B is a random sentence (NotNext). The NSP component predicts the likelihood of B following A, relying on a specific matrix $W_{nsp}$ and the [CLS] token's hidden vector. For tasks like sentiment analysis, one might use a sentence such as "The ambiance of the restaurant was cozy and inviting," and assess if the sentiment is positive by juxtaposing it with prompts like "The sentiment of this sentence is positive." and "The sentiment of this sentence is negative.", comparing their "IsNext" probabilities. This approach allows labels to correspond with phrases of varying lengths, crucial for extracting information from radiology reports, which often contain findings described in multiple words.

NSP-BERT is optimized for classifying individual sentences, as demonstrated in the pre-training task 3. However, radiology reports consist of multiple sentences, posing a challenge for its application. Furthermore, NSP-BERT capitalizes on the logical progression found in narrative texts, where the sequence of ideas or events aids in making predictions. Contrarily, radiology reports primarily present factual details without a narrative flow, diminishing the method's effectiveness in such contexts.

## 3 SMP-BERT Framework

### 3.1 Section Matching Prediction

To overcome these challenges, we propose the Section Matching Prediction (SMP) task, designed specifically for analyzing radiology reports. These reports typically contain structured sections, notably "Findings" and "Impression". The "Findings" segment provides detailed observations from radiological examinations, while the "Impression" segment offers crucial observations and their summarized interpretations. SMP, inspired by the Next Sentence Prediction approach, considers "Findings" as the first segment and "Impression" as the follow-up. During training, "Impression" sections are accurately matched with their "Findings" counterparts half of the time (Match), and mismatched the rest (NotMatch).

Let $\mathcal{M}$ denote the model trained on our radiology reports. The model is trained on the SMP task where $x^F$ and $x^I$ represent the findings and impression sections, respectively. The model's input takes the following form:

$x_{input} = $ [CLS]$x_i^F$[SEP]$x_i^I$[EOS]

Let $q_{\mathcal{M}}(n_k|x_i^F, x_i^I)$ denotes the output probability from the model's SMP head based on the input, where $n \in \{$Match, NotMatch$\}$. The scores $s$ are computed by: $s = W_{smp}(\text{Tanh}(Wh_{\text{[CLS]}} + b))$ where $h_{\text{[CLS]}}$ represents the hidden vector of the special token [CLS] and $W_{smp}$ is the SMP head matrix. The output probability is calculated using the softmax function:

$$q_{\mathcal{M}}(n_k|x_i^F, x_i^I) = \frac{\exp s(n_k|x_i^F, x_i^I)}{\sum_n \exp s(n|x_i^F, x_i^I)}$$

This training process, optimized by a cross-entropy loss function, allows the model to discern and assess the logical link between these report sections effectively. During inference, we can leverage this learned ability to construct prompts that specifically target the presence or absence of findings in our reports.

### 3.2 Inference with SMP-BERT

In the inference stage, SMP-BERT leverages its pre-trained understanding of the connection between "Findings" and "Impression" sections. We substitute the "Impression" section with a prompt corresponding to the presence/absence of a clinical finding. By analyzing both the "Findings" section and the prompt, SMP-BERT assigns a higher probability to "Match"" when the prompt aligns with the content of the "Findings" section. The input for inference is formulated as: $x_{input} = $ [CLS]$x_i^F$[SEP]$p^j$[EOS]. Here, $p^j$ represents the prompt corresponding to the j'th label (presence/absence of a finding).

The template $\mathcal{T}$ combines the report's findings section ($x_i^F$) with generalized prompt: $\mathcal{T}(x) = $ [CLS] $x^F$ [SEP] There {is/isn't} {finding} in the {organ} [EOS]. This approach maps labels to prompts of varying lengths. A verbalizer function $f : \mathcal{Y} \to \mathcal{P}$ associates each label $y^j \in \mathcal{Y}$ with its corresponding prompt $p^j \in \mathcal{P}$. For example, let $p^j = $ "There is narrowed lumen in the Ileum" and $p^k = $ "There is **not** narrowed lumen in the Ileum" then, the prediction for report $x_i$ regarding narrowed lumen in the Ileum would be argmax $(q_{\mathcal{M}}(\text{Match}|x_i^F, p^k), q_{\mathcal{M}}(\text{Match}|x_i^F, p^j))$.

| Pre-Training Task | Pre-Training Example | Inferece Example |
|---|---|---|
| Masked Language Modeling (MLM) | ... cat sky ... / MLM head / [CLS] The [MASK] chased the mouse. | [CLS] The ambiance of the restaurant was cozy and inviting. → The sentiment of this sentence is [MASK]. / positive negative / MLM head |
| Next Sentence Prediction (NSP) | IsNext NotNext / NSP head / [CLS] [Sentence A] [SEP] [Sentence B] | IsNext NotNext / MLM head / [CLS] The ambiance of the restaurant was cozy and inviting. → [SEP] The sentiment of this sentence is **positive**. / [SEP] The sentiment of this sentence is **negative**. |
| Section Matching Prediction (SMP) | Match NotMatch / SMP head / [CLS] [Findings Section] [SEP] [Impression Section] | Match NotMatch / SMP head / [CLS] The abdominal MRI shows thickening of the bowel wall extending over 6 cm in the terminal ileum, with marked mucosal hyperenhancement. Additionally, a 3 cm segment in the proximal descending colon exhibits mild wall thickening and moderate stenosis. No evidence of fistula, abscess, or significant upstream dilation is noted. The remaining bowel loops appear normal. Liver, spleen, pancreas, and kidneys are unremarkable. → [SEP] There **is** bowel wall thickening in the Ileum / [SEP] There **is not** bowel wall thickening in the Ileum. |

Figure 3: SMP-BERT Methodology - This figure illustrates three pre-training tasks and how they can be used for text classification through prompt learning. Using MLM (token-level) for inference requires "cloze question" prompts and a verbalizer function to convert labels into single-token answers (e.g., "positive"/"negative"). Using NSP (sentence-level) is more simple. While it allows prompts of varying lengths, it's still limited to single-sentence classification. Our novel SMP solves it by pre-training on matching whole sections (multiple sentence level). Then, replace the "Impression" section with a prompt about the presence/absence of a finding.

## 3.3 SMP-tuning

The SMP-tuning process is visualized in Figure 4 and conducted similarly to the approach of NSP-tuning from NSP-BERT (Sun et al., 2022).

Generally, this process is a continuation of the SMP pre-training just given annotated reports we use the prompts instead of actual "Impression" sections. Given a sample $i$ with its reference label $y_i^+$, we define a positive instance as $(\mathcal{T}(x_i, y_i^+), \texttt{Match})$ and for each label $y_i^-$ that does not match the reference label, we define negative instances as $\{(T(x_i, y_i^-), \texttt{NotMatch})\}_{y_i^- \in Y \setminus \{y_i^+\}}$, where $Y$ is the set of all possible labels. This constructed data sums up to (n_samples*n_phenotypes*n_labels) instances and then used to fine-tune the model, leveraging the initialized weights from the SMP pre-training phase.

## 4 Experiments

### 4.1 Data

This study's dataset consists of radiology reports from three medical institutions, spanning 2010 to 2023. This dataset contains 9,683 free-text reports (one for each visit) for 8093 distinct patients. Since this dataset is confidential, no study has used it to assess the performance of any model. Ethics approval was obtained from the Shaare Zedek Medical Center Institutional Review Board (Helsinki)

committee.

For this study, a subset of 700 reports were manually annotated for the presence or absence of certain phenotypes in various organs according to the Consensus Recommendations of the American Gastroenterological Association and the Society for Abdominal Radiology (Bruining et al., 2018). The annotations focused on the following organs: organs jejunum, ileum, cecum, colon, sigmoid, and rectum. Specific findings annotated included bowel wall thickening, hyper-enhancement, pre-stenotic dilatation, narrowed lumen, restricted diffusion, and comb sign. Since our radiology reports are in the form of free text, we segmented them into "Findings" and "Impression" sections using keywords like "In summary:".

### 4.2 Experimental Setup

We divided the dataset into three distinct sets using a multi-label stratification (Sechidis et al., 2011): training (300 reports), validation (100 reports), and test (300 reports) as illustrated in Figure 5. This stratification was crucial to maintain representative distributions of labels across the sets, considering the significant class imbalance present in the majority of labels.

Our goal was to compare the performance of our method against standard fine-tuning and assess the advantages of adding the SMP-tuning step on top of the zero-shot approach.

Figure 4: SMP-tuning - Fine-tuning SMP-BERT by generating a negative and a positive instance for every annotated sample and every label. The true label is "There is finding ..." so the negative instance is paired with "There is not finding ..."



Figure 5: Flowchart of study design - The flowchart outlines the sequence of processing steps from data acquisition to model evaluation. It visualizes the progression from the initial collection of MRI and CT Hebrew radiology reports, through the stages of manual annotation and multi-label stratification, culminating in the pre-training/training of the different models.

The foundation of our models is the Hebrew RoBERTa (HeRo) model (Shalumov and Haskey, 2023), initially pre-trained on the HeDC4 corpus, a comprehensive Hebrew language corpus. We further pre-trained the model on all our radiology reports using the Masked Language Modeling (MLM) task, since there are no other open medical large corpora for Hebrew.

We conducted experiments using three models:

- **Standard Fine-tuning**: This model was fine-tuned directly for multi-label classification for all phenotypes.

- **SMP-BERT Zero-Shot**: This model was further pre-trained on all radiology reports using the SMP task. Inference was executed using the SMP-BERT methodology mentioned in the Inference section.

- **SMP-BERT + tuning**: Like the zero-shot model, this model underwent pre-training with the SMP task on all radiology reports. Additionally, it was trained further using SMP-tuning to optimize its performance.

In addition, we assessed the impact of training set size: The models were trained on datasets of varying sizes (50 to 300 reports) to analyze how the amount of training data affects their performance and ability to generalize to unseen data. We further conducted an ablation study to asses the contributions of MLM and SMP pre-training tasks to the model's performance.

Our initial goal was to compare our method with open-source generative LLMs like Llama 2. However, currently available open-source LLMs are not optimized for low-resource languages such as Hebrew, which made the comparison infeasible.

Due to the inherent class imbalance in the dataset, where most labels have a low number of positive samples, we primarily evaluated the models using the F1-score alongside the AUC metric. The F1-score considers both precision and recall, making it well-suited for imbalanced datasets. Additionally, we reported the Interquartile Range (IQR) along with the scores to provide insight into the variability and distribution of model performance across different labels.

All experiments were conducted using a single NVIDIA RTX A6000 GPU, with each experiment taking approximately 1-3 hours.

**Hyper-parameters**

For SMP-BERT + tuning, we train 6 epochs on the constructed dataset ($300 * 36 * 2 = 21600$). For standard Fine Tuning, we trained 120 epochs on the original data (300). For both we set learning rate as 2e-5 with linear decay and the batch size is 24.

## 5 Results

To account for the inherent class imbalance in our dataset, we focused our analysis on phenotypes with at least 10 positive samples, ensuring the reliability of our findings.

Our evaluation across three distinct model configurations highlighted the superior performance of the SMP-BERT + tuning approach in extracting phenotypic information from CD radiology reports. The SMP-BERT + tuning model achieved the highest median AUC of 0.99 (IQR 0.98-0.99), outperforming the Standard Fine-tuning model's median AUC of 0.94 (IQR 0.92-0.96) and the SMP-BERT Zero-Shot model's median AUC of 0.88 (IQR 0.81-0.91). For F1-score evaluations, the SMP-BERT + tuning model again leads with a median score of 0.84 (IQR 0.76-0.94), which is substantially higher than the scores of the Standard Fine-tuning model (0.34, IQR 0.22-0.85) and the SMP-BERT Zero-Shot model (0.58, IQR 0.55-0.62). A comprehensive breakdown of these results, including F1 and AUC scores for individual phenotypes, is detailed in the accompanying Table 1.

Further analysis presented in Figure 7 of model performance relative to the count of positive instances exhibited the strength of SMP-BERT + tuning, particularly for labels with sparse positives in the training set. For example, with only 19 positive cases for "Rectum Bowel Wall Thickening," SMP-BERT + tuning achieved a significantly higher F1-score (0.74) compared to the standard model (0.1). This demonstrates its superior ability to generalize well from limited data.

However, both models performed well when dealing with abundant positive instances. For example, with 137 positives for "Ileum Bowel Wall Thickening" (almost half the dataset), both models achieved good results, with SMP-BERT + tuning maintaining a decent gap (F1-score 0.97 vs. 0.915 for the standard model).

The graph shown in Figure 7 indicates that the performance gap between the models decreases with an increase in the number of positive instances. This suggests that while SMP-BERT + tuning shines with limited data, it still performs better when more data is available.

We also analyzed how the size of the training set impacts model performance. As shown in Figure 6, the SMP-BERT + tuning model exhibits superior adaptability. Notably, it achieves good performance even with limited training data (50-100 samples). The Standard Fine-tuning model exhibits a trend of broadening IQRs and decrease of median score. This could suggests an improving performance for common phenotypes (like Ileum Bowel Wall Thickening) but potentially decreasing performance for rarer ones due to increased data imbalance.

**Ablation Study**

As evidenced by Table 2, both pre-training tasks, MLM and SMP, significantly contribute to optimizing the performance of SMP-BERT. Moreover, it appears that standard fine-tuning benefits from the inclusion of the SMP task.

## 6 Discussion

This study examined the efficacy of SMP-BERT, a novel prompt-learning approach, in extracting detailed information from Hebrew radiology reports of CD patients. Our results reveal that SMP-BERT, especially the fine-tuned version (SMP-BERT + tuning), significantly outperforms the standard fine-tuning approach, , achieving an improvement of 49% in median F1 score and 5% in median AUC.

Our study highlights the significant improvement of SMP-BERT + tuning, achieving superior F1-scores and AUCs compared to standard fine-tuning across all analyzed phenotypes. Notably, the model performs well even with a low amount of annotated data. This improvement is particularly notable for rarer phenotypes, demonstrating the model's ability to handle imbalanced datasets, a common challenge in the medical domain. This robustness is crucial for advancing research in CD and other conditions with diverse clinical presentations.

Furthermore, this study contributes to the growing exploration of prompt learning for NLP tasks in healthcare. Unlike traditional fine-tuning approaches, which require substantial labeled data, SMP-BERT leverages pre-training on the "Section Matching Prediction" task and further SMP-tuning

Figure 6: Median F1 scores and IQRs for SMP-BERT + tuning and Standard fine-tuning trained on different training set sizes.



Figure 7: This line chart plots the F1 scores against the number of positive instances of all phenotypes in the dataset (300 total).

to achieve exceptional performance even with limited data. This opens exciting possibilities for applying prompt learning in scenarios with limited annotated data, imbalanced data, or low-resource languages, pushing the boundaries of NLP applications in healthcare.

## References

Lisa C. Adams, Daniel Truhn, Felix Busch, Avan Kader, Stefan M. Niehues, Marcus R. Makowski, and Keno K. Bressem. 2023. Leveraging gpt-4 for post hoc transformation of free-text radiology reports into structured reporting: A multilingual feasibility study. *Radiology*, 307(4).

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

David H. Bruining, Ellen M. Zimmermann, Edward V. Loftus, William J. Sandborn, Cary G. Sauer, and Scott A. Strong. 2018. Consensus recommendations for evaluation, interpretation, and utilization of computed tomography and magnetic resonance enterography in patients with small bowel crohn's disease. *Radiology*, 286(3):776–799.

Matthias A. Fink, Klaus Kades, Arved Bischoff, Martin Moll, Merle Schnell, Maike Küchler, Gregor Köh-

ler, Jan Sellner, Claus Peter Heussel, Hans-Ulrich Kauczor, Heinz-Peter Schlemmer, Klaus Maier-Hein, Tim F. Weber, and Jens Kleesiek. 2022. Deep learning–based assessment of oncologic outcomes from natural language processing of structured radiology reports. *Radiology: Artificial Intelligence*, 4(5).

Phillip Gu, Oreen Mendonca, Dan Carter, Shishir Dube, Paul Wang, Xiuzhen Huang, Debiao Li, Jason H Moore, and Dermot P B McGovern. 2024. Ai-luminating artificial intelligence in inflammatory bowel diseases: A narrative review on the role of ai in endoscopy, histology, and imaging for ibd. *Inflammatory Bowel Diseases*.

Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. 2019. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *arXiv:1904.05342*.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023a. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35.

Qianchu Liu, Stephanie Hyland, Shruthi Bannur, Kenza Bouzid, Daniel Castro, Maria Wetscherek, Robert Tinn, Harshita Sharma, Fernando Pérez-García, Anton Schwaighofer, Pranav Rajpurkar, Sameer Khanna, Hoifung Poon, Naoto Usuyama, Anja Thieme, Aditya Nori, Matthew Lungren, Ozan Oktay, and Javier Alvarez-Valle. 2023b. Exploring the boundaries of gpt-4 in radiology. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Laura Maria Minordi, Antonio Bevere, Alfredo Papa, Luigi Larosa, and Riccardo Manfredi. 2022. Ct and mri evaluations in crohn's complications: A guide for the radiologist. *Academic Radiology*, 29(8):1206–1227.

| Organ-Finding | SMP-BERT + tuning | SMP-BERT Zero-Shot | Standard Fine-Tuning | prevalence |
|---|---|---|---|---|
| ileum-bowel wall thickening | **0.97/1.0** | 0.85/0.91 | 0.92/0.98 | 44% |
| ileum-enhancement | **0.95/0.99** | 0.77/0.86 | 0.86/0.95 | 36% |
| ileum-narrowed lumen | **0.96/1.0** | 0.79/0.92 | 0.85/0.97 | 19% |
| ileum-dilatation | **0.96/1.0** | 0.62/0.87 | 0.86/0.96 | 18% |
| ileum-comb sign | **0.9/0.99** | 0.48/0.81 | 0.84/0.97 | 15% |
| ileum-restricted diffusion | **0.94/0.99** | 0.78/0.91 | 0.9/**0.99** | 16% |
| colon-bowel wall thickening | **0.84/0.98** | 0.58/0.88 | 0.62/0.93 | 12% |
| colon-enhancement | **0.92/0.99** | 0.57/0.88 | 0.59/0.95 | 9% |
| colon-comb sign | **0.86/1.0** | 0.18/0.74 | 0.33/0.94 | 3% |
| colon-restricted diffusion | **0.76/0.98** | 0.33/0.94 | 0.29/0.91 | 3% |
| rectum-bowel wall thickening | **0.74/0.96** | 0.56/0.89 | 0.1/0.96 | 6% |
| rectum-enhancement | **0.76/0.98** | 0.59/0.78 | 0.22/0.89 | 5% |
| sigmoid-bowel wall thickening | **0.75/0.97** | 0.55/0.77 | 0.3/0.9 | 10% |
| sigmoid-enhancement | **0.7/0.98** | 0.58/0.89 | 0.34/0.88 | 7% |
| sigmoid-comb sign | **0.53/0.98** | 0.31/0.78 | 0.17/0.93 | 3% |
| cecum-bowel wall thickening | **0.77/0.98** | 0.62/0.89 | 0.12/0.93 | 5% |
| cecum-enhancement | **0.82/0.99** | 0.56/0.93 | 0.0/0.92 | 3% |

Table 1: Performance comparison. Values are F1/AUC scores for each model across different phenotypes. The Prevalence column indicates the percentage of test samples in which the phenotype is present.

| Method | MLM | SMP | F1-Score | AUC |
|---|---|---|---|---|
| SMP-BERT + tuning | ✓ | ✓ | **0.84 [0.76,0.94]** | **0.99 [0.98,0.99]** |
| | ✓ | ✗ | 0.75 [0.59,0.87] | 0.97 [0.96,0.98] |
| | ✗ | ✓ | 0.73 [0.67,0.89] | 0.97 [0.95,0.98] |
| | ✗ | ✗ | 0.42 [0.26,0.57] | 0.94 [0.92,0.96] |
| Standard Fine-tuning | ✓ | ✓ | **0.55 [0.35,0.86]** | **0.96 [0.95,0.98]** |
| | ✓ | ✗ | 0.34 [0.22,0.85] | 0.94 [0.92,0.96] |
| | ✗ | ✓ | 0.15 [0.0,0.72] | 0.85 [0.82,0.91] |
| | ✗ | ✗ | 0.12 [0.0,0.61] | 0.83 [0.78,0.88] |

Table 2: Ablation Study on Pre-training Tasks.

Ali Mozayan, Alexander R. Fabbri, Michelle Maneevese, Irena Tocino, and Sophie Chheang. 2021. Practical guide to natural language processing for radiology. *RadioGraphics*, 41(5):1446–1453.

Pritam Mukherjee, Benjamin Hou, Ricardo B. Lanfredi, and Ronald M. Summers. 2023. Feasibility of using the privacy-preserving large language model vicuna for labeling radiology reports. *Radiology*, 309(1).

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Timo Schick and Hinrich Schütze. 2021a. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics.

Timo Schick and Hinrich Schütze. 2021b. It's not just size that matters: Small language models are also few-shot learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.

Konstantinos Sechidis, Grigorios Tsoumakas, and Ioannis Vlahavas. 2011. *On the Stratification of Multi-label Data*, page 145–158. Springer Berlin Heidelberg.

Vitaly Shalumov and Harel Haskey. 2023. Hero: Roberta and longformer hebrew language models. *arXiv preprint arXiv:2304.11077*.

Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Ng, and Matthew Lungren. 2020. Combining automatic labelers and expert annotations for accurate radiology report labeling using bert. In *Proceedings of the 2020 Conference on Empirical Meth-*

*ods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.

Yi Sun, Yu Zheng, Chao Hao, and Hangping Qiu. 2022. NSP-BERT: A prompt-based few-shot learner through an original pre-training task —— next sentence prediction. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3233–3250, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Ali S. Tejani, Yee S. Ng, Yin Xi, Julia R. Fielding, Travis G. Browning, and Jesse C. Rayan. 2022. Performance of multiple pretrained bert models to automate and accelerate data annotation for large datasets. *Radiology: Artificial Intelligence*, 4(4).

Joana Torres, Saurabh Mehandru, Jean-Frédéric Colombel, and Laurent Peyrin-Biroulet. 2017. Crohn's disease. *The Lancet*, 389(10080):1741–1755.

Sinong Wang, Han Fang, Madian Khabsa, Hanzi Mao, and Hao Ma. 2021. Entailment as few-shot learner.

An Yan, Julian McAuley, Xing Lu, Jiang Du, Eric Y. Chang, Amilcare Gentili, and Chun-Nan Hsu. 2022. Radbert: Adapting transformer-based language models to radiology. *Radiology: Artificial Intelligence*, 4(4).

# Context Aggregation with Topic-focused Summarization for Personalized Medical Dialogue Generation

**Zhengyuan Liu, Siti Umairah Md Salleh, Pavitra Krishnaswamy, Nancy F. Chen**
Institute for Infocomm Research (I²R), A*STAR, Singapore
`{liu_zhengyuan,nfychen}@i2r.a-star.edu.sg`

## Abstract

In the realm of dialogue systems, generated responses often lack personalization. This is particularly true in the medical domain, where research is limited by scarce available domain-specific data and the complexities of modeling medical context and persona information. In this work, we investigate the potential of harnessing large language models for personalized medical dialogue generation. In particular, to better aggregate the long conversational context, we adopt topic-focused summarization to distill core information from the dialogue history, and use such information to guide the conversation flow and generated content. Drawing inspiration from real-world telehealth conversations, we outline a comprehensive pipeline encompassing data processing, profile construction, and domain adaptation. This work not only highlights our technical approach but also shares distilled insights from the data preparation and model construction phases.

## 1 Introduction

Medical dialogue systems hold significant potential for improving the efficiency of clinical workflows (Xu et al., 2021). As a specialized form of task-oriented dialogue, medical dialogue typically involves the completion of multiple tasks, including diagnosis, question answering, and consultation (Althoff et al., 2016; Tian et al., 2019; Xia et al., 2020; Gupta et al., 2020). There has been significant progress in this research field of the dialogue system in past years with the development of contextualized representation learning and neural language generation (Xu et al., 2019; Palanica et al., 2019). However, the general-purpose conversational interactive systems are proven to be inadequate, as they cannot adapt their responses to the unique medical histories and the diverse user preferences and personalities (Li et al., 2016; Mazaré et al., 2018). Personalized dialogue systems, tailored to the specific needs and characteristics of dif-



Figure 1: One dialogue example for "physical activity customized coaching" based on the personalized medical dialogue generation.

ferent users, can potentially bridge this gap (Ghosh et al., 2018; Schloss and Konam, 2020). By leveraging patient profiles, such as medical records, demographic information, and previous interactions, the personalized systems can facilitate more nuanced, empathetic, and context-aware conversations. This level of personalization not only enhances patient engagement and satisfaction, but also has the potential to improve healthcare outcomes by fostering adherence to treatment plans and providing tailored health education.

In this work, we conduct a case study on a clinical conversation scenario. Because of the chronic nature of diabetes and its associated complications, it requires constant attention and regular follow-up operation (Piette et al., 2000; Lawson et al., 2005). In practice, nurses schedule calls with patients to track their compliance status and health condition, provide general education, and customized coaching and lifestyle advice (Piette et al., 2001; Kivelä et al., 2014). To facilitate the communication process and deliver more efficient health management, the follow-up calls are organized according to a

medical protocol and telecarers adjust the conversation topics based on the patient's lifestyle management status and medication records (Kirkman et al., 1994; Taylor et al., 2003). This renders the follow-up call a representative use case for personalized dialogue generation. For instance, customized coaching is an effective patient education method (Kivelä et al., 2014), and its sub-topics are strongly correlated to the patient profile (as the example shown in Figure 1). The challenges of developing a personalized medical dialogue system come from three fundamental aspects: the lack of domain-specific data (Zhou et al., 2022); the complexity of modeling medical context and persona information (Liu et al., 2022a); and how to extensively evaluate the system (Abbasian et al., 2023). Moreover, due to the verbal nature of human spoken dialogues, the follow-up calls are often lengthy by covering various topics, which results in a low information density. The noisy long context also poses challenges for modeling and generation. We thus propose and adopt topic-focused summarization to distill and aggregate core information of the dialogue context, and use such information to guide the subsequent conversation flow and content generation.

In practice, to bootstrap the data-driven approaches, we construct a sample set derived from human spoken conversations, and we leverage the advancements in Large Language Models (LLMs) (Ouyang et al., 2022; Touvron et al., 2023b) for developing the dialogue system, which have demonstrated their exceptional language understanding and generation capabilities in the medical domain (Singhal et al., 2023). We add user profile information to produce personalized conversation, and improve the generation coherence based on topic-level context aggregation. Experiments show that our proposed method can substantially improve the generation quality, especially in the long context setting. This work not only highlights the technical approach but also shares distilled insights from the data preparation and model construction phases.

## 2 Related Work

**Medical Dialogue Generation** Medical dialogue systems aim to provide medical services for patients (Xu et al., 2021). As one specialized form of a task-oriented dialogue system, many previous studies focus on making diagnostic predictions after gathering patients' information of symptoms (Wei et al., 2018; Xu et al., 2019; Zhou et al., 2021), and healthcare counseling (Cao et al., 2019; Shen et al., 2020). Data-driven approaches and methods are proposed and applied for medical dialogue generation upon the development of large-scale medical dialogue datasets such as MedDialog (Zeng et al., 2020) and MedDG (Liu et al., 2022a), and the scarcity of domain-specific data still poses this task as a low-resource challenge (Lin et al., 2021).

**Personalized Dialogue Systems** One-size-fits-all approaches to human-machine communication have shown limitations in accommodating the diverse needs, preferences, and contexts of individual users. By contrast, personalized dialogue systems (Li et al., 2016; Mazaré et al., 2018) offer the potential to transcend these limitations by tailoring interactions to unique characteristics and requirements, thus raising much research interest. In particular, improving the modeling of persona or user information is one of the key points, and there are different approaches proposed in previous studies, such as explicitly utilizing pre-defined persona attributes to generate conditional responses (Qian et al., 2018; Olabiyi et al., 2019), constructing user embeddings to enhance personalized dialogue generation (Li et al., 2016; Chan et al., 2019), and building implicit user information from dialogue history (Al-Rfou et al., 2016; Ma et al., 2021).

**Language Models as Conversational Agent** Leveraging pre-trained language backbones for building conversation agents has seen remarkable progress recently (Liao et al., 2023), and the recent large language models have demonstrated impressive capabilities in both open-domain and task-oriented scenarios (Zhang et al., 2020; Thoppilan et al., 2022). Instruction tuning is one efficient and effective way to enable the conversational capabilities of large language models, such as Alpaca (Taori et al., 2023) and Vicuna (Chiang et al., 2023). It has been proved that using reinforcement learning with human feedback can further optimize language models for human-machine interaction, and the LLMs not only take conversation in a human-like manner, but also can do task solving and complex reasoning (Ouyang et al., 2022). Furthermore, LLMs demonstrate strong language understanding and generation capabilities in various downstream tasks that require certain domain knowledge (Wang et al., 2022; Hendrycks et al., 2020) (even in the zero-shot setting), which benefits from their large-scale pre-training (Touvron et al., 2023a).

| Intent | Topic Type | Example |
|---|---|---|
| Information Gathering | Identification, Medical Experience, Appointments, Programme, Vitals, Insulin, Hyper/Hypo Incident, Base Compliance | [**Topic**: Vitals] **Nurse**: Can you tell me your blood sugar level four hours after dinner? **Patient**: If I remember correctly, it was around 13.4. **Nurse**: And what about your post-dinner reading? **Patient**: Ah, yes. After dinner, it was around 23 to 24, if I'm not wrong. |
| General/Customized Coaching | Self-Monitoring, Diet Management, Insulin, Physical Activity, General Education | [**Topic**: Diet Management] **Nurse**: From a dietary perspective, do you have any issues? **Patient**: no no **Nurse**: Are you okay with your diet? **Patient**: Yes, I'm fine. **Nurse**: Okay, good. A bit difficult, but you have to control it. **Patient**: I know, I have to be disciplined for my own health. |
| Other | Introduction, Social Chatting, Financial and Social Aid | [**Topic**: Social Chatting] **Nurse**: Never mind, this computer is taking a while to respond. **Patient**: Okay, Okay. **Nurse**: We'll have to wait for a bit. **Patient**: Ok, no problem. |

Table 1: List of the dialogue topics and their intent categorization.

# 3 Personalized Dialogue Generation: Data Preparation & Refinement

In this work, we conduct a case study on personalized follow-up calls for diabetes patients. Diabetes is a chronic metabolic disorder characterized by abnormal glucose regulation, and effective management of diabetes is essential to mitigate its associated complications and improve patients' overall quality of life (Lawson et al., 2005). In practical use cases, the general-purpose messages may not adequately address the unique needs of individual patients. For example, customized coaching of physical activity should take into account factors such as the patient's age, comorbidities, lifestyle, and psychosocial aspects. By recognizing the heterogeneity of diabetes patients and offering tailored coaching interventions, it is useful for improving health management.

## 3.1 Raw Data Collection and Statistics

The raw data are extracted from call recordings of diabetes health management conversations (Liu et al., 2023) and fully anonymized.[1] Speech transcribers are employed for manual speech-to-text conversion to ensure quality. Speaker roles (e.g., nurse, patient, caregiver) are added to each utterance, and the informal and spontaneous styles of spoken dialogues such as back-channeling, hesitation, and repetition are preserved. The dialogue segmentation and topic categorization are manually



Figure 2: Feature visualization of segment embeddings via t-SNE. The colored points denote topically coherent segments labeled with different topics.

performed.[2] Our linguistic annotators are familiar with clinical conversations, and have finished a training session on diabetes health management. Topic categories are built on the medical protocol refined by the healthcare provider. Moreover, there have been interactions for the corpus construction, where we collect feedback from nurses, refine the annotation scheme, and update the whole corpus.

The transcribed dataset contains 856 transcripts. Depending on the patient's medical history and phases of the healthcare programme, nurses schedule their follow-up calls differently, and this results in length and topic variation. We obtain the segment representations from an unsupervised sentence embedding model (Gao et al., 2021), and use t-SNE (Van der Maaten and Hinton, 2008) to illustrate their distribution in a 2-dimensional space. As shown in Figure 2, dialogue utterances in different topics are semantically diverse and distinct. Moreover, there are two major types of dialogue

---

[1]This research study was approved by the SingHealth and A*STAR Institutional Review Boards. Participants enrolled in the healthcare programme consented to use of anonymized versions of their data for research.

[2]All dialogue examples in this manuscript are dummy data for demonstration purposes.

Figure 3: One dummy example of the spoken language conversion. Sentences are normalized and adjacent utterances with the same speaker are combined.



Figure 4: One dummy example of the patient profile. The basic information and summary from information gathering topics are collected.

intent: information gathering and general and customized coaching. As shown in Table 1, there are four topics that are strongly related to customized coaching: physical activity, diet management, insulin, and self-monitoring, which usually shows a strong dependency on the dialogue context, as nurses will adjust the dialogue content based on the patient's response and feedback.

## 3.2 Spoken Language Conversion

While both human-human and human-machine medical conversations are task-oriented and topically organized, they demonstrate distinct linguistic characteristics, especially from the lexical and syntactic perspectives (Bernsen et al., 1996). More specifically, compared with real-world spoken dialogues, there is much less informal and colloquial wording in the human-machine interaction (Hill et al., 2015). Directly training on the raw transcripts will result in issues such as verbose sentences, unnecessary repetition, and incomplete utterances. Therefore, to improve the formality and readability of machine-generated responses, we conduct a spoken language conversion on the transcribed samples. As the example shown in Figure 3, there are three basic pre-processing steps: (1) We adopt an off-the-shelf text normalization model to process the utterances (Liu et al., 2022b). The colloquial sentences are paraphrased and the grammar errors are corrected. (2) We further normalize the utterances by reducing other common spoken language features such as repetition, pauses, and fillers. (3) To construct the turn-by-turn interaction for human-machine conversation, adjacent utter-

ances with the same speaker are combined.[3] In our corpus preparation, we observe that the normalization step brings substantial changes in most utterances, and the processed sample set is significantly distinct from the raw dialogue data.

## 3.3 Patient Profile Construction

Considering each patient's health condition and personal preferences, telecarers adjust their health management advice and provide general and customized coaching (Piette et al., 2000; Lawson et al., 2005). For instance, when discussing the type and frequency of physical activity, nurses should ask patients who have hypoglycemia symptoms to pay more attention to their sugar levels during exercise. Therefore, a feature-rich profile should include both basic demographic information, and up-to-date health condition of patients. To this end, aside from the basic information (e.g., age, gender, scheduled call phase) extracted from a structured database,[4] we also collect the key discussed points from the information gathering topics, as shown in Figure 4. In our clinical data, the gathered information from each follow-up call is recorded in a human-written summary. When such manually collected information is not available, automated approaches such as entity and event extraction can also be used for information extraction.

## 4 Context Aggregation via Topic-focused Summarization

Due to the complexity and verbose nature of human spoken dialogues (Sacks et al., 1978), and the necessity to cover multiple topics in clinical follow-up calls, nurse-to-patient conversations are

---

[3]Since our raw data contain topic-level annotation, we conduct the normalization process on each topic segment.

[4]Both language and structured data are fully anonymized, without any identifiable personal information.

Figure 5: One dummy dialogue example in two topics. Frames indicate topically-coherent segments, and their corresponding label is highlighted.

often lengthy and thus characterized by lower information density than other document formats. For instance, in our transcribed calls, the maximum, median, and minimum utterance numbers are 1996, 221, and 21, respectively; the maximum, median, and minimum number of words are 16701, 1684, and 70, respectively. Nearly 5% samples (at the 95% quantile) are comprised of more than 800 utterances (6000 words). This requires models to precisely capture the core information from the long dialogue context and poses challenges for dialogue systems in both modeling and generation. In this work, we propose and adopt topic-focused summarization, to distill and aggregate the salient pieces from a noisy dialogue context. The refined context is then leveraged to guide the subsequent generation, and improve relevance and coherence. More specifically, we leverage the large language models to generate dialogue summaries for each dialogue snippet about a certain topic, and concatenate them as the history context. We conduct the following steps to build samples for training the data-driven approach:

### 4.1 Topic Segmentation and Categorization

First, each dialogue is processed with topic segmentation and topic categorization, as shown in Figure 5. This step is to parse the conversation into coherent segments, and helps identify the underlying structure of the dialogue. Here we use the manual annotated information in both the training and testing process: each training sample is to generate one coherent dialogue segment with a topic label and previous dialogue context, and it ends with a '*<topic-end>*' token for boundary modeling and a topic label of next segment prediction, which is a supervised approach for the dialogue topic modeling.



Figure 6: One dummy example of topic-focused summarization. The corresponding topic label is in brackets.

### 4.2 Topic-focused Summarization

For each identified segment, we then distill the core information by using a dialogue summarization model. In our preliminary study, we found that prompting large language models can produce reasonable dialogue summaries in the clinical scenario. We thus employ a state-of-the-art open model (i.e., Mistral-7B-Instruct-v0.2) for this step.[5] As shown in Figure 6, the summarizer is able to capture salient spans in the dialogue, and generate a concise version. Moreover, to better incorporate the dialogue topic information (Liu et al., 2019), we add their corresponding topic label before each summary.

### 4.3 Dialogue Generation Integration

The generated summaries serve as the historical context for the dialogue system. Since there is more than one topic segment in the conversation, we concatenate all summaries as one context and feed it into the system for subsequent generations. The response generation process is informed by a concentrated version of the dialogue history, emphasizing relevance and topic coherence. This enables the system to generate responses that are not only contextually appropriate but also enriched with the distilled essence of the prior conversation.

## 5 Personalized Dialogue Generation: Training & Evaluation

### 5.1 Task Definition

In a multi-turn human-machine conversation, we define $C_i$ as the profile of the user $i$, and at a turn $t$, $U_t$ is the user input and $S_t$ is the system's response.

---

[5]The user prompt for the summarization step is *"Given the following nurse-patient dialogue about <topic-label>, please write a concise summary: <dialogue-content>."*

Figure 7: Overview of the pipeline for training and inference with personalized medical dialogue generation.

Basically, for modeling the dialogue history, all previous turns are concatenated and fed to the system as input: $H = [U_0, S_0, U_1, S_1, ..., U_{t-1}, S_{t-1}]$. In our framework of personalized dialogue generation with context aggregation, the user profile $C$ and topic-focused summaries $H_{summary}$ are also part of context information. Therefore, at a turn $t$, the system's response $S_t$ is conditioned on profile information $C_i$, summarized context $H_{summary}$, in-topic context $H_{topic}$ and user's current utterance $U_t$, which are concatenated as a single sequence. To allow for handling descriptive profiles, we retain the profile $C_i$ in the form of natural language text, in contrast to previous studies that encode the profile features via one-hot encoding and limit the model's accessibility to various features.

## 5.2 Adapting LLMs as Conversation Agents

Large language models have been shown to achieve remarkable performance across a variety of natural language tasks. Aside from their versatile capabilities of language understanding and generation where expert knowledge is not required, LLMs also show impressive results in medical document processing and decision support, and obtained comparable scores in medical examinations to human (Singhal et al., 2023). By learning from large volumes of text data to predict the subsequent tokens, LLMs with the auto-regressive framework can generate coherent, fluent, and reasonable responses to diverse prompts, and they are adopted as the

conversation agents via in-context learning and instruction tuning (Chiang et al., 2023). To leverage the large-scale language backbone and adapt it to our domain-specific use case, we conduct experiments on some representative large language models, such as LLaMA (Touvron et al., 2023b) and Mistral (Jiang et al., 2023) on the profile-aware dialogue samples[6], and improve the efficiency of the training process from data and model perspective.

### 5.2.1 Parameter-Efficient Training

One major challenge of utilizing LLMs is the high demand for computational resources for adaptive training. To fine-tune LLMs in a low-resource setting, here we employ parameter-efficient approaches: Low-rank adaption (LoRA) (Hu et al., 2021) and QLoRA (Dettmers et al., 2024). Previous studies show that the over-parameterized models in fact reside on a low intrinsic dimension. Compared with full-parameter training, LoRA and QLoRA update to the weight matrices with a low-rank matrix factorization, and significantly reduces the number of trainable parameters, and speeds up training with little impact on the final performance.

### 5.2.2 Dialogue-level Efficient Training

Given one multi-turn dialogue sample, at the fine-tuning stage, generally, only the system responses are used for loss calculation and weight updating. In practice, if we split a $n$-turn dialogue into $n$

---

[6]All open models used in this work are only for research use. We follow their corresponding license in our experiments.

| Model Type | BLEU-2 | BLEU-3 | ROUGE-1 | ROUGE-2 | ROUGE-L | SimCSE |
|---|---|---|---|---|---|---|
| LLaMA-2 7B | 4.314 | 2.517 | 12.75 | 2.500 | 13.09 | 28.90 |
| + Utterance Normalization | 5.752 | 3.521 | 17.04 | 4.052 | 18.24 | 41.89 |
| + Context Aggregation | 7.087 | 4.533 | 18.49 | 4.183 | 20.36 | 44.40 |
| LLaMA-2-Chat 7B | 4.530 | 2.788 | 13.00 | 2.553 | 13.18 | 28.46 |
| + Utterance Normalization | 7.849 | 5.205 | 18.84 | 6.001 | 21.05 | 42.95 |
| + Context Aggregation | 9.313 | 7.344 | 20.27 | 6.492 | 21.96 | 44.38 |
| LLaMA-2-Chat 13B | 4.526 | 2.625 | 12.78 | 2.711 | 13.85 | 29.77 |
| + Utterance Normalization | 8.160 | 5.596 | 20.99 | 5.205 | 23.88 | 45.13 |
| + Context Aggregation | 10.53 | 7.227 | 22.77 | 6.544 | 26.16 | 49.03 |
| Mistral-7B | 4.434 | 2.406 | 13.06 | 2.501 | 14.28 | 30.06 |
| + Utterance Normalization | 8.782 | 6.441 | 19.74 | 6.353 | 21.75 | 42.97 |
| + Context Aggregation | 11.29 | 8.248 | 21.68 | 10.11 | 25.34 | 48.98 |
| Mistral-7B-Instruct-v0.2 | 4.878 | 2.957 | 14.34 | 2.901 | 13.28 | 28.86 |
| + Utterance Normalization | 7.942 | 5.341 | 18.96 | 6.541 | 21.88 | 46.24 |
| + Context Aggregation | 11.76 | 8.358 | 22.36 | 9.783 | 26.40 | 53.19 |

Table 2: Experimental results with automated evaluation metrics on topically-coherent dialogue generation.

turn-level samples, the learning step increases by a factor of $n$. To improve training efficiency, here we leverage the properties of causal language models since each token only depends on its precedent tokens. Therefore, we feed the entire dialogue sequence to the decoder-only model, and mask out the user utterances, and compute the loss of all system responses in parallel.

### 5.2.3 Balanced Data Sampling

Since the sample number of customized coaching is limited, we mixed dialogue segments from other topics for training data augmentation. The frequency distribution of different topics is imbalanced. For instance, compared with the topic "oral medication", the "general education" is more frequently discussed and presents a larger utterance number. When fine-tuning the language backbone, a diverse and balanced sample set can bring higher performance, we thus construct the training set by sampling a balanced ratio at the topic level.

## 6 Experiments and Results

### 6.1 Experimental Setting

The processed conversational data (5.0K topic-level dialogue samples) are used for training, and we randomly select 10% for validation and testing (0.5K samples) respectively. The maximum length of the dialogue sequence is set at 2048. *AdamW* optimizer is used with a learning rate of 1e-5, the batch size with gradient accumulation is set at 64, and the epoch number is 5. Best checkpoints are selected based on validation results using cross-entropy loss. Models are imple-

mented with PyTorch[7] and HuggingFace Transformers[8]. Parameter-efficient fine-tuning is applied with PEFT (Mangrulkar et al., 2022), and the rank $k$ in LoRA adaptation is set at 16. Following previous work, we add the projection layers of the Transformer network to the LoRA training process, and the trainable parameter sizes of LLaMA-2-7B/Mistral-7B and LLaMA-2-13B are 2.32M and 3.63M, respectively. All experiments are run on a single Nvidia A100 GPU with 40G memory.

### 6.2 Evaluation Metrics

Following previous work (Shen et al., 2020), we use two lexical automated evaluation metrics: BLEU (BLEU-2 and BLEU-4) and ROUGE (ROUGE-1, ROUGE-2 and ROUGE-L) (Lin, 2004), as well as the embedding-based metrics SimCSE (Gao et al., 2021). All reported scores are rescaled to percentage values. For each topically coherent dialogue segment ended with '*<topic-end>*', we calculate the averaged evaluation scores of each nurse's utterance. Speaker role tokens (e.g., *Nurse*, *Patient*) and model-generated special tokens (e.g., *</s>*, *[INST]*) are not included.

### 6.3 Evaluation Results & Analysis

We use a hold-out test set to evaluate the generated nurse responses. In our experiments, we indicate gold topic labels for model comparison. Since personalized dialogue generation is mainly for delivering customized education or consultation, we thus focus on evaluating the four customized coaching

---
[7]https://pytorch.org
[8]https://github.com/huggingface/transformers

| Model Type | BLEU-2 | BLEU-3 | ROUGE-1 | ROUGE-2 | ROUGE-L | SimCSE |
|---|---|---|---|---|---|---|
| LLaMA-2-Chat 7B | 7.012 | 4.336 | 17.79 | 3.848 | 19.98 | 41.80 |
| - Context Aggregation | 5.997 | 3.409 | 16.67 | 3.101 | 17.42 | 35.75 |
| - Patient Profile | 4.473 | 3.166 | 11.24 | 2.630 | 11.02 | 29.85 |
| LLaMA-2-Chat 7B | 8.334 | 5.554 | 18.40 | 5.760 | 21.24 | 43.20 |
| - Context Aggregation | 6.019 | 3.757 | 17.65 | 3.221 | 19.65 | 38.28 |
| - Patient Profile | 4.509 | 3.190 | 11.53 | 2.981 | 11.63 | 30.42 |
| Mistral-7B | 9.636 | 7.022 | 21.53 | 7.319 | 23.14 | 48.65 |
| - Context Aggregation | 6.914 | 5.065 | 15.68 | 4.751 | 18.47 | 35.72 |
| - Patient Profile | 5.303 | 3.682 | 12.52 | 2.673 | 13.79 | 34.78 |
| Mistral-7B-Instruct-v0.2 | 10.19 | 7.282 | 21.31 | 7.520 | 24.19 | 48.30 |
| - Context Aggregation | 6.062 | 4.252 | 15.28 | 3.808 | 17.30 | 36.12 |
| - Patient Profile | 5.136 | 3.508 | 12.06 | 2.351 | 13.51 | 35.11 |

Table 3: Ablation study on the context aggregation via topic-focused summarization at the inference stage.

topics: self-monitoring, diet management, insulin, and physical activity.

### 6.3.1 Dialogue Generation Evaluation

Table 2 shows the results of dialogue generation by training the representative open LLMs (e.g., LLaMA, Mistral). Here we report evaluation results of modeling training with our proposed enhancements: the human spoken dialogue data refinement (i.e., utterance normalization) and context modeling and aggregation (i.e., utilizing topic-focused summarization). As shown in Table 2, the generation quality benefits a lot from adopting utterance normalization on all tested models and at all metrics. This is because human conversations contain many spoken linguistic features such as fillers, thus training on the original noisy spoken data affects the generation quality significantly, models tend to produce less meaningful and fluent sentences. Therefore, to build reasonable human-machine conversational interaction, it is necessary to include the normalization step in the spoken dialogue samples. On the other hand, compared with other language generation tasks such as machine translation, the overall evaluation scores of dialogue response generation are at a low level, this is mainly due to the utterance diversity in the nurse-patient conversations.

Moreover, adding context aggregation with topic-focused summarization also significantly improves the scores, demonstrating its effectiveness of coherent personalizing response generation. Considering the scoring alignment between lexicon-based and embedding-based metrics, the overall evaluation ranks are consistent across the three tested metrics: BLEU, ROUGE, and SimCSE. Upon the summarization process, the historical context length can be reduced to 20% of the

original length, with dense information in a formal wording. This is also beneficial for the model to capture important features to organize the subsequent generations. Surprisingly, in our experimental setting, we observe that instruction-tuned models (e.g., LLaMA-2-Chat, Mistral-7B-Instruct) did not show substantial gain over the pre-training foundation models, and scores become even lower in some metrics (e.g., BLEU-2, ROUGE-2) when training with utterance normalization. As LLMs contain massive prior knowledge from large-scale pre-training, both model types could achieve the same dialogue modeling and generation capabilities after domain-specific adaptation on one downstream task.

### 6.3.2 Leveraging LLMs as Evaluator

Recent work shows that the LLMs can be used as evaluators for various NLP tasks, and present a high correlation with human preference (Li et al., 2023). Here we use GPT3.5-turbo for the automatic evaluation. We feed generated utterances from our trained Mistral-7B-Instruct-v0.2, and compare the vanilla model with our model upon normalization and context aggregation, by predicting which response is better. We sampled 30 utterances for evaluation, and the winning rate of the enhanced model is 0.80, demonstrating the effectiveness of our proposed methods.

### 6.3.3 Ablation Study on Context Aggregation

We conduct an ablation study on the context aggregation of topic-focused summarization. In our preliminary experiment, we observe that the first three utterances from the nurse of each topically-coherent dialogue segment show more dependency on the historical context, due to the explicit topic shift (e.g., from symptom checking to customized

coaching of insulin). Therefore, at the inference stage, we collect the first-3 generated utterances of each topic, and compared models with and without adding the aggregated summaries. As shown in Table 3, all evaluation scores drop significantly when the historical summaries are removed, for all tested models (e.g., LLaMA, Mistral). This demonstrates that nurse dynamically change their topic during the conversation, the topic-specified questions in certain topics depend on the information they collect from the patient.

### 6.3.4 Ablation Study on Patient Profile

We conduct an additional ablation study on the patient profile. Following the previous step, at the inference stage, we still collect the first-3 generated utterances of each topic, and compared models with and without adding the patient profile information. As shown in Table 3, the generation performance for all tested models (e.g., LLaMA, Mistral) drops significantly when no profile is provided. For instance, in the topic *'Diet'* we observe that models tend to generate common questions (e.g., *"how is your diet?"*) when there is no profile and dialogue context. In comparison, models can ask more targeted questions (e.g., *"What about your sugar intake? Do you consume sweetened beverages?"*), which are more informative, especially at the beginning of each topic segment.

## 7 Conclusion

In this work, we investigated the feasibility and effectiveness of leveraging language language models for personalized medical dialogue generation. We conducted a case study on healthcare follow-up calls for diabetes management. Inspired by real-world conversations, we built a data preparation and refinement pipeline for spoken conversation processing, user profile construction, and proposed topic-focused summarization to distill and aggregate the historical context. To exploit the potential of LLMs, we applied efficient model training methods for domain adaptation. Our experimental results showed that context aggregation via topic-focused summarization is beneficial for long-context modeling and coherent generation.

## Limitations

The data and model used in this work are in English, thus to apply the approach to other languages, it will require training data on the specified language or using multilingual language backbones. While our proposed methods are general, when adopt them to other conversational data, in-domain annotation is required to obtain reliable results. Moreover, the hallucination made by large language models is an open problem, and the system generations in clinical scenarios still need human verification and intervention if necessary.

## Ethics and Impact Statement

We acknowledge that all of the co-authors of this work are aware of the provided ACL Code of Ethics and honor the code of conduct. The in-domain samples used in this work are fully anonymized. The original data are collected under consent for academic research purposes. Our proposed framework and methodology in general do not create a direct medical implication, and are intended to be used to improve the model accuracy and robustness for downstream applications.

## References

Mahyar Abbasian, Elahe Khatibi, Iman Azimi, David Oniani, Zahra Shakeri Hossein Abad, Alexander Thieme, Zhongqi Yang, Yanshan Wang, Bryant Lin, Olivier Gevaert, et al. 2023. Foundation metrics: Quantifying effectiveness of healthcare conversations powered by generative ai. *arXiv preprint arXiv:2309.12444*.

Rami Al-Rfou, Marc Pickett, Javier Snaider, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2016. Con-

versational contextual cues: The case of personalization and history for response ranking. *arXiv preprint arXiv:1606.00372*.

Tim Althoff, Kevin Clark, and Jure Leskovec. 2016. Large-scale analysis of counseling conversations: An application of natural language processing to mental health. *Transactions of the Association for Computational Linguistics*, 4:463–476.

Niels Ole Bernsen, Hans Dybkjær, and Laila Dybkjær. 1996. Cooperativity in human-machine and human-human spoken dialogue. *Discourse processes*, 21(2):213–236.

Jie Cao, Michael Tanana, Zac Imel, Eric Poitras, David Atkins, and Vivek Srikumar. 2019. Observing dialogue in therapy: Categorizing and forecasting behavioral codes. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5599–5611, Florence, Italy. Association for Computational Linguistics.

Zhangming Chan, Juntao Li, Xiaopeng Yang, Xiuying Chen, Wenpeng Hu, Dongyan Zhao, and Rui Yan. 2019. Modeling personalization in continuous space for response generation via augmented wasserstein autoencoders. 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 1931–1940.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910.

Shameek Ghosh, Sammi Bhatia, and Abhi Bhatia. 2018. Quro: facilitating user symptom check using a personalised chatbot-oriented dialogue system. *Stud Health Technol Inform*, 252:51–56.

Itika Gupta, Barbara Di Eugenio, Brian Ziebart, Aiswarya Baiju, Bing Liu, Ben Gerber, Lisa Sharp, Nadia Nabulsi, and Mary Smart. 2020. Human-human health coaching via text messages: Corpus, annotation, and analysis. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 246–256.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.

Jennifer Hill, W Randolph Ford, and Ingrid G Farreras. 2015. Real conversations with artificial intelligence: A comparison between human–human online conversations and human–chatbot conversations. *Computers in human behavior*, 49:245–250.

Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

M Sue Kirkman, Morris Weinberger, Pamela B Landsman, Gregory P Samsa, E Anne Shortliffe, David L Simel, and John R Feussner. 1994. A telephone-delivered intervention for patients with niddm: effect on coronary risk factors. *Diabetes care*, 17(8):840–846.

Kirsi Kivelä, Satu Elo, Helvi Kyngäs, and Maria Kääriäinen. 2014. The effects of health coaching on adult patients with chronic diseases: a systematic review. *Patient education and counseling*, 97(2):147–157.

Margaret L Lawson, Nini Cohen, Christine Richardson, Elaine Orrbine, and Ba' Pham. 2005. A randomized trial of regular standardized telephone contact by a diabetes nurse educator in adolescents with poor diabetes control. *Pediatric Diabetes*, 6(1):32–40.

Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and William B Dolan. 2016. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003.

Minzhi Li, Taiwei Shi, Caleb Ziems, Min-Yen Kan, Nancy Chen, Zhengyuan Liu, and Diyi Yang. 2023. Coannotating: Uncertainty-guided work allocation between human and large language models for data annotation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1487–1505.

Lizi Liao, Grace Hui Yang, and Chirag Shah. 2023. Proactive conversational agents in the post-chatgpt world. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3452–3455.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Shuai Lin, Pan Zhou, Xiaodan Liang, Jianheng Tang, Ruihui Zhao, Ziliang Chen, and Liang Lin. 2021. Graph-evolving meta-learning for low-resource medical dialogue generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 13362–13370.

Wenge Liu, Jianheng Tang, Yi Cheng, Wenjie Li, Yefeng Zheng, and Xiaodan Liang. 2022a. Meddg: an entity-centric medical consultation dataset for entity-aware medical dialogue generation. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 447–459. Springer.

Zhengyuan Liu, Angela Ng, Sheldon Lee, Ai Ti Aw, and Nancy F Chen. 2019. Topic-aware pointer-generator networks for summarizing spoken conversations. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 814–821. IEEE.

Zhengyuan Liu, Shikang Ni, Aiti Aw, and Nancy Chen. 2022b. Singlish message paraphrasing: A joint task of creole translation and text normalization. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3924–3936.

Zhengyuan Liu, Siti Umairah Md Salleh, Hong Choon Oh, Pavitra Krishnaswamy, and Nancy Chen. 2023. Joint dialogue topic segmentation and categorization: A case study on clinical spoken conversations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 185–193.

Zhengyi Ma, Zhicheng Dou, Yutao Zhu, Hanxun Zhong, and Ji-Rong Wen. 2021. One chatbot per person: Creating personalized chatbots based on implicit user profiles. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*, pages 555–564.

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.

Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. Training millions of personalized dialogue agents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2775–2779, Brussels, Belgium. Association for Computational Linguistics.

Oluwatobi Olabiyi, Anish Khazane, Alan Salimov, and Erik Mueller. 2019. An adversarial learning framework for a persona-based multi-turn dialogue model. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*,

pages 1–10, Minneapolis, Minnesota. Association for Computational Linguistics.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Adam Palanica, Peter Flaschner, Anirudh Thommandram, Michael Li, and Yan Fossat. 2019. Physicians' perceptions of chatbots in health care: cross-sectional web-based survey. *Journal of medical Internet research*, 21(4):e12887.

John D Piette, Morris Weinberger, Frederic B Kraemer, and Stephen J McPhee. 2001. Impact of automated calls with nurse follow-up on diabetes treatment outcomes in a department of veterans affairs health care system: a randomized controlled trial. *Diabetes care*, 24(2):202–208.

John D Piette, Morris Weinberger, Stephen J McPhee, Connie A Mah, Fredric B Kraemer, and Lawrence M Crapo. 2000. Do automated calls with nurse follow-up improve self-care and glycemic control among vulnerable patients with diabetes? *The American journal of medicine*, 108(1):20–27.

Qiao Qian, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Assigning personality/profile to a chatting machine for coherent conversation generation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 4279–4285.

Harvey Sacks, Emanuel A Schegloff, and Gail Jefferson. 1978. A simplest systematics for the organization of turn taking for conversation. In *Studies in the organization of conversational interaction*, pages 7–55. Elsevier.

Benjamin Schloss and Sandeep Konam. 2020. Towards an automated soap note: classifying utterances from medical conversations. In *Machine Learning for Healthcare Conference*, pages 610–631. PMLR.

Siqi Shen, Charles Welch, Rada Mihalcea, and Verónica Pérez-Rosas. 2020. Counseling-style reflection generation using generative pretrained transformers with augmented context. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 10–20.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. *Nature*, pages 1–9.

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.

C Barr Taylor, Nancy Houston Miller, Kelly R Reilly, George Greenwald, Darby Cunning, Allison Deeter, and Liana Abascal. 2003. Evaluation of a nurse-care management system to improve outcomes in patients with complicated diabetes. *Diabetes care*, 26(4):1058–1063.

Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.

Yuanhe Tian, Weicheng Ma, Fei Xia, and Yan Song. 2019. ChiMed: A Chinese medical corpus for question answering. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 250–260, Florence, Italy. 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, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(11).

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.

Zhongyu Wei, Qianlong Liu, Baolin Peng, Huaixiao Tou, Ting Chen, Xuan-Jing Huang, Kam-Fai Wong, and Xiang Dai. 2018. Task-oriented dialogue system for automatic diagnosis. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 201–207.

Yuan Xia, Jingbo Zhou, Zhenhui Shi, Chao Lu, and Haifeng Huang. 2020. Generative adversarial regularized mutual information policy gradient framework for automatic diagnosis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 1062–1069.

Lin Xu, Qixian Zhou, Ke Gong, Xiaodan Liang, Jianheng Tang, and Liang Lin. 2019. End-to-end knowledge-routed relational dialogue system for automatic diagnosis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7346–7353.

Lu Xu, Leslie Sanders, Kay Li, James CL Chow, et al. 2021. Chatbot for health care and oncology applications using artificial intelligence and machine learning: systematic review. *JMIR cancer*, 7(4):e27850.

Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, Hongchao Fang, Penghui Zhu, Shu Chen, and Pengtao Xie. 2020. MedDialog: Large-scale medical dialogue datasets. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9241–9250, Online. Association for Computational Linguistics.

Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and William B Dolan. 2020. Dialogpt: Large-scale generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278.

Meng Zhou, Zechen Li, Bowen Tan, Guangtao Zeng, Wenmian Yang, Xuehai He, Zeqian Ju, Subrato Chakravorty, Shu Chen, Xingyi Yang, Yichen Zhang, Qingyang Wu, Zhou Yu, Kun Xu, Eric Xing, and Pengtao Xie. 2021. On the generation of medical dialogs for COVID-19. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 886–896, Online. Association for Computational Linguistics.

Yue Zhou, Barbara Di Eugenio, Brian Ziebart, Lisa Sharp, Bing Liu, Ben Gerber, Nikolaos Agadakos, and Shweta Yadav. 2022. Towards enhancing health coaching dialogue in low-resource settings. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 694–706, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

# Evaluating Lexicon Incorporation for Depression Symptom Estimation

**Kirill Milintsevich**[1,2] and **Gaël Dias**[1] and **Kairit Sirts**[2]

[1]Normandie Univ, UNICAEN, ENSICAEN, CNRS, GREYC, France
[2]Institute of Computer Science, University of Tartu, Estonia
{first_name}.{last_name}@{unicaen.fr[1]|ut.ee[2]}

## Abstract

This paper explores the impact of incorporating sentiment, emotion, and domain-specific lexicons into a transformer-based model for depression symptom estimation. Lexicon information is added by marking the words in the input transcripts of patient-therapist conversations as well as in social media posts. Overall results show that the introduction of external knowledge within pre-trained language models can be beneficial for prediction performance, while different lexicons show distinct behaviours depending on the targeted task. Additionally, new state-of-the-art results are obtained for the estimation of depression level over patient-therapist interviews.

## 1 Introduction

Considerable interest has emerged in using natural language processing to unobtrusively infer one's mental health condition (Chancellor and De Choudhury, 2020). A majority of studies have focused on predicting major depressive disorder (MDD) either as a symptom-based estimation (Yadav et al., 2020; Milintsevich et al., 2023) or a binary classification problem (Burdisso et al., 2023; Xezonaki et al., 2020). Both clinically motivated research initiatives and social media studies have emerged. In the latter case, Twitter (Zhang et al., 2023a), Reddit (Gupta et al., 2022) and depression-related forums (Yao et al., 2021) have fostered attention. In the former case, recorded patient-therapist conversations are transcribed and associated with self-assessment depression questionnaires, such as PHQ-8 (Kroenke et al., 2009) or BDI (Beck et al., 1988).

The DAIC-WOZ dataset (Gratch et al., 2014) has mostly been studied within the context of clinical research. Different works have been proposed to automatically infer depression level on this dataset: multi-modal (Qureshi et al., 2019; Wei et al., 2022)

**Illustration of the lexicon-based input marking**

a) i'm pretty much good because see by me being a bus operator you run into circumstances and situations you gotta remain calm and still remain professional at the same time

b) i'm @ pretty @ much @ good @ because see by me being a bus operator you run into circumstances and situations you gotta remain @ calm @ and still remain professional at the same time

c) i'm @ pretty @ much @ good @ because see by me being a bus operator you run into circumstances and situations you gotta remain @ calm @ and still remain @ professional @ at the same @ time @

Table 1: Example of input marking. Text a) is the original text without markings, b) and c) show text with terms from AFINN and NRC lexicons.

and text-based architectures (Li et al., 2023; Agarwal et al., 2022). The PRIMATE dataset (Gupta et al., 2022) has also received recent attention within the context of early symptom prediction on social media posts. The most comprehensive work on this dataset is proposed by Zhang et al. (2023a), which defines a context- and PHQ-aware transformer-based architecture.

People with MDD have shown increased use of negative emotional words and decreased use of positive emotional words (Rude et al., 2004; Savekar et al., 2023). In this line, Xezonaki et al. (2020) and Qureshi et al. (2020) used feature-level and task fusion of emotion and sentiment knowledge and showed improved performance for depression estimation. However, these works, along with other studies on social media mental health data (Zhang et al., 2023b), have used pre-transformer era neural architectures. Recent state-of-the-art approaches that rely on transformer-based pre-trained language models (PLMs) have not explored external knowledge fusion (Milintsevich et al., 2023).

In this paper, we investigate whether pre-trained language models could benefit from

322

| Lexicon | PHQ-8 | Train | Dev | Test |
|---------|-------|-------|-----|------|
| AFINN | $\geq 10$ | 8.4 | 7.6 | 8.0 |
| | $< 10$ | 8.2 | 7.6 | 7.9 |
| NRC | $\geq 10$ | 7.6 | †6.8 | †7.1 |
| | $< 10$ | 7.7 | †7.6 | †7.6 |
| SDD | $\geq 10$ | †0.6 | 0.4 | 0.5 |
| | $< 10$ | †0.4 | 0.3 | 0.4 |

Table 2: Proportion of marked words for each lexicon over the DAIC-WOZ. Reported values are in percentage. † shows if the difference between the depressed and non-depressed populations is statistically significant.

the introduction of emotional, sentimental, and domain-specific external knowledge from the lexicons: AFINN (Nielsen, 2011), NRC (Mohammad and Turney, 2013) and SDD (Yazdavar et al., 2017). Introducing this external knowledge into a transformer-based model is feature-level and is achieved by modifying the input with specific markers that highlight spans of text, as shown in Table 1, inspired by the works of Wang et al. (2021) and Zhou and Chen (2022). This approach does not require any modification to the model's architecture, such as changing attention mechanism (Li et al., 2021; Wang et al., 2022) or adding new layers (Bai et al., 2022); it also keeps the model's vocabulary unchanged unlike Zhong and Chen (2021).

Results on the DAIC-WOZ dataset show that the performance of transformer-based models is impacted by the added lexicon information (especially sentiment), and new state-of-the-art values can be obtained from the combination of the three lexicons. However, such results are less expressive for the PRIMATE dataset, with slight improvements induced by the introduction of external information. Overall, the improvement in predicting particular symptoms evidences that lexicon information can be helpful, provided that its content closely corresponds to the targeted task.

## 2 Methodology

**Data.** In this work, we use two depression datasets: DAIC-WOZ (Gratch et al., 2014) and PRIMATE (Gupta et al., 2022). The DAIC-WOZ dataset contains 189 clinical interviews in a dialogue format. Each interview has two actors: a human-controlled virtual therapist and a participant. The dataset is distributed in pre-determined splits, such that 107 interviews are used for training, 35 for validation, and 47 for testing. Each interview



Figure 1: Overview of the model architecture. $U_i^N$ stands for $i$-th utterance of $N$-th input. *Symptom Scores* are $||L||$ real numbers, where $||L||$ is the number of symptoms to predict.

in the dataset is accompanied with a PHQ-8 assessment, which consists of eight questions inquiring about symptoms. Each question is scored from 0 to 3 on a Likert scale, and the total PHQ score ranging from 0 to 24 is the sum of the eight symptom scores. According to a standard cutoff score of 10, the interviews can be divided into diagnostic classes, where subjects with PHQ-8 total score $< 10$ are considered non-depressed, and those with score $\geq 10$ are categorized as depressed. The eight listed symptoms are: LOI (lack of interest), DEP (feeling down), SLE (sleeping disorder), ENE (lack of energy), EAT (eating disorder), LSE (low self-esteem), CON (concentration problem), MOV (hyper/lower activity).

The PRIMATE dataset is based on Reddit posts from depression-related communities, or subreddits, in which people describe their health conditions. A total of 2003 posts were manually annotated with binary labels for each individual symptom from the PHQ-9 (Kroenke et al., 2001), each label signifying whether the corresponding symptom is discussed in the post or not. PHQ-9 has the same first eight symptoms as PHQ-8 and one additional SUI (suicidal thoughts). The data was labeled by five crowd workers and verified by a mental health professional. The dataset is not presplit into the train, validation, and test sets, so we randomly take 1601, 201, and 201 posts for each split accordingly.

**Model architecture.** To encode the interview transcripts, we adopt the hierarchical model from (Milintsevich et al., 2023). In their model, the interview is first split utterance-by-utterance, with each utterance processed by a word-level encoder.

| Model | LOI | DEP | SLE | ENE | EAT | LSE | CON | MOV | PHQ-8 |
|---|---|---|---|---|---|---|---|---|---|
| BERT | $0.56_{\pm.05}$ | $\mathbf{0.63}_{\pm.02}$ | $0.77_{\pm.05}$ | $0.87_{\pm.04}$ | $\mathbf{0.81}_{\pm.03}$ | $0.78_{\pm.06}$ | $0.74_{\pm.01}$ | $0.34_{\pm.01}$ | $4.38_{\pm.21}$ |
| +SDD | $0.70_{\pm.02}$ | $0.88_{\pm.05}$ | $0.94_{\pm.05}$ | $0.94_{\pm.04}$ | $1.00_{\pm.07}$ | $0.97_{\pm.04}$ | $0.87_{\pm.02}$ | $0.34_{\pm.00}$ | $5.60_{\pm.18}$ |
| +AFINN | $\mathbf{0.50}_{\pm.03}$ | $0.70_{\pm.03}$ | $0.79_{\pm.03}$ | $0.81_{\pm.04}$ | $0.85_{\pm.03}$ | $0.72_{\pm.02}$ | $0.77_{\pm.02}$ | $0.34_{\pm.00}$ | $4.56_{\pm.22}$ |
| +NRC | $\mathbf{0.50}_{\pm.03}$ | $0.66_{\pm.05}$ | $\mathbf{0.73}_{\pm.05}$ | $0.77_{\pm.03}$ | $0.81_{\pm.05}$ | $0.71_{\pm.07}$ | $\mathbf{0.73}_{\pm.05}$ | $0.34_{\pm.00}$ | $\mathbf{4.31}_{\pm.18}$ |
| +ALL | $\mathbf{0.50}_{\pm.04}$ | $0.69_{\pm.03}$ | $0.81_{\pm.12}$ | $\mathbf{0.74}_{\pm.06}$ | $0.81_{\pm.07}$ | $0.69_{\pm.05}$ | $0.74_{\pm.03}$ | $0.34_{\pm.00}$ | $4.56_{\pm.42}$ |
| MeBERT | $0.59_{\pm.02}$ | $0.64_{\pm.06}$ | $0.91_{\pm.05}$ | $0.92_{\pm.04}$ | $0.89_{\pm.04}$ | $0.71_{\pm.02}$ | $0.71_{\pm.04}$ | $0.35_{\pm.01}$ | $4.71_{\pm.23}$ |
| +SDD | $0.69_{\pm.07}$ | $0.72_{\pm.08}$ | $0.89_{\pm.07}$ | $0.92_{\pm.02}$ | $0.93_{\pm.07}$ | $0.85_{\pm.07}$ | $0.78_{\pm.06}$ | $0.34_{\pm.00}$ | $5.07_{\pm.38}$ |
| +AFINN | $0.48_{\pm.04}$ | $0.62_{\pm.02}$ | $0.71_{\pm.05}$ | $0.78_{\pm.04}$ | $0.79_{\pm.03}$ | $0.70_{\pm.03}$ | $0.74_{\pm.03}$ | $0.34_{\pm.00}$ | $4.27_{\pm.22}$ |
| +NRC | $0.60_{\pm.05}$ | $0.68_{\pm.03}$ | $0.71_{\pm.05}$ | $0.78_{\pm.04}$ | $0.80_{\pm.08}$ | $0.74_{\pm.02}$ | $0.71_{\pm.05}$ | $0.34_{\pm.00}$ | $4.35_{\pm.26}$ |
| +ALL | $\mathbf{0.44}_{\pm.06}$ | $\mathbf{0.55}_{\pm.04}$ | $\mathbf{0.63}_{\pm.06}$ | $\mathbf{0.72}_{\pm.07}$ | $\mathbf{0.69}_{\pm.03}$ | $\mathbf{0.67}_{\pm.04}$ | $\mathbf{0.67}_{\pm.03}$ | $0.34_{\pm.00}$ | $\mathbf{3.59}_{\pm.31}$ |
| SOTA | $0.53_{\pm.05}$ | $\mathbf{0.55}_{\pm.03}$ | $0.75_{\pm.07}$ | $\mathbf{0.64}_{\pm.03}$ | $0.81_{\pm.05}$ | $\mathbf{0.62}_{\pm.02}$ | $0.83_{\pm.04}$ | $0.44_{\pm.02}$ | $3.78_{\pm.13}$ |

Table 3: Results for the DAIC-WOZ test set. The mean MAE and standard deviation are reported for five runs. The best MAE for each symptom is **in bold**. SOTA means current state-of-the-art results in the literature (Milintsevich et al., 2023).

All utterance representations are then concatenated into one sequence, later processed by an utterance-level encoder. In the end, the classification head produces a real number in the range from 0 to 3 for each symptom. Several changes are made to the original architecture to gain training efficiency. First, the BiLSTM utterance-level encoder is replaced with a randomly initialized 4-layer 12-head transformer encoder. Second, we change the way the input data is represented. In the original model, each utterance of the interview is encoded separately by a word-level encoder. This is far from optimal since most of the utterances are short (<10 tokens), thus, a lot of computation is wasted on padding tokens. Instead, the utterances are concatenated into one input text separated by the [SEP] special token. This way, the number of passes through the encoder is reduced from the number of utterances $K$ to $\bar{K}$, defined as in Equation 1, where $|U_i|$ is the number of tokens in an utterance and $m$ is the maximum input length of the word-level encoder.

$$\bar{K} = \left\lceil \frac{\sum (|U_i| + 1)}{m} \right\rceil \quad (1)$$

In practice, it reduces the number of word-level encoder passes by $\sim 40$ times for each input. After, we perform the *Mean* [SEP] *pooling* on the tokens representing each utterance to get the final utterance representation. The overview of the model architecture is presented in Figure 1.

**Lexicons.** To incorporate the external knowledge into the model, we use three lexicons: AFINN (Nielsen, 2011), NRC (Mohammad and Turney, 2013), and SDD (Yazdavar et al., 2017).

AFINN is a sentiment lexicon that includes a list of 2,477 terms manually rated for the sentiment valence with a value between $-5$ (negative) and $+5$ (positive). Nielsen (2011) used Twitter postings together with different word lists as a source for the lexicon. NRC is a word-emotion association lexicon that is a list of 14,182 words and their associations with eight basic emotions (anger, fear, anticipation, trust, surprise, sadness, joy, and disgust) and two sentiments (negative and positive). Mohammad and Turney (2013) compiled terms from Macquarie Thesaurus (Bernard, 1986), WordNet Affect Lexicon (Strapparava and Valitutti, 2004), and General Inquirer (Stone et al., 1966) and labeled them with the help of crowd-sourced workers. SDD is a part of the Social-media Depression Detector and is a lexicon of more than 1,620 depression-related words and phrases created in collaboration with a psychologist clinician.

**Input marking.** In particular, we employ the technique proposed by Zhou and Chen (Zhou and Chen, 2022) to identify and annotate the lexicon words in the input text. It involves marking a lexicon word using the "@" token on either side (see Table 1 for examples). We chose the "@" token for marking since it is not present in the data but included in the model's vocabulary. This way, the pre-trained model's architecture remains unchanged[1]. The proportion of marked words within the DAIC-WOZ is illustrated in Table 2, where the statistical test is Student's t-test with p-value $< 0.05$.

---

[1]Typed marking strategies that include emotion and sentiment values have also been tested and provided no additional insights compared to the simple input marking.

Figure 2: Average predicted values for depressed and non-depressed patients of the DAIC-WOZ test set.

| Model | LOI | DEP | SLE | ENE | EAT | LSE | CON | MOV | SUI |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| BERT | $\mathbf{0.59}_{\pm.03}$ | $\mathbf{0.65}_{\pm.03}$ | $0.81_{\pm.01}$ | $0.62_{\pm.02}$ | $0.75_{\pm.06}$ | $0.60_{\pm.02}$ | $\mathbf{0.65}_{\pm.01}$ | $0.81_{\pm.01}$ | $0.82_{\pm.01}$ |
| +SDD | $0.58_{\pm.03}$ | $0.62_{\pm.02}$ | $0.81_{\pm.01}$ | $\mathbf{0.64}_{\pm.03}$ | $0.74_{\pm.03}$ | $\mathbf{0.63}_{\pm.03}$ | $0.63_{\pm.03}$ | $\mathbf{0.82}_{\pm.02}$ | $0.82_{\pm.01}$ |
| +AFINN | $0.57_{\pm.03}$ | $0.60_{\pm.03}$ | $0.80_{\pm.02}$ | $0.62_{\pm.02}$ | $0.76_{\pm.02}$ | $0.59_{\pm.03}$ | $0.64_{\pm.01}$ | $0.81_{\pm.02}$ | $\mathbf{0.83}_{\pm.01}$ |
| +NRC | $0.55_{\pm.04}$ | $0.62_{\pm.04}$ | $\mathbf{0.82}_{\pm.01}$ | $0.60_{\pm.02}$ | $0.79_{\pm.04}$ | $0.59_{\pm.03}$ | $0.61_{\pm.04}$ | $0.80_{\pm.01}$ | $0.82_{\pm.02}$ |
| +ALL | $0.56_{\pm.05}$ | $0.63_{\pm.02}$ | $0.79_{\pm.02}$ | $0.61_{\pm.02}$ | $\mathbf{0.80}_{\pm.02}$ | $0.58_{\pm.03}$ | $0.61_{\pm.01}$ | $\mathbf{0.82}_{\pm.01}$ | $0.82_{\pm.02}$ |
| MEBERT | $\mathbf{0.58}_{\pm.03}$ | $0.58_{\pm.02}$ | $0.82_{\pm.02}$ | $0.62_{\pm.01}$ | $0.78_{\pm.03}$ | $0.60_{\pm.04}$ | $0.62_{\pm.03}$ | $\mathbf{0.82}_{\pm.01}$ | $0.84_{\pm.01}$ |
| +SDD | $0.53_{\pm.04}$ | $\mathbf{0.60}_{\pm.02}$ | $\mathbf{0.83}_{\pm.01}$ | $0.62_{\pm.02}$ | $0.79_{\pm.01}$ | $0.60_{\pm.02}$ | $0.61_{\pm.03}$ | $0.81_{\pm.02}$ | $\mathbf{0.86}_{\pm.01}$ |
| +AFINN | $0.57_{\pm.03}$ | $0.55_{\pm.04}$ | $\mathbf{0.83}_{\pm.01}$ | $0.62_{\pm.02}$ | $0.79_{\pm.01}$ | $\mathbf{0.63}_{\pm.02}$ | $0.58_{\pm.02}$ | $0.81_{\pm.02}$ | $0.85_{\pm.02}$ |
| +NRC | $0.57_{\pm.03}$ | $0.58_{\pm.03}$ | $0.82_{\pm.02}$ | $\mathbf{0.63}_{\pm.03}$ | $0.79_{\pm.02}$ | $\mathbf{0.63}_{\pm.01}$ | $0.61_{\pm.03}$ | $0.80_{\pm.02}$ | $0.85_{\pm.01}$ |
| +ALL | $0.56_{\pm.03}$ | $0.59_{\pm.04}$ | $0.80_{\pm.02}$ | $0.62_{\pm.02}$ | $\mathbf{0.80}_{\pm.02}$ | $0.61_{\pm.01}$ | $\mathbf{0.63}_{\pm.02}$ | $\mathbf{0.82}_{\pm.02}$ | $0.84_{\pm.01}$ |

Table 4: Results for the PRIMATE test set. The mean macro-F1 score is reported for five runs. The best macro-F1 for each symptom is **in bold**. As standard splits are not provided, we cannot present SOTA results. As standard splits are not provided, we cannot present SOTA results.

**Experimental setup.** We used two pre-trained models in the word-level encoder of our architecture: BERT-Base model (Devlin et al., 2018) and MentalBERT (Ji et al., 2022). We refer to them as **BERT** and **MeBERT** further on. Both models share the same architecture; however, BERT was pre-trained on general domain data, while MeBERT used mental health-related data, mostly based on Reddit. Each model is finetuned with the same hyperparameters (mostly following Mosbach et al., 2020) and different input markings. For example, the BERT+SDD model uses BERT as a pre-trained model and SDD lexicon for input marking. +ALL models use a union of all three lexicons. All models are trained with a mini-batch size of 16, Py-Torch realization of AdamW optimizer (Loshchilov and Hutter, 2017) with a learning rate of $2 \cdot 10^{-5}$ and linear scheduler with a warm-up ratio of 0.1. For the word-level PLMs, only their attention layers are finetuned. The utterance-level encoder is randomly initialized based on the transformer encoder architecture with the following hyperparameters: 4 layers, 12 attention heads, hidden dimensions of encoder and pooler layers of 768, intermediate hidden dimension of 1536. The rest of the

hyperparameters follow the default BertConfig from the HuggingFace Transformers library (Wolf et al., 2020). For the DAIC-WOZ dataset, results are evaluated with micro-averaged mean absolute error (MAE). Symptom-based errors are calculated for each symptom individually. PHQ-8 score is obtained by summing the eight symptom scores, and MAE for PHQ-8 is calculated on this summation. We evaluate results on the PRIMATE dataset with a macro-averaged F1 score.

## 3 Results and Discussion

Table 3 shows the results for the DAIC-WOZ test set. For the BERT model, the lexicon-based input marking brings slight overall improvement when AFINN or NRC lexicons are introduced. Most notably, the NRC input marking shows improved or equal MAE for all symptom scores except DEP. The combination of all lexicons is marginally beneficial overall, and results have deteriorated with the exclusive introduction of the SDD lexicon. On the other hand, for the MeBERT model, the combination of all the lexicons produces the best results overall, both symptom-wise and for the global PHQ-8 score. Furthermore, both AFINN and NRC

lexicons improve the prediction for the MeBERT model, similar to the BERT model. Also, when only the SDD lexicon is used for input marking, the model shows worse performance than the baseline setting.

Figure 2 depicts a more detailed overview of the best-performing models: BERT+NRC and MeBERT+ALL. Additionally, we finetune the +Rand version of both BERT and MeBERT to verify if the improvement comes only from the input marking by randomly marking 8% of the words in each interview. From the results, the improvement for the BERT+NRC model comes from the non-depressed population. MeBERT+All model, however, improves for both depressed and non-depressed populations and is less sensitive to the marking bias. Interestingly, +Rand models show some improvement for the non-depressed population, suggesting that input markings alone act as a regularizer.

Table 4 shows the results for the PRIMATE test set. Contrary to the results from Table 3, introducing external knowledge does not clearly improve performances. The models that use the lexicon input marking show signs of improvement for some symptoms, but it is largely inconsistent. Unlike for the DAIC-WOZ, the SDD-based input marking provides the best F1 score for three symptoms, both for BERT and MentalBERT models, while the benefits of AFINN and NRC are limited or absent and spread over symptoms.

The results from the DAIC-WOZ show that PLMs can indeed benefit from the introduction of external knowledge about the sentiment and emotional value of the words. Surprisingly, the introduction of the depression-specific lexicon had the opposite effect. We hypothesize that two reasons could cause it. First, as seen in Table 2, SDD covers less than 0.5% of words in the interview, almost 15 times less than AFINN and NRC. Thus, the introduced signal might be too weak for the model to learn. Second, the SDD lexicon was based on Twitter data, while DAIC-WOZ contains transcripts of real conversations. From our observations, the people describe their problems more explicitly in their social media posts. At the same time, DAIC-WOZ conversations are more generally themed, and the PHQ-8 scores are based on the person's self-assessment test rather than the conversations themselves. This brings us back to the conceptual difference between the DAIC-WOZ and PRIMATE datasets. While the first one aims at establishing the link between the underlying person's mental condition and their speech, the latter one sets a goal of detecting whether a particular symptom is mentioned in the text. In addition, the PRIMATE dataset is annotated by layman crowd workers, and the labels are not consistent and contain inevitable mistakes (Milintsevich et al., 2024). This might explain the reason behind the greater impact of the AFINN and NRC lexicons for modeling the DAIC-WOZ dataset.

# 4 Conclusion

This paper targets lexicon incorporation in transformer-based models for symptom-based depression estimation. The external information is supplied through a marking strategy, which avoids any modification to the model's architecture. The set of endeavoured experiments shows that introducing sentimental, emotional and/or domain-specific lexicons can correlate with overall performance improvement if adapted to the targeted task[2].

## Limitations

The main limitation in automated clinical mental health assessment with natural language processing is the difficulty of acquiring and accessing large quantities of data. DAIC-WOZ and PRIMATE are rare exceptions as it is publicly available and clinically verified. However, DAIC-WOZ, in particular, suffers from a small number of data points that makes it hard to train and validate hypotheses, as both validation and test sets are particularly small. As a consequence, this piece of research requires further validation on a larger body of clinical data.

## Ethical Considerations

We acknowledge the potential ethical aspects of the work that studies the methods to unobtrusively detect someone's mental health status. Here, we are using publicly available datasets collected for research purposes. Also, the lexicons we use are publicly available and have not been composed based on private confidential material. If such a system that could predict the presence of depression symptoms based on actual clinical interviews would be deployed in practice, it would require the informed consent of all participants involved

---

[2]Source code is available here: https://github.com/501Good/dialogue-classifier.

as well as the understanding of the validity boundaries of such systems, meaning that the predictions of such systems cannot replace the assessment of trained clinicians, but rather assist them in their activities.

## Acknowledgements

## References

Navneet Agarwal, Gaël Dias, and Sonia Dollfus. 2022. Agent-based splitting of patient-therapist interviews for depression estimation. In *PAI4MH @ 36th Conference on Neural Information Processing Systems (NeurIPS)*, New Orleans, USA.

Jiangang Bai, Yujing Wang, Hong Sun, Ruonan Wu, Tianmeng Yang, Pengfei Tang, Defu Cao, Mingliang Zhang1, Yunhai Tong, Yaming Yang, Jing Bai, Ruofei Zhang, Hao Sun, and Wei Shen. 2022. Enhancing self-attention with knowledge-assisted attention maps. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 107–115, Seattle, United States. Association for Computational Linguistics.

Aaron T Beck, Robert A Steer, and Margery G Carbin. 1988. Psychometric properties of the beck depression inventory: Twenty-five years of evaluation. *Clinical psychology review*, 8(1):77–100.

J.R.L. Bernard. 1986. *The MacQuarrie Thesaurus: The Book of Words*. Macquarie Library.

Sergio Burdisso, Esaú Villatoro-Tello, Srikanth Madikeri, and Petr Motlicek. 2023. Node-weighted graph convolutional network for depression detection in transcribed clinical interviews. In *INTERSPEECH*.

Stevie Chancellor and Munmun De Choudhury. 2020. Methods in predictive techniques for mental health status on social media: a critical review. *NPJ digital medicine*, 3(1):43.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Stratou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-Philippe Morency. 2014. The distress analysis interview corpus of human and computer interviews. In *Ninth International Conference on Language Resources and Evaluation (LREC)*, pages 3123–3128, Reykjavik, Iceland. European Language Resources Association (ELRA).

Shrey Gupta, Anmol Agarwal, Manas Gaur, Kaushik Roy, Vignesh Narayanan, Ponnurangam Kumaraguru, and Amit Sheth. 2022. Learning to automate follow-up question generation using process knowledge for depression triage on Reddit posts. In *Eighth Workshop on Computational Linguistics and Clinical Psychology (CLPSY)*, pages 137–147, Seattle, USA. Association for Computational Linguistics.

Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria. 2022. MentalBERT: Publicly Available Pretrained Language Models for Mental Healthcare. In *Proceedings of LREC*.

Kurt Kroenke, Robert L Spitzer, and Janet BW Williams. 2001. The phq-9: validity of a brief depression severity measure. *Journal of general internal medicine*, 16(9):606–613.

Kurt Kroenke, Tara W. Strine, Robert L. Spitzer, Janet B.W. Williams, Joyce T. Berry, and Ali H. Mokdad. 2009. The PHQ-8 as a measure of current depression in the general population. *Journal of Affective Disorders*, 114(1-3):163–173.

Mingzheng Li, Xiao Sun, and Meng Wang. 2023. Detecting depression with heterogeneous graph neural network in clinical interview transcript. *IEEE Transactions on Computational Social Systems*.

Zhongli Li, Qingyu Zhou, Chao Li, Ke Xu, and Yunbo Cao. 2021. Improving BERT with syntax-aware local attention. In *Findings of the Association for Computational Linguistics (ACL-IJCNLP)*, pages 645–653. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Kirill Milintsevich, Kairit Sirts, and Gaël Dias. 2023. Towards automatic text-based estimation of depression through symptom prediction. *Brain Informatics*, 10(1):1–14.

Kirill Milintsevich, Kairit Sirts, and Gaël Dias. 2024. Your model is not predicting depression well and that is why: A case study of PRIMATE dataset. In *Proceedings of the 9th Workshop on Computational Linguistics and Clinical Psychology (CLPsych 2024)*, pages 166–171, St. Julians, Malta. Association for Computational Linguistics.

Saif M. Mohammad and Peter D. Turney. 2013. Crowdsourcing a word-emotion association lexicon. *Computational Intelligence*, 29(3):436–465.

Marius Mosbach, Maksym Andriushchenko, and Dietrich Klakow. 2020. On the stability of fine-tuning bert: Misconceptions, explanations, and strong baselines. *arXiv preprint arXiv:2006.04884*.

Finn Årup Nielsen. 2011. A new anew: Evaluation of a word list for sentiment analysis in microblogs. *arXiv preprint arXiv:1103.2903*.

Syed Arbaaz Qureshi, Gael Dias, Mohammed Hasanuzzaman, and Sriparna Saha. 2020. Improving depression level estimation by concurrently learning emotion intensity. *IEEE Computational Intelligence Magazine*, 15(3):47–59.

Syed Arbaaz Qureshi, Sriparna Saha, Mohammed Hasanuzzaman, and Gaël Dias. 2019. Multitask representation learning for multimodal estimation of depression level. *IEEE Intelligent Systems*, 34(5):45–52.

Stephanie Rude, Eva-Maria Gortner, and James Pennebaker. 2004. Language use of depressed and depression-vulnerable college students. *Cognition & Emotion*, 18(8):1121–1133.

Anbu Savekar, Shashikanta Tarai, and Moksha Singh. 2023. Structural and functional markers of language signify the symptomatic effect of depression: A systematic literature review. *European Journal of Applied Linguistics*, 11(1):190–224.

Philip J Stone, Dexter C Dunphy, and Marshall S Smith. 1966. The general inquirer: A computer approach to content analysis.

Carlo Strapparava and Alessandro Valitutti. 2004. WordNet affect: an affective extension of WordNet. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal. European Language Resources Association (ELRA).

University of Tartu. 2018. UT rocket.

Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. 2021. K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters. In *Findings of the Association for Computational Linguistics (ACL-IJCNLP)*, pages 1405–1418. Association for Computational Linguistics.

Shanshan Wang, Zhumin Chen, Zhaochun Ren, Huasheng Liang, Qiang Yan, and Pengjie Ren. 2022. Paying more attention to self-attention: Improving pre-trained language models via attention guiding. *arXiv preprint arXiv:2204.02922*.

Ping-Cheng Wei, Kunyu Peng, Alina Roitberg, Kailun Yang, Jiaming Zhang, and Rainer Stiefelhagen. 2022. Multi-modal depression estimation based on sub-attentional fusion. In *European Conference on Computer Vision Workshops (ECCVW)*, pages 623–639.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 38–45, Online. Association for Computational Linguistics.

Danai Xezonaki, Georgios Paraskevopoulos, Alexandros Potamianos, and Shrikanth Narayanan. 2020. Affective Conditioning on Hierarchical Attention Networks Applied to Depression Detection from Transcribed Clinical Interviews. In *INTERSPEECH*, pages 4556–4560.

Shweta Yadav, Jainish Chauhan, Joy Prakash Sain, Krishnaprasad Thirunarayan, Amit Sheth, and Jeremiah Schumm. 2020. Identifying depressive symptoms from tweets: Figurative language enabled multitask learning framework. In *28th International Conference on Computational Linguistics (COLING)*, pages 696–709, Barcelona, Spain.

Xiaoxu Yao, Guang Yu, Jingyun Tang, and Jialing Zhang. 2021. Extracting depressive symptoms and their associations from an online depression community. *Computers in human behavior*, 120:106734.

Amir Hossein Yazdavar, Hussein S Al-Olimat, Monireh Ebrahimi, Goonmeet Bajaj, Tanvi Banerjee, Krishnaprasad Thirunarayan, Jyotishman Pathak, and Amit Sheth. 2017. Semi-supervised approach to monitoring clinical depressive symptoms in social media. In *Proceedings of the 2017 IEEE/ACM international conference on advances in social networks analysis and mining 2017*, pages 1191–1198.

Tianlin Zhang, Kailai Yang, Hassan Alhuzali, Boyang Liu, and Sophia Ananiadou. 2023a. Phq-aware depressive symptoms identification with similarity contrastive learning on social media. *Information Processing & Management*, 60(5):103417.

Tianlin Zhang, Kailai Yang, Shaoxiong Ji, and Sophia Ananiadou. 2023b. Emotion fusion for mental illness detection from social media: A survey. *Information Fusion*, 92:231–246.

Zexuan Zhong and Danqi Chen. 2021. A frustratingly easy approach for entity and relation extraction. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 50–61. Association for Computational Linguistics.

Wenxuan Zhou and Muhao Chen. 2022. An improved baseline for sentence-level relation extraction. In *2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (AACL-IJCNLP)*, pages 161–168. Association for Computational Linguistics.

# Semi-automatic Construction of a Word Complexity Lexicon for Japanese Medical Terminology

**Soichiro Sugihara**[1]    **Tomoyuki Kajiwara**[1]    **Takashi Ninomiya**[1]
**Shoko Wakamiya**[2]    **Eiji Aramaki**[2]
[1]Ehime University    {sugihara@ai., kajiwara@, ninomiya@}cs.ehime-u.ac.jp
[2]Nara Institute of Science and Technology    {wakamiya, aramaki}@is.naist.jp

## Abstract

We construct a word complexity lexicon for medical terms in Japanese. To facilitate communication between medical practitioners and patients, medical text simplification is being studied. Medical text simplification is a natural language processing task that paraphrases complex technical terms into expressions that patients can understand. However, in contrast to English, where this task is being actively studied, there are insufficient language resources in Japanese. As a first step in advancing research on medical text simplification in Japanese, we annotate the 370,000 words from a large-scale medical terminology lexicon with a five-point scale of complexity for patients.

## 1 Introduction

Communication between medical practitioners and patients is important to facilitate understanding of the diagnosis and agreement on a treatment plan (Ha and Longnecker, 2010). One of the factors that make communication difficult in the medical field is the difference in expertise between medical practitioners and patients. In particular, since many medical terms are difficult for patients to understand, medical practitioners are expected to paraphrase them into simple expressions to make them easier to understand.

To solve this problem, medical text simplification (Leroy and Endicott, 2012; Joseph et al., 2023; Yang et al., 2023) has been studied, mainly in English. However, there is a lack of available lexicons and corpora for medical text simplification in Japanese. In this study, as a first step to tackle Japanese medical text simplification, we construct a complexity lexicon for medical terms.

We first recruited 40 annotators, who were not medical practitioners via crowdsourcing to survey word complexity for 10,000 medical terms. As a

| Complexity | Medical Terminology |
|---|---|
| 1 (Simple) | めまい (Dizzy) |
| 2 | 感電死 (Electrocution) |
| 3 | 若年性脱毛症 (Premature Alopecia) |
| 4 | 後天性てんかん (Acquired Epilepsy) |
| 5 (Complex) | 掌蹠膿疱症性骨関節炎 (Pustulotic Arthro-Osteitis) |

Table 1: Examples of Japanese medical terminology.

result, we found that the number of unknown medical terms decreased with age and that men tended to be unaware of medical terms related to pregnancy and childbirth, among other characteristics observed for each of the attributes of the annotators. Furthermore, we trained a complexity estimation model for medical terms using machine learning with features such as character types, word frequencies, and word embeddings, and achieved higher performance than existing methods. Finally, as shown in Table 1, we estimated the word complexity for 370,000 disease names and symptom expressions from a large-scale medical terminology lexicon in Japanese[1] (Ito et al., 2018). Our word complexity lexicon will be available[2] upon publication of this paper.

## 2 Related Work

Large-scale word complexity lexicons in English have been constructed using two approaches. One is to estimate word complexity using the log ratio of the probability of word occurrence in the normal and simple corpora (Pavlick and Nenkova, 2015). The other is to manually annotate word complexity for a subset of the vocabulary and train a word complexity estimation model using these annotations (Pavlick and Callison-Burch, 2016;

---

[1] https://sociocom.naist.jp/manbyou-dic/
[2] https://github.com/EhimeNLP/J-MeDic-Complexity

Maddela and Xu, 2018). In Japanese, the former approach cannot be applied because of the unavailability of a large-scale corpus written in simple language. Therefore, this study takes the latter approach to construct a word complexity lexicon.

In Japanese, a domain-independent word complexity estimation model has been proposed that employs character types, word frequencies, and word embeddings as features (Kajiwara et al., 2020). For word complexity estimation specific to the medical domain, a method that takes into account the number of characters and morphemes has been proposed (Yamamoto et al., 2019). Similar to these previous studies, we train a machine learning-based word complexity estimator.

## 3 Word Complexity Annotation

### 3.1 Crowdsourcing

To train the word complexity estimation model, we asked non-medical practitioners to annotate the complexity of medical terms. These medical terms are 10,000 terms randomly selected from the top 30,000 terms with the most reliable terminology in a large-scale lexicon of disease names in Japanese[1] (Ito et al., 2018).

For diversity of annotators, eight groups were formed based on a combination of age (20s, 30s, 40s, and 50s) and gender (male and female), with five annotators per group, for a total of 40 annotators recruited. For the crowdsourcing service, we used Lancers[3] and paid the annotators 1 JPY per word (1,000 JPY per hour).

The annotators assigned each word the following a five-point scale of complexity.

1. I use this term in my daily conversation.

2. I have used this terminology.

3. I can understand what this term means.

4. I have seen or heard this term but do not know what it means.

5. I do not know what this term means and have never seen or heard of it.

To improve quality, two levels of filtering were applied to the annotators. First, we requested a small annotation of 300 words. We reviewed the responses and asked only those who had no problems to annotate the remaining 9,700 words. In

Figure 1: Distribution of complexity by age and gender.

addition, after all 10,000 words were annotated, inter-annotator agreement was calculated for each group of age and gender. Annotators with a Quadratic Weighted Kappa (QWK) (Cohen, 1968) of less than 0.3 with someone in the group were excluded and new annotators were recruited.

### 3.2 Analysis

We analyze characteristics by age and gender based on our complexity annotations. Figure 1 shows the distribution of complexity labels by age and gender. In their 20s and 30s, only about 10% of medical terms are understood. As they get older, the number of medical terms they don't know decreases. However, even in their 50s, more than 70% of medical terms cannot be understood.

Next, we observe examples of medical terms that are known above a certain age. All annotators know "しゃっくり" (hiccups) and "かぜ" (cold) used in daily conversation, while only annotators in their 40s or older or 50s know "食道ポリープ" (esophageal polyp) and "大腿骨骨折" (femur fracture) which tend to increase in patients as they get older. These imply that our complexity annotations reflect age-specific characteristics.

Finally, we observe examples of medical terms that certain groups do not know. Young men in their 30s and younger seem to be unfamiliar with some of the medical terms related to pregnancy and childbirth, such as "異常胎位" (abnormal fetal presentation) and "早発卵巣不全" (premature ovarian failure). These imply that our complexity annotations reflect gender-specific characteristics.

## 4 Word Complexity Estimation

We train a machine learning-based word complexity estimation model in addition to the three ba-

sic features used in the previous study (Yamamoto et al., 2019), with three proposed features. As in previous studies (Yamamoto et al., 2019; Kajiwara et al., 2020), we use the support vector machine (SVM) model[4] for machine learning.[5]

## 4.1 Basic Features

**Character Types** These features represent the types of characters (hiragana, katakana, kanji, numbers, and alphabetic characters) that make up a medical term. It consists of the following 15 dimensions: binary features (5 dimensions) that represent the presence or absence of each character type, integer features (5 dimensions) that represent the number of characters for each character type, and integer features (5 dimensions) that represent the maximum number of consecutive characters for each character type.

**Number of Morphemes** This is one-dimensional integer feature that represents how many morphemes a medical term is composed of. Medical terms are tokenized with MeCab[6] (IPADIC) (Kudo et al., 2004) and the number of morphemes is counted.

**Character/Morpheme Frequencies** These features are the frequencies of the letters and morphemes that make up the medical term in the corpus. Six types of frequency information are used as the features: the total, average, maximum, and minimum frequencies of morphemes in the medical term, as well as the frequency of the first morpheme and the frequency of the last morpheme. Japanese Wikipedia was used as the corpus, and MeCab was used as the morphological analyzer. Note that frequencies are used logarithmically, but as in previous study (Yamamoto et al., 2019), when the frequency is 0, 0 is used instead of log 0. These features are obtained not only in morpheme units but also in character units, for a total of a 12-dimensional real number of features.

## 4.2 Proposed Features

**PF1: Frequencies on Web Corpus** We count frequencies of characters and morphemes similar to basic features on the CC-100[7] (Conneau et al., 2020), a large-scale Web corpus. These are 12-dimensional real number of features, same as the basic features. Counting frequencies on multiple corpora is known to contribute to the word complexity estimation (Kajiwara and Komachi, 2018). However, as mentioned earlier, this study does not use the Twittr and the BCCWJ corpora used in previous study (Yamamoto et al., 2019), so a large-scale Web corpus is employed instead.

**PF2: Word Frequencies** In contrast to previous study (Yamamoto et al., 2019), we also count the frequency of medical terms in word units without segmentation. This is implemented by extending MeCab's morphological analysis with a Japanese disease lexicon[8] (Ito et al., 2018). We count word frequencies in each of the Wikipedia and CC-100 corpora, logarithmize them, and use them as two-dimensional real number features.

**PF3: Word Embeddings** We also employ word embeddings, which has been used in previous study (Kajiwara et al., 2020). We use pre-trained fastText[9] (Bojanowski et al., 2017). If a medical term consists of multiple morphemes, each of those vectors is averaged and used as a 300-dimensional real number of features.

## 5 Experiments and Results

We train and evaluate word complexity estimation models using complexity annotations for 10,000 medical terms.

## 5.1 Experiments

**Dataset** We average the complexity labels obtained from 40 annotators and round them to integers to define a five-point scale of gold complexity labels for 10,000 medical terms. Since this task is an ordinal classification, we use accuracy and QWK (Cohen, 1968) as evaluation metrics. As shown in Table 2, the training and evaluation dataset were randomly split at a ratio of 9:1 for our experiments. Since our dataset is unbalanced, we

---

[4] We also experimented with neural networks, but the SVM model achieved higher performance.

[5] As one of the features, previous study (Yamamoto et al., 2019) employed word frequencies counted on Twitter. However, we do not use this feature because changes in Twitter's API restrictions have made this counting difficult. Furthermore, word frequencies from the Balanced Corpus of Contemporary Written Japanese (BCCWJ) (Maekawa et al., 2010) are not used in this study, since previous study (Yamamoto et al., 2019) reported that these word frequencies were not effective.

[6] https://taku910.github.io/mecab/

[7] https://data.statmt.org/cc-100/
[8] https://sociocom.naist.jp/j-meddic-for-mecab/
[9] https://fasttext.cc/docs/en/crawl-vectors.html

331

| Labels | 1 | 2 | 3 | 4 | 5 | Total |
|---|---|---|---|---|---|---|
| Train | 33 | 100 | 341 | 2,650 | 5,876 | 9,000 |
| Test | 4 | 11 | 38 | 294 | 653 | 1,000 |
| Total | 37 | 111 | 379 | 2,944 | 6,529 | 10,000 |

Table 2: Number of terms per complexity.

adjusted the label ratios in both training and evaluation datasets to be equal by stratified splitting.[10]

**Model** For word complexity estimation model, a multi-class classification model was implemented using SVM (RBF kernel) in scikit-learn (1.3.2)[11] (Pedregosa et al., 2011). The hyperparameters C and gamma were selected from $\{1, 5, 10, 50, 100\}$ and $\{0.0001, 0.0005, 0.001, 0.05, 0.1\}$, respectively, and the combination with the highest QWK was selected by grid search with a five-fold cross-validation.[12] The features were standardized.[13]

**Comparative Methods** We compare the proposed method to two types of baselines. One is a simple baseline that always outputs the most frequent class, label 5. The other is a baseline that uses only the basic features of Section 4.1, which replicates the previous study (Yamamoto et al., 2019). Our method uses the proposed features of Section 4.2 in addition to the basic features.

## 5.2 Results

Table 3 shows the experimental results. Existing method using only basic features does not perform well enough, as it is equivalent in accuracy to a baseline that always outputs the most frequent labels. The proposed method significantly improved performance over these baselines by 14 points in accuracy and 28 points in QWK.

To clarify the effectiveness of each of the proposed features, an ablation analysis was performed to remove one of the proposed features from the proposed method. The fact that both accuracy and QWK decrease when any of the features are ex-

|  | Accuracy | QWK |
|---|---|---|
| Baseline | 0.653 | - |
| Basic features | 0.653 | 0.456 |
| Proposed method | **0.793** | **0.732** |
| Proposed method w/o PF1 | 0.782 | 0.729 |
| Proposed method w/o PF2 | 0.785 | 0.695 |
| Proposed method w/o PF3 | 0.718 | 0.612 |
| Only PF1 | 0.658 | 0.483 |
| Only PF2 | 0.612 | 0.444 |
| Only PF3 | 0.768 | 0.660 |

Table 3: Experimental results of word complexity estimation.

cluded shows that all of our proposed features are useful. Note that the performance decreases significantly when PF3 is excluded, suggesting that word embeddings are a particularly important feature. When each of the proposed features was used alone, PF1 alone outperformed the baselines, revealing that frequency features on a large-scale Web corpus are also useful for estimating the complexity of medical terminology.

## 6 Conclusion

In this study, we trained a word complexity estimation model based on word complexity annotations of 10,000 Japanese medical terms by 40 non-medical practitioners. Our word complexity annotations revealed that even though the number of unknown medical terms decreases with increasing age, more than 70% of medical terms are difficult to understand, even for those in their 50s. Experiments on word complexity estimation revealed that features of word frequencies and word embeddings obtained from a large-scale Web corpus are useful. Finally, we developed a word complexity estimator for Japanese medical terms that can classify five levels of complexity with about 80% accuracy, and released a word complexity lexicon[2] covering about 370,000 Japanese medical terms.

Although this study focused on disease and symptom names in Japanese, our future work includes the application of complexity estimation to more diverse medical terminology, such as drug names and names of human body parts. Note that the "word complexity" in this study was judged by the patients themselves. Even if the patients themselves consider it to be simple, it is possible that medical misunderstandings may have occurred.

## Acknowledgements

## References

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *TACL*, 5:135–146.

Jacob Cohen. 1968. Weighted Kappa: Nominal Scale Agreement Provision for Scaled Disagreement or Partial Credit. *Psychological Bulletin*, 70(4):213–220.

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

Jennifer Fong Ha and Nancy Longnecker. 2010. Doctor-Patient Communication: A Review. *Ochsner Journal*, 10(1):38–43.

Kaoru Ito, Hiroyuki Nagai, Taro Okahisa, Shoko Wakamiya, Tomohide Iwao, and Eiji Aramaki. 2018. J-MeDic: A Japanese Disease Name Dictionary Based on Real Clinical Usage. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, pages 2365–2369.

Sebastian Joseph, Kathryn Kazanas, Keziah Reina, Vishnesh Ramanathan, Wei Xu, Byron Wallace, and Junyi Jessy Li. 2023. Multilingual Simplification of Medical Texts. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16662–16692.

Tomoyuki Kajiwara and Mamoru Komachi. 2018. Complex Word Identification Based on Frequency in a Learner Corpus. In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 195–199.

Tomoyuki Kajiwara, Daiki Nisihara, Tomonori Kodaira, and Mamoru Komachi. 2020. Language Resources for Japanese Lexical Simplification. *Journal of natural language processing*, 27(4):189–210.

Taku Kudo, Kaoru Yamamoto, and Yuji Matsumoto. 2004. Applying Conditional Random Fields to Japanese Morphological Analysis. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 230–237.

Gondy Leroy and James E. Endicott. 2012. Combining NLP with Evidence-Based Methods to Find Text Metrics Related to Perceived and Actual Text Difficulty. In *Proceedings of the 2nd ACM SIGHIT International Health Informatics Symposium*, pages 749–754.

Mounica Maddela and Wei Xu. 2018. A Word-Complexity Lexicon and A Neural Readability Ranking Model for Lexical Simplification. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3749–3760.

Kikuo Maekawa, Makoto Yamazaki, Takehiko Maruyama, Masaya Yamaguchi, Hideki Ogura, Wakako Kashino, Toshinobu Ogiso, Hanae Koiso, and Yasuharu Den. 2010. Design, Compilation, and Preliminary Analyses of Balanced Corpus of Contemporary Written Japanese. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation*, pages 1483–1486.

Ellie Pavlick and Chris Callison-Burch. 2016. Simple PPDB: A Paraphrase Database for Simplification. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 143–148.

Ellie Pavlick and Ani Nenkova. 2015. Inducing Lexical Style Properties for Paraphrase and Genre Differentiation. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 218–224.

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, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12(85):2825–2830.

Hideya Yamamoto, Kaoru Ito, and Eiji Aramaki. 2019. Fukugougo no Kouseiso Jouhou wo Kouryo Shita Byoumei Nannido no Suitei (Estimation Methods for Medical Term's Difficulty Utilizing Information on Constituents). In *Proceedings of the 25th Association for Natural Language Processing*, pages 1495–1498. (in Japanese).

Ziyu Yang, Santhosh Cherian, and Slobodan Vucetic. 2023. Data Augmentation for Radiology Report Simplification. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1922–1932.

# TEAM MIPAL at MEDIQA-M3G 2024:
# Large VQA Models for Dermatological Diagnosis

**Hyeonjin Kim**   **MIN KYU KIM**   **Jae Won Jang**
**KiYoon Yoo**   **Nojun Kwak**[*]
Seoul National University
{peaceful1,alsrb7000,pert0407,961230,nojunk}@snu.ac.kr

## Abstract

This paper describes the methods used for the NAACL 2024 workshop MEDIQA-M3G shared task (wai Yim et al., 2024a) for generating medical answers from image and query data for skin diseases. MedVInT-Decoder (Zhang et al., 2023b), LLaVA (Liu et al., 2024), and LLaVA-Med (Li et al., 2024) are chosen as base models. Finetuned with the task dataset on the dermatological domain, MedVInT-Decoder achieved a BLEU score of 3.82 during competition, while LLaVA and LLaVA-Med reached 6.98 and 4.62 afterward, respectively.

## 1 Introduction

The advancement of telecommunication technologies and the increasing demand for healthcare services have accelerated the demand for remote disease diagnosis and treatment. However, existing medical-related multimodal problems have predominantly focused on general diseases or radiology image analysis. In this task, abnormal skin images along with a single conversational query from the patient are provided as inputs. The goal is to utilize a multimodal model to identify the patient's condition and generate appropriate responses from the physician tailored to the patient's situation.

To address this problem scenario, we finetuned three multimodal VQA models, MedVInT-Decoder (Zhang et al., 2023b), LLaVA (Liu et al., 2024), LLaVA-Med (Li et al., 2024). When finetuned with the task dataset (wai Yim et al., 2024b), the BLEU scores of the MedVInT-Decoder, LLaVA, and LLaVA-Med were 3.82, 6.98, 4.62.The results for LLaVA and LLaVA-Med were submitted after the challenge. Since the released train dataset was small, we further explored ways to augment this data. Specifically, we crawled skin disease images online and synthesized query-response pairs using GPT-3.5. However, models

trained on the synthetic data reached a BLEU score lower than 1. The reason for such failure is discussed in Section 5.

## 2 Related Works

### 2.1 Visual Question-Answering

Multimodal models that target Visual Question-Answering tasks (VQA) are mostly consisted of a vision encoder, a text encoder and a decoder that decodes the encoded image and text at once. Some models use ViT as the vision encoder (Yu et al., 2022; Chen et al., 2022; Liu et al., 2024) while others employ ConvNet such as ResNet-50 (Wang et al., 2021). The encoded visual features are subsequently processed through the projection layer, where they are transformed into the word embedding space. The language instructions along with the projected image features are concatenated and inputted to the language model decoder to generate the output texts.

### 2.2 VQA on Medical Domain

MedVInT (Zhang et al., 2023b) uses pretrained ResNet-50 from PMC-CLIP (Lin et al., 2023) as the vision encoder, and PMC-LLaMA (Wu et al., 2023) as the language model. The model is then pretrained on a large dataset for VQA tasks on medical domain. MedVInT comes in two different forms: one uses encoder-based language model and the other uses decoder-based one as the generator. MedVInT-Decoder seems to generate more human-like answers in our experiments and is chosen as a base model.

LLaVA (Liu et al., 2024) employs ViT-L/14 from CLIP (Radford et al., 2021) to encode images, and Vicuna (Chiang et al., 2023) to encode and generate texts. LLaVA-Med (Li et al., 2024) is a LLaVA baseline model finetuned on medical dataset of 178k text queries and 61k images across X-ray, MRI, histology, gross pathology and CT domains.

---

[*]Corresponding Author

| |
|---|
| **encounter_id:** ENC00966 |
|  |
| **query_title:** "" |
| **query_content_en:** "Patient is 46 years old, female. The problem came up at the back of the body after wearing old clothes. Symptom felt: acute itchiness, small amount of discharge leaking from the affected area. It started small, but outburst to patches after washing with warm water. It is not improving after taking Clarityne. Now apply topically Disong camphor thin cream (A Chinese herbal antiinflammatory, antiallergic ointment). Orally taking Ketotifen." |
| **Reference:** "Based on the medical history and picture, it should be Allergic dermatitis. Use antiallergy treatment." <br> **MedVInT workshop:** "It is a case of allergic dermatitis." <br> **LLaVA 29 diseases:** "Scabies. Treatment involves killing mites and eggs with medicated cream or pill." <br> **LLaVA 324 diseases:** "Malassezia folliculitis. Treat with antifungal agents. Consider oral antifungal medications like ketoconazole." <br> **LLaVA workshop\*:** "Dermatitis" |

Table 1: Example of generated responses using our pipeline.

## 3 Method

### 3.1 Dataset Preprocessing

The chosen base models take a single image-query pair as their input and generate a single response as their output. The task dataset has multiple images for input and multiple possible responses per query. This requires a selection process to match the model input structure.

For each query, we sort the responses by the reliability of the authors, which are determined by their level of expertise. Then, we pair each of the responses with a single image. For example, if a query had 5 images and 3 different answers, three image-answer pairs are created for the query. Since some responses were exactly identical, we remove the duplicated responses to give diversity.

Some authors of the responses are ranked with low reliability score. One can consider removing these from the dataset to improve the validity of the train set. In this paper, we chose not to remove such answers to keep the dataset as large as possible. After the whole process, we obtained 2101 triplets of image, query and response. We train only on the English data.

### 3.2 Synthetic Data Generation

To augment the limited number of training samples, we attempted to generate synthetic data. Two different datasets were collected from two sources: one with 29 classes of most common diseases (Furue et al., 2011; Li et al., 2023; Zaidi and Lanigan, 2010) and another containing up to 324 dis-

eases (Atlas). The two datasets are denoted as '29 diseases' and '324 diseases' each from below. Associated queries attached to the images are generated with GPT-3.5 to imitate the answers in the given dataset.

### 3.3 Finetuning MedVInT

We finetune the original MedVInT-Decoder model with the processed dataset. Many training options were tested to find the optimal epochs, batch size, learning rate and early stopping. Then the same tests were executed again with other pretrained language models such as BioMedGPT (Zhang et al., 2023a) and MedAlpaca (Han et al., 2023).

The given training set is not large compared to the model size. The composition of the dataset is also noisy as uncertain and unreliable responses were included as well to maintain the size as much as possible. Such noisiness may have caused the slow convergence and prevented the training loss to converge to a lower number. In most cases, the training loss was over 2.5 which is quite large when considering that the metric used to evaluate loss in the model was Cross Entropy Loss. Training for more than 3 epochs results in overfitting because of the small dataset size.

### 3.4 Finetuning LLaVA variants

We adopt the pretrained LLaVA model as the baseline and subsequently conduct finetuning of both the projection matrix and the language model using low ranked adaption (LoRA) (Hu et al., 2021). We train the LLaVA for one epoch with a batch size

| Train Data | Model | BLEU | BERT |
|---|---|---|---|
| - | MedVInT | 0.91 | 0.82 |
| | LLaVA | 0.93 | 0.84 |
| | LLaVA-Med | 1.35 | 0.85 |
| MEDIQA-M3G | MedVInT | 3.82 | 0.87 |
| | LLaVA* | 6.98 | 0.86 |
| | LLaVA-Med* | 4.62 | 0.84 |
| Synthetic | LLaVA (29 diseases) | 0.92 | 0.84 |
| | LLaVA (324 diseases) | 0.98 | 0.85 |

Table 2: BLEU score and BERTscore of models fine-tuned on different datasets. The scores for LLaVA workshop and LLaVA-Med workshop was attained after the competition. All scores were measured on the final test set.

of 8 and a gradient accumulation step of 16. The LoRA hyperparameter $r$ was set to 128 and $\alpha$ was set to 256. The learning rates were set to 2e-5 for the projection layer and 2e-4 following the original configuration. We similarly inputted single image per query as the model was not finetuned on multiple images. Throughout the finetuning process, which spanned 10 epochs, we utilized the validation dataset to select the checkpoint from the epoch with the lowest evaluation loss. We employed the pretrained LLaVA-Med model, which is a version of LLaVA that has been finetuned on a medical dataset (Li et al., 2024). Finetunning was done for 10 epochs with a batch size of 8 and gradient accumulation step of 16.

When training only on the task data, we denote it by "MEDIQA-M3G". The results for training additionally on the synthetic data is denoted by "Synthetic".

## 4 Results

Examples of the generated responses are provided in Table 1 and summarized results are in Table 2. All models and checkpoints are evaluated using the official test set only.

MedVInT-Decoder shows low BLEU score of 0.91 when the inference is made directly on the test set without any finetuning. During training, MedVInT-Decoder reaches the lowest validation error after around 3 epochs and starts to show signs of overfitting afterwards. Early stopping is introduced to make use of such pattern.

The BLEU score measured with the validation set was the highest when trained with learning rate of 4e-6, batch size 16 and epochs 10. The train-

ing process stopped early at epoch 3. The BLEU score was 4.48 on the validation set and 3.82 on the test set. Although not impressively high, the increased values prove that actual learning has been conducted.

In case of LLaVA, inferencing with the vanilla model without any finetuning yielded a BLEU score of 0.93 which is similar to that of MedVInT. Upon finetuning the model with 29 diseases set and 324 diseases set, the BLEU scores on the validation dataset showed 0.97 and 4.51, respectively. The latter seemed promising, but did not meet the expected performance when inferenced on the test set.

LLaVA-Med scoreed a BLEU score 1.35 on zero shot inference. Afterwards, when finetuning on the task data, LLaVA-Med achieved a BLEU score of 3.82 and BERT score 0.84.

## 5 Discussion

**The Necessity of Finetuning** Although there already exist models pretrained with medical data, they mostly failed to give high-quality answers without further finetuning on the domain specific data. Two observations were made to explain this phenomenon.

One reason would be the difference of the image domain. Both PMC-VQA and LLaVA-Med training set are composed of professional images such as X-Ray, MRI and CT. The raw photos are not many in these sets, and it becomes even scarcer when limited to skin diseases. Therefore the task images would have been regarded as new and unfamiliar to the pretrained models.

Another possible reason would be the difference in the text domain. MedVInT-Decoder trained with PMC-VQA has learned to give short answers for most of the time. PMC-VQA is consisted of simple yes or no questions, or those that can be answered with one or two vocabularies. This may not have been sufficient when it comes to diagnosis and prescription tasks that require the machine to give long answers. Also, non-negligible amount of texts used for training is structured with certain formats or may be excerpts from academic texts. This is in contrast to the task dataset which is mostly written in casual spoken language.

**Low Performance When Trained with Crawled Data** Upon observing lower-than-expected performance following finetuning of the model with additional crawled data, efforts were directed towards

---
*After-challenge submission

enhancing LLaVA's performance after the competition. Inspired by the impressive performance demonstrated by MedVInT upon finetuning with solely the workshop-provided training data, a similar approach was applied to LLaVA and LLaVA-Med. LLaVA in this method yielded BLEU and BERT scores of 6.98 and 0.86, respectively, representing the highest scores achieved on the test dataset as shown is Table 2. This outcome underscores LLaVA's superior performance compared to MedVInT and LLaVA-Med. Furthermore, it suggests potential disparities between the dataset synthetically generated by ChatGPT and the real-world data. Lastly, utilizing BLEU and BERT scores as metrics implies that achieving similar linguistic nuances may contribute to superior performance, rather than merely focusing on the accuracy of individual predictions.

## 6 Conclusion

We propose our submission to the MEDIQA-M3G shared task for generating medical responses to multimodal queries. In our study, three existing models MedVInT, LLaVA and LLaVA-Med are finetuned using the competition dataset along with synthetically generated dataset. Their performance are evaluated using BLEU and BERT score. Our results indicate that utilizing only the task dataset leads to substantial improvements in both models, reaching the BLEU score of 3.82 with MedVInT-Decoder and ranked second in the English section of the workshop. After the competition, LLaVA finetuned with the workshop dataset achieved the highest BLEU score of 6.98 and lastly finetuned LLaVA-Med model with workshop dataset performed a BLEU score of 4.62.

## References

Dermatology Atlas. DermatologyAtlas. https://www.atlasdermatologico.com.br. Accessed: 2024-04-04.

Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. 2022. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6.

Masutaka Furue, Souji Yamazaki, Koichi Jimbow, Tetsuya Tsuchida, Masayuki Amagai, Toshihiro Tanaka, Kayoko Matsunaga, Masahiko Muto, Eishin Morita, Masashi Akiyama, et al. 2011. Prevalence of dermatological disorders in japan: a nationwide, cross-sectional, seasonal, multicenter, hospital-based study. *The Journal of dermatology*, 38(4):310–320.

Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bressem. 2023. Medalpaca–an open-source collection of medical conversational ai models and training data. *arXiv preprint arXiv:2304.08247*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2024. Llavamed: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36.

Huanyu Li, Peng Zhang, Zikun Wei, Tian Qian, Yiqi Tang, Kun Hu, Xianqiong Huang, Xinxin Xia, Yishuang Zhang, Haixing Cheng, et al. 2023. Deep skin diseases diagnostic system with dual-channel image and extracted text. *Frontiers in Artificial Intelligence*, 6.

Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-clip: Contrastive language-image pre-training using biomedical documents. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 525–536. Springer.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Wen wai Yim, Asma Ben Abacha, Velvin Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen, and Martin Krallinger. 2024a. Overview of the mediqa-m3g 2024 shared task on multilingual and multimodal medical answer generation. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wen wai Yim, Velvin Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904*.

Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Further fine-tuning llama on medical papers. *arXiv preprint arXiv:2304.14454*.

Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*.

Zohra Zaidi and Sean W Lanigan. 2010. *Dermatology in clinical practice*. Springer Science & Business Media.

Kai Zhang, Jun Yu, Zhiling Yan, Yixin Liu, Eashan Adhikarla, Sunyang Fu, Xun Chen, Chen Chen, Yuyin Zhou, Xiang Li, et al. 2023a. Biomedgpt: A unified and generalist biomedical generative pre-trained transformer for vision, language, and multimodal tasks. *arXiv preprint arXiv:2305.17100*.

Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023b. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*.

# MediFact at MEDIQA-M3G 2024: Medical Question Answering in Dermatology with Multimodal Learning

**Nadia Saeed**

Computational Biology Research Lab
Department of Computer Science
National University of Computer and Emerging Sciences (NUCES-FAST)
Islamabad, Pakistan
i181606@nu.edu.pk

## Abstract

The MEDIQA-M3G 2024 challenge necessitates novel solutions for Multilingual & Multimodal Medical Answer Generation in dermatology (wai Yim et al., 2024a). This paper addresses the limitations of traditional methods by proposing a weakly supervised learning approach for open-ended medical question-answering (QA). Our system leverages readily available MEDIQA-M3G images via a VGG16-CNN-SVM model, enabling multilingual (English, Chinese, Spanish) learning of informative skin condition representations. Using pretrained QA models, we further bridge the gap between visual and textual information through multimodal fusion. This approach tackles complex, open-ended questions even without predefined answer choices. We empower the generation of comprehensive answers by feeding the ViT-CLIP model with multiple responses alongside images. This work advances medical QA research, paving the way for clinical decision support systems and ultimately improving healthcare delivery. [1]

## 1 Introduction

Dermatological telemedicine consultations, while offering a promising solution for remote diagnosis and treatment, face hurdles due to limitations in capturing subtle visual details and the inability to physically examine lesions. This can lead to miscommunication, such as difficulties in describing the texture or progression of lesions, which can hinder the development of effective treatment plans (Elsner, 2020; Hwang et al., 2024; Mehraeen et al., 2023). However, recent advancements in image-text learning, like Vision Transformer (ViT) for image captioning and Contrastive Language-Image Pre-Training (CLIP) for aligning text and image representations, offer promising avenues to bridge this gap (Yin et al., 2022; Li et al., 2021).

Existing approaches to teledermatology consultations have limitations. Traditional consumer health question-answering systems primarily focus on textual data, neglecting the valuable information within visual details (Abacha et al., 2019b). This limits their ability to understand the nuances of skin conditions often best captured visually. Visual question-answering efforts have mainly targeted radiology images, overlooking the crucial context provided by clinical text (Abacha et al., 2019a). While recent advancements in deep learning have shown promise in lesion classification for dermatology (Li et al., 2022), these approaches often focus on specific image types and cannot integrate textual information, essential for a holistic understanding of a patient's condition. While some research explores combining clinical text and images for specific dermatology tasks, such as melanoma risk assessment, they haven't addressed open-ended question answering (Groh et al., 2022; Lin et al., 2023).

This research tackles these limitations by introducing a novel framework for multilingual and multimodal query response generation in clinical dermatology. Our system leverages the power of multimodal fusion, which combines information from different sources. In this case, the sources are textual and visual: textual clinical context and user queries in multiple languages, along with user-uploaded images. This work introduces Medifact-M3G, a framework for tackling uncertainties in medical question answering for dermatology shown in Figure 1. Medifact-M3G prepares the data and assigns weights to potential answers, considering their relevance and trustworthiness (Section a). It then uses a powerful image analysis tool to extract key features from skin condition images (Section b). By combining these features with text analysis, Medifact-M3G leverages multiple powerful models to generate informative answers to medical questions (Sections c and d). This framework has the potential to improve the accuracy

---

[1] Fine-tuned models and Code avaliable: https://github.com/NadiaSaeed/MediFact-M3G-MEDIQA-2024

Figure 1: MediFact-M3G Framework: From Uncertain Data to Informed Answers

and reliability of AI-powered diagnosis systems in telemedicine, ultimately assisting healthcare professionals in providing better diagnoses and treatment plans. This research addresses the following key questions: 1) Can feature fusion from weakly supervised learning techniques effectively support open-ended medical question answering in dermatology? 2) Can a Medifact-M3G fine-tuned model trained solely on the MEDIQA-M3G training dataset adequately capture similarities and relatedness for unseen samples? 3) How can contrastive learning be seamlessly integrated with Medifact-M3G to quantify uncertainty in response generation for ambiguous queries and limited content information?

## 2 Methodology

Our response generation system for the MEDIQA-M3G 2024 task tackles the challenge of limited labeled data while aiming to generate informative responses to user queries about dermatological conditions (wai Yim et al., 2024a). This methodology leverages several key steps, as illustrated in the accompanying MediFact-M3G framework shown in Figure 1.

### 2.1 Data Preprocessing and Response Weighting

We begin by ensuring the quality of the raw data through techniques like handling missing values, text cleaning, and formatting consistency. This establishes a clean and consistent foundation for subsequent model training.

Next, a weighting function assigns scores to each response based on the author's expertise (e.g., medical doctor) and response completeness. This guides the model to prioritize learning from the most effective responses during training, ultimately improv-

340

Figure 2: Example of Original Text "MEDIQA-M3G" and System Output "MediFact-M3G

ing the quality of the generated responses.

## 2.2 Weakly Supervised Learning for Image Representation: Addressing Data Limitations

While large, labeled datasets are ideal for training robust response generation models in dermatology, ethical considerations, and data access limitations often restrict their availability. To address this challenge, we employed a weakly supervised learning approach that leverages the available data effectively.

Our approach utilizes a pre-trained Convolutional Neural Network (CNN), specifically VGG16, to extract high-level features from the dermato-logical images. These features capture the visual characteristics relevant to diagnosis (Desai et al., 2021). We then use a Support Vector Machine (SVM) classifier to learn the relationship between the extracted image features and the high-quality textual responses associated with labeled image-response pairs. The SVM essentially learns to map images to their most relevant textual descriptions (Chandra and Bedi, 2021).

This weakly supervised approach allows us to overcome limitations in labeled data. The SVM generalizes the learned relationship between labeled image-response pairs to unlabeled images. By incorporating the information gleaned from the textual responses, this process enriches the im-

age representations learned by the VGG16 model, even without explicit labels for each unlabeled image. These enriched image representations (English, Chinese, and Spanish languages) capture the semantic meaning associated with the images, providing valuable information for the response generation model during training. Additionally, for comparison purposes, we evaluated the performance of Inception and ResNet models in place of VGG16 to determine the most effective CNN architecture for this task (Zheng et al., 2021; Zhou et al., 2021).

## 2.3 Multi-Model Response Generation with Feature Fusion

This step focuses on generating responses to user queries. We employ a multi-model approach that combines pre-trained question-answering (QA) models with the image representation learned from the weakly supervised approach described in Section 2 (Cortiz, 2022). Due to limitations in the performance and availability of non-English language models, this step focuses on English responses.

A comprehensive feature vector for each query-response pair is created by combining the following elements:

- The user's query itself.

- Relevant textual content (e.g., patient demographics).

- The image representation learned from the weakly supervised approach (Section 2).

We utilize two pre-trained English models:

### 2.3.1 Extractive QA Model

This model retrieves relevant answer passages from a text corpus (potentially including high-quality responses) that directly address the user's query (Guo et al., 2023; Clark et al., 2020; He et al., 2021).

### 2.3.2 Abstractive QA Model

This model goes beyond retrieval and generates a new, comprehensive response. It incorporates information from various sources (textual features, extracted passages) and potential reasons over the information to provide a more informative answer (Lewis et al., 2019).

This multi-model approach offers the advantage of combining factual grounding from the extractive model with flexible response generation from the abstractive model, while also incorporating visual

information through the image features. This ultimately leads to more accurate and informative responses within the teledermatology domain.

## 2.4 Response Selection with Contrastive Learning

Selecting the most informative response for a query-image pair, especially in non-English settings, requires a robust approach. We leverage CLIP, a contrastive learning model adept at learning relationships between image and text embeddings (Li et al., 2021). CLIP utilizes a Vision Transformer (ViT) (Section 2) to extract high-dimensional image features and a separate text encoder for potential responses (Yin et al., 2022). We employ CLIP in two key settings: First, CLIP receives the ViT-extracted image embedding and multiple response lists (English, Spanish, Chinese). It calculates the cosine similarity between each response embedding (in a specific language) and the image embedding. The response with the highest similarity (closest semantic relationship) is chosen for that language. Second, CLIP focuses on the relationship between the image and English responses from pre-trained QA models (Section 3). It assesses the cosine similarity between the image embedding and the selected English response embedding. Google Translate then converts this English response to Spanish and Chinese for user convenience, acknowledging potential translation inaccuracies (Taira et al., 2021).

## 3 Experimental Setup and Results

We evaluated our model's capability in addressing the problem of clinical dermatology multimodal query response generation. This evaluation was conducted within the Shared Task of MEDIQA-M3G 2024, which focuses on multilingual and multimodal medical answer generation (wai Yim et al., 2024a). As illustrated in Figure 2, each sample in the task comprised k medical images related to dermatological conditions, a textual query describing the user's skin concern, and its content. Additionally, the ground truth for each sample included multiple possible responses with corresponding scores. Leveraging the framework outlined in Figure 1, our Medifact-M3G model was employed to generate answers in three languages for each sample.

### 3.1 Dataset

The MEDIQA-M3G dataset is divided into training (842 instances), validation (56 instances), and

test (100 instances) sets, with each set available in Chinese, English, and Spanish versions (wai Yim et al., 2024b). While non-English training sets are machine-translated, validation and test sets are human-translated for accuracy. Each instance is represented as a JSON object containing a unique encounter ID, a list of image IDs, the query title and content in the specific language, and author information from a separate CSV file. Participants are expected to generate responses in JSON format, including a unique encounter ID and a list of generated responses for the specified languages. Participation in all language evaluations is optional, with empty strings allowed for non-participating languages.

## 3.2 Evaluation Metrics

Our system's performance was evaluated using official available evaluation program of MEDIQA-M3G [2]. metrics commonly employed in Natural Language Generation (NLG) tasks. DeltaBLEU and BERTScore were chosen for this assessment [cite]. DeltaBLEU measures the similarity between a generated response and reference responses by considering n-gram (sequence of n words) overlap but weighs these n-grams based on human judgment. BERTScore, on the other hand, focuses on the semantic similarity between the generated response and the references, taking the maximum score from any available reference response. The evaluation script processed instances across three languages (English, Spanish, and Chinese).

## 3.3 Result

In this study, we evaluated our approach using the Mediqa-M3G framework, employing three different feature extraction models while maintaining consistency in other aspects of the setup. These models included SVMs with default sklearn settings and pre-trained CNN architectures like from the Keras library. The evaluation results are summarized in Table 1.

Table 1 displays the evaluation results for two setups of the MediFact-M3G framework. In the first setup, denoted as VGG16-Individual, separate VGG16 models were trained for each language, yielding individual scores for each language. In the second setup, the best-performing VGG16 model output, which was observed to be the Chinese language model, was utilized to translate responses

into English and Spanish languages following the MediFact-M3G framework. While the translated version of MediFact-M3G showed slight improvement in BERT_Score, the Deltableu score performed better in the individual setup for Spanish language responses.

Additionally, it's worth noting our performance in the MEDIQA-M3G 2024 shared task, where we achieved 7th rank in English language response generation, and 3rd rank in Chinese and Spanish language response generation, out of a total of 75 participants. These rankings underscore the effectiveness of our approach across different languages and its competitiveness in challenging benchmark tasks (wai Yim et al., 2024a).

## 3.4 Discussion

The results presented here were obtained after rigorous testing in a challenging setting, providing insights into the performance of different feature extraction models within the MediFact-M3G framework. VGG16-Translated demonstrated significant improvements over VGG16-Individual, underscoring the effectiveness of data translation in enhancing translation quality. The evaluation results are summarized in Table 2.

After replacing the VGG16 models with ResNet and SqueezNet in MediFact-M3G framework, we obtained the following evaluation results as shown in Table 2. SqueezNet demonstrated exceptional proficiency in Chinese translations, achieving the highest Deltableu scores across all languages. On the other hand, although ResNet exhibited slightly lower Deltableu scores, its competitive performance across all languages highlights its versatility in handling various translation tasks. These findings underscore the critical role of selecting appropriate feature extraction models tailored to specific language requirements and task objectives, ultimately enhancing the effectiveness of the MediFact-M3G framework in addressing medical query challenges.

## 4 Future Work

In the future, we plan to conduct further experiments to explore the robustness and scalability of our approach across larger and more diverse datasets. Additionally, we aim to investigate the integration of domain-specific ontologies and medical terminologies to enhance the semantic understanding and accuracy of our system. Furthermore,

| Model | | Deltableu | | | BERT_Score | | |
|---|---|---|---|---|---|---|---|
| | | en | zh | es | en | zh | es |
| MediFact-M3G | VGG16-Individual | 0.588 | 4.503 | **0.918** | 0.837 | **0.771** | 0.804 |
| | VGG16-Translated | **0.717** | **4.503** | 0.823 | **0.842** | 0.763 | **0.809** |

Table 1: Scores for Response Generation Approaches on MEDIQA-M3G Testing Dataset (submitted at the competition)

| Model | | Deltableu | | | BERT_Score | | |
|---|---|---|---|---|---|---|---|
| | | en | zh | es | en | zh | es |
| MediFact-M3G | VGG16-Individual | 0.588 | 4.503 | **0.918** | **0.845** | 0.763 | 0.806 |
| | VGG16-Translated | 0.717 | 4.503 | 0.823 | 0.842 | 0.763 | **0.809** |
| | ResNet | 0.565 | **6.457** | 0.542 | 0.837 | **0.771** | 0.804 |
| | SqueezNet | **0.744** | 2.125 | 0.641 | 0.841 | 0.702 | 0.808 |

Table 2: Scores for Response Generation Approaches on MEDIQA-M3G Testing Dataset (after the competition)

we are interested in exploring novel techniques for handling multi-turn dialogue scenarios, allowing our system to engage in more natural and interactive conversations with users. Additionally, we plan to collaborate with medical professionals to validate the clinical relevance and effectiveness of our approach in real-world healthcare settings. By addressing these challenges, we hope to continue advancing the field of medical question-answering and contribute to the development of more practical and clinically useful systems.

# References

Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Joey Liu, Dina Demner-Fushman, and Henning Müller. 2019a. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. *CLEF (working notes)*, 2(6).

Asma Ben Abacha, Chaitanya Shivade, and Dina Demner-Fushman. 2019b. Overview of the mediqa 2019 shared task on textual inference, question entailment and question answering. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 370–379.

Mayank Arya Chandra and SS Bedi. 2021. Survey on svm and their application in image classification. *International Journal of Information Technology*, 13(5):1–11.

Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators.

Diogo Cortiz. 2022. Exploring transformers models for emotion recognition: A comparision of bert, distilbert, roberta, xlnet and electra. In *Proceedings of the 2022 3rd International Conference on Control, Robotics and Intelligent System*, pages 230–234.

Padmashree Desai, Jagadeesh Pujari, C Sujatha, Arinjay Kamble, and Anusha Kambli. 2021. Hybrid approach for content-based image retrieval using vgg16 layered architecture and svm: an application of deep learning. *SN Computer Science*, 2(3):170.

Peter Elsner. 2020. Teledermatology in the times of covid-19–a systematic review. *JDDG: Journal Der Deutschen Dermatologischen Gesellschaft*, 18(8):841–845.

Matthew Groh, Caleb Harris, Roxana Daneshjou, Omar Badri, and Arash Koochek. 2022. Towards transparency in dermatology image datasets with skin tone annotations by experts, crowds, and an algorithm. *Proceedings of the ACM on Human-Computer Interaction*, 6(CSCW2):1–26.

Muzhe Guo, Muhao Guo, Edward T Dougherty, and Fang Jin. 2023. Msq-biobert: Ambiguity resolution to enhance biobert medical question-answering. In *Proceedings of the ACM Web Conference 2023*, pages 4020–4028.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.

Jonathan K Hwang, Natalia Pelet Del Toro, George Han, Dennis H Oh, Trilokraj Tejasvi, and Shari R Lipner. 2024. Review of teledermatology: lessons learned from the covid-19 pandemic. *American Journal of Clinical Dermatology*, 25(1):5–14.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.

Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. 2021. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. In *International Conference on Learning Representations*.

Zhouxiao Li, Konstantin Christoph Koban, Thilo Ludwig Schenck, Riccardo Enzo Giunta, Qingfeng Li, and Yangbai Sun. 2022. Artificial intelligence in dermatology image analysis: current developments and future trends. *Journal of clinical medicine*, 11(22):6826.

Zhihong Lin, Donghao Zhang, Qingyi Tao, Danli Shi, Gholamreza Haffari, Qi Wu, Mingguang He, and Zongyuan Ge. 2023. Medical visual question answering: A survey. *Artificial Intelligence in Medicine*, page 102611.

Esmaeil Mehraeen, SeyedAhmad SeyedAlinaghi, Mohammad Heydari, Amirali Karimi, Abdollah Mahdavi, Mehrnaz Mashoufi, Arezoo Sarmad, Peyman Mirghaderi, Ahmadreza Shamsabadi, Kowsar Qaderi, et al. 2023. Telemedicine technologies and applications in the era of covid-19 pandemic: A systematic review. *Health informatics journal*, 29(2):14604582231167431.

Breena R Taira, Vanessa Kreger, Aristides Orue, and Lisa C Diamond. 2021. A pragmatic assessment of google translate for emergency department instructions. *Journal of General Internal Medicine*, 36(11):3361–3365.

Wen wai Yim, Asma Ben Abacha, Velvin Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen, and Martin Krallinger. 2024a. Overview of the mediqa-m3g 2024 shared task on multilingual and multimodal medical answer generation. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wen wai Yim, Velvin Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. 2022. A-vit: Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10809–10818.

Xiangtao Zheng, Binqiang Wang, Xingqian Du, and Xiaoqiang Lu. 2021. Mutual attention inception network for remote sensing visual question answering. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–14.

Changjian Zhou, Jia Song, Sihan Zhou, Zhiyao Zhang, and Jinge Xing. 2021. Covid-19 detection based on image regrouping and resnet-svm using chest x-ray images. *Ieee Access*, 9:81902–81912.

# MediFact at MEDIQA-CORR 2024: Why AI Needs a Human Touch

**Nadia Saeed**

Computational Biology Research Lab
Department of Computer Science
National University of Computer and Emerging Sciences (NUCES-FAST)
Islamabad, Pakistan
i181606@nu.edu.pk

## Abstract

Accurate representation of medical information is crucial for patient safety, yet artificial intelligence (AI) systems, such as Large Language Models (LLMs), encounter challenges in error-free clinical text interpretation. This paper presents a novel approach submitted to the MEDIQA-CORR 2024 shared task (Ben Abacha et al., 2024a), focusing on the automatic correction of single-word errors in clinical notes. Unlike LLMs that rely on extensive generic data, our method emphasizes extracting contextually relevant information from available clinical text data. Leveraging an ensemble of extractive and abstractive question-answering approaches, we construct a supervised learning framework with domain-specific feature engineering. Our methodology incorporates domain expertise to enhance error correction accuracy. By integrating domain expertise and prioritizing meaningful information extraction, our approach underscores the significance of a human-centric strategy in adapting AI for healthcare.[1]

## 1 Introduction

Accurately identifying pathogens from textual descriptions of symptoms is crucial in effective healthcare management (Qian and Morral, 2022). However, existing datasets often present significant challenges that hinder reliable inferences and accurate pathogen identification, especially for rare diseases with limited data availability (Wang et al., 2021; Qian and Morral, 2022).

One major challenge lies in the inherent linguistic ambiguities present within these descriptions. Synonyms, homonyms, and polysemy (words with multiple meanings) can lead to confusion and misinterpretations (Karabacak and Margetis, 2023). For example, the term "fever" could indicate a wide range of illnesses, making it difficult to pinpoint the specific pathogen without additional context. Additionally, the distribution of diagnostic and pathogen information within the data can be imbalanced, with some diseases being vastly over-represented compared to others. This imbalance can skew the model's performance and hinder its ability to accurately identify pathogens for less frequently encountered diseases (Thirunavukarasu et al., 2023; Wang et al., 2021).

Furthermore, incorporating sensitive diagnostic data for training LLMs raises significant ethical concerns regarding patient privacy and authorization requirements (Kelly, 2002). Moreover, pre-trained LLMs often learn from vast amounts of generic text data, which might not be tailored to the specific domain of pathogenic research (Qian and Morral, 2022). This lack of domain-specific knowledge can hinder their ability to capture the nuances of rare disease entities and the intricate relationships between textual descriptions and underlying pathogens (Thirunavukarasu et al., 2023; Chanda et al., 2022).

Existing approaches to medical text correction have explored various techniques, including rule-based systems like MetaMap (which utilizes predefined rules to map terms to standardized medical concepts) and machine learning algorithms like RNN-based models (trained to identify and correct errors based on patterns learned from training data) (Chanda et al., 2022; Kumar et al., 2021; Minaee et al., 2021). However, these methods often struggle with the complexity of medical terminology, the inherent ambiguities of natural language, and the limitations of rule-based systems in capturing the ever-evolving nuances of medical language (Qian and Morral, 2022).

While recent advancements in LLMs have shown promise in various natural language processing tasks like text correction, their application in medical diagnostics necessitates careful considera-

---

[1]Code is available: https://github.com/NadiaSaeed/MediFact-MEDIQA-CORR-2024

tion due to the sensitivity of the data and the need for domain-specific knowledge. Existing LLM-based medical text correction approaches primarily address basic issues like typos and grammatical errors (Thirunavukarasu et al., 2023; Lee et al., 2022). However, they often fall short in addressing patient hallucinations, which can introduce factual errors and lead to misdiagnosis (Wang et al., 2023). Additionally, fine-tuning these models on relevant datasets often yields limited improvements, with models producing generic corrections instead of medically accurate ones (Lee et al., 2022).

This paper aims to present a methodology for automatically correcting single-word errors in clinical notes, submitted to the MEDIQA-CORR 2024 shared task (Ben Abacha et al., 2024a). The approach utilizes supervised learning with tailored feature engineering for the medical domain, emphasizing meaningful information extraction from clinical text data. Two distinct strategies are employed: an extractive question-answering (QA) approach for observed error-correction pairs and an abstractive QA approach for unobserved relations. This framework addresses the following important research questions: 1) How can domain expertise be further integrated into the model to improve its accuracy and ability to explain its reasoning? 2) How can this approach be effectively utilized to assist human reviewers in the process of medical record correction, potentially improving efficiency and accuracy? 3) What ethical considerations are involved in using AI for automatic error correction in healthcare settings, such as potential bias, transparency, and accountability?

## 2 Methodology

This paper introduces MediFact-CORR QA, a data-efficient approach for one-word error correction in clinical text paragraphs. MediFact-CORR QA leverages a two-stage process combining weakly supervised learning with pre-trained models to address labeled medical text data limitations.

### 2.1 Error Sentence Identification with Weak Supervision Motivation

MediFact-CORR QA, an innovative framework, employs weakly-supervised learning to discern distinctive patterns in clinical errors within textual data. The process involves analyzing paired paragraphs, each comprising an error-laden version and its corrected counterpart, with the error ex-plicitly annotated. Utilizing Support Vector Machines (SVMs) (Jamaluddin and Wibawa, 2021), the framework effectively discriminates between accurate and erroneous sentences within the clinical domain as shown in Figures 1 and 2 respectively.

This methodology capitalizes on the inherent information within error sentences, thereby mitigating the necessity for extensive labeled datasets. Moreover, the model not only indicates the presence of an error but also precisely identifies the erroneous sentence's location when applicable. Initially training separate SVMs for error and correct sentences, the model's efficacy during testing is indirectly enhanced by the utilization of supervised training labels. Consequently, MediFact-CORR QA proficiently tags erroneous sentences based on acquired patterns from the paired training data.

### 2.2 Error Correction with Extractive QA

Furthermore, in the process of generating correct sentences, MediFact-CORR QA relies on the inherent structure of the training data and adopts an extractive QA methodology. A notable feature of the MEDIQA-CORR dataset is the existence of paragraph pairs, where one contains an error and the other presents the corrected version (Ben Abacha et al., 2024b). Leveraging this characteristic, MediFact-CORR QA focuses on these error-correction pairs. When identifying sentences as erroneous in Step 1, we apply fuzzy matching between them and their corresponding corrected counterparts from the training data. This fuzzy matching helps to annotate the error information and correct information accurately and efficiently. Through this process, we can locate the most probable correct sentence by finding the matched pair of paragraphs, as they closely resemble each other. Extractive QA proves advantageous in scenarios where the answer can be directly extracted from a given text source. In our context, since the corrected sentence is already present within the training data, MediFact-CORR QA efficiently identifies it through similarity matching. This approach stands out for its data efficiency and effectiveness. Figure 3 depicts the framework where matched paragraph pairs are considered, with one containing error information and the other representing the correct information. This behavior of our dataset is crucial for the extractive QA model, as it allows us to utilize the inherent information within the content. This information is then positioned using

Figure 1: MediFact-CORR: Framework of the Correct SVM model



Figure 2: MediFact-CORR: Framework of the Error SVM model

the previously trained SVM models.

## 2.3 Error Correction with Abstractive QA

Recognizing that not all errors will have corresponding corrected versions in the training data, MediFact-CORR QA employs a pre-trained question-answering (QA) model specifically tailored for unanswerable questions (Lewis et al., 2019). Sentences identified as erroneous in Step 1 but lacking a match in the training data are directed to this pre-trained model. Trained on a vast corpus of text and questions, this model can generate potential corrections for unseen errors by analyzing contextual relationships between words within the erroneous sentence. Pre-trained QA models, having been trained on extensive datasets, excel at handling unseen information and complex language (Cortiz, 2022). Consequently, MediFact-CORR QA can address errors not explicitly present in the training data, thereby enhancing its robustness and generalizability. To illustrate, Figure 4 depicts the

framework's step where sentences lacking matched pairs in the training data are passed through the pretrained QA model for potential corrections (Cortiz, 2022).

By integrating weakly-supervised error detection with extractive QA for observed corrections, and leveraging a pre-trained QA model for unseen errors, MediFact-CORR QA provides a data-efficient solution for error correction in clinical text. This approach is particularly valuable in contexts where access to large labeled medical text data is limited.

## 3 Experimental Setup and Results

This section details the experimental setup and evaluates the performance of our two-stage model for one-word error correction in clinical text paragraphs.

### 3.1 Dataset

The MEDIQA-CORR 2024 shared tasks that employed a dataset of clinical texts from the MS and

Figure 3: MediFact-CORR: Framework of the Error Correction with Extractive QA



Figure 4: MediFact-CORR: Framework of the Error Correction with Abstractive QA

UW collections (Ben Abacha et al., 2024b). The training set (MS collection) comprised 2,189 texts. Validation sets contained 574 texts from MS and 160 texts from UW. Each text along with the split sentences, Error sentence, and its index, and the corresponding correct sentence, is also given with an error flag. The testing set (MS and UW collection) comprised 925 texts. MEDIQA-CORR 2024 shared tasks comprise three challenging tasks to perform, 1) Error flag prediction, 2) Index of the error sentence detection, and 3) Generate correct sentence.

## 3.2 Evaluation Metrics

The evaluation has been performed using the available program file by the MEDIQA-CORR 2024 [2]. In performance evaluation following metrics include AggregateScore, R1F score, BERTSCORE, BLEURT, and AggregateC (Yuan et al., 2021; Sel-

lam et al., 2020). *AggregateScore* serves as an overarching metric, consolidating various aspects of model performance, while *R1F* score measures the effectiveness of error correction by considering precision, recall, and F1 score. Additionally, *AggregateC* provides a composite metric summarizing model performance across different dimensions. We also evaluate the model's ability to accurately identify sentences containing errors and pinpoint the precise location of these errors within sentences.

## 3.3 Results

The models underwent rigorous evaluation across various metrics, including error flag accuracy, error sentence detection accuracy, and Natural Language Generation (NLG) performance. Evaluation was conducted on the validation sets of the MEDIQA_CORR 2024 dataset (Ben Abacha et al., 2024b). Our experimental setup involved training the SVM models using a combination of both train-

---

[2]MEDIQA-CORR evaluation code: `https://github.com/abachaa/MEDIQA-CORR-2024`

ing and validation sets. These trained models are now available in our GitHub repository [3].

For the abstraction QA model utilized in the experiment, we leveraged the BART model to answer questions of diagnosing expected medical conditions from provided text (Lewis et al., 2019).

Our performance in the tasks was notably obtained scores out of 106 participants shown in Table 1 (Ben Abacha et al., 2024a). In Task 1 for Error Flags Accuracy, we secured the 2nd rank. For Task 2, which focused on Error Sentence Detection Accuracy, we attained the 8th rank. Task 3 evaluated the Aggregate Score for NLG, where we achieved the 14th rank. Overall, these results underscore the effectiveness of our two-stage model for one-word error correction in clinical text paragraphs, surpassing the performance of the provided baseline model. By integrating error flag prediction, precise sentence extraction, and NLG techniques, we present a promising approach to enhancing the quality and reliability of clinical text data.

## 4 Discussion

Large Language Models (LLMs) have shown remarkable success in various natural language processing tasks, but their application in medical text correction faces unique challenges (Thirunavukarasu et al., 2023; Wu et al., 2022). Our approach tackles the challenging task of correcting one-word errors in clinical text paragraphs. Unlike LLMs that rely solely on statistical patterns learned from vast amounts of text data, our approach utilizes features specifically tailored to the medical context. This allows the model to leverage domain knowledge and prioritize terms. The example demonstrating the limitations of LLMs and the strengths of SVMs with TF-IDF can be added as a separate paragraph in the same section, following the current paragraph.

Example paragraph: *'A 5-year-old male presents with complaints of a painful mouth/gums, and vesicular lesions on the lips and buccal mucosa for the past 4 days. He is unable to eat or drink due to the pain and reports muscle aches. Vital signs: T 39.1°C, HR 110, BP 90/62 mmHg, RR 18, SpO2 99%. Physical examination reveals vesicular lesions on the tongue, gingiva, and lips, with some ruptured and ulcerated, and palpable cervical and submandibular lymphadenopathy. Patient is diag-*

*nosed with an [MASK] infection.'*

While a fine-tuned DistillBERT model predicted a general term like 'goat' or 'Highlander' (Wu et al., 2022). On the other side, our SVM model trained with TF-IDF utilizes domain knowledge through feature weights (Quach et al., 2023). Features like 'vesicular lesions', 'lips', and 'gingiva' receive high weights, guiding the model towards the medically accurate prediction of 'HSV-1' due to its alignment with the clinical context."

Our journey focused on error detection and correction within clinical text data. While Transformer-based models are powerful, their limitations in interpretability, data requirements, and over-fitting prompted us to explore an alternative: SVMs with TF-IDF features. Unlike many models, SVMs offer valuable insights through feature weights (Campbell and Ying, 2022). Features were designed to recognize specific medical terms, abbreviations, and entities like drug names, diagnoses, and anatomical locations. Rules and patterns observed in common errors were translated into features (Quach et al., 2023). Features captured aspects like sentence structure, negation markers, and temporal inconsistencies, which can indicate factual errors like incorrect dates or inconsistent medication names.

The provided dataset posed a unique challenge due to pre-defined sentence indices that deviated from standard newline ("\n") splitting (Ben Abacha et al., 2024b). To address this challenge, we compared detected errors' content with the dataset's available sentences. The index reported in the "Error sentence index" column was predicted as the starting digit of the most similar sentence. Therefore, we must recognize that inherent dataset issues influenced our final score. These challenges underscore the significance of high-quality data for training machine learning models.

In our submission, we investigated three key outcomes in an alternative setting. In the first and second scenarios, utilizing a QA model instead of the static correction model of SVM resulted in an improved R1F score from 0.342 to 0.454. This enhancement underscores the effectiveness of employing a QA model for error correction tasks. Moreover, the accuracy of error sentence detection significantly increased from 0.39 to 0.6 by utilizing the starting digit of the most similar sentence in the pre-defined index of sentences within given samples. This improvement stemmed from addressing an index problem, specifically by selecting the in-

| Model | R1F | BERT | BLEURT | AggScore | AggC | Error Flag | Error Sentence |
|---|---|---|---|---|---|---|---|
| MediFact_CORR | 0.454 | 0.444 | 0.439 | 0.446 | 0.535 | 0.737 | 0.600 |

Table 1: Performance on error correction tasks, including error flags accuracy and error sentence detection accuracy (submitted at the competition).

dex from the upper part of the sentence. Table 2 provides a summary of these findings.

This research demonstrates the effectiveness of combining human expertise and AI through feature engineering in a supervised learning approach. While SVMs offer interpretability and efficiency, human collaboration remains crucial for optimal performance in complex domains like healthcare (Campbell and Ying, 2022). This collaboration paves the way for improved error detection and correction in clinical text data, ultimately leading to better patient care.

## 5 Future Work

Our initial success with SVMs for pathogen identification in clinical text data paves the way for further exploration using LLMs. However, LLMs pose unique challenges. Data scarcity, particularly in the specific medical domain, could be a significant hurdle (Wang et al., 2023). Limited data restricts the use of a separate validation set. Future work will explore acquiring more data and data augmentation to enhance model generalizability. Techniques like data augmentation and transfer learning from pre-trained medical LLMs might be crucial to overcome this limitation.

Ethical considerations are paramount, and mitigating biases within the training data is essential. Furthermore, ensuring interpretability through techniques like attention mechanisms is vital for trust and acceptance in healthcare settings.

Finally, for practical implementation, we need to explore computationally efficient LLM architectures or develop task-specific models focused on pathogen identification. Continuous evaluation through techniques like active learning and performance monitoring will be crucial for maintaining a robust, ethical, and interpretable system in real-world clinical text analysis.

## References

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the*

*6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Colin Campbell and Yiming Ying. 2022. *Learning with support vector machines*. Springer Nature.

Ashis Kumar Chanda, Tian Bai, Ziyu Yang, and Slobodan Vucetic. 2022. Improving medical term embeddings using umls metathesaurus. *BMC Medical Informatics and Decision Making*, 22(1):114.

Diogo Cortiz. 2022. Exploring transformers models for emotion recognition: A comparision of bert, distilbert, roberta, xlnet and electra. In *Proceedings of the 2022 3rd International Conference on Control, Robotics and Intelligent System*, pages 230–234.

M Jamaluddin and Adhi Dharma Wibawa. 2021. Patient diagnosis classification based on electronic medical record using text mining and support vector machine. In *2021 International Seminar on Application for Technology of Information and Communication (iSemantic)*, pages 243–248. IEEE.

Mert Karabacak and Konstantinos Margetis. 2023. Embracing large language models for medical applications: opportunities and challenges. *Cureus*, 15(5).

Curly Kelly. 2002. Hipaa compliance: Lessons from the repeal of hawaii's patient privacy law. *Journal of Law, Medicine & Ethics*, 30(2):309–312.

A Sampath Kumar, Leta Tesfaye Jule, Krishnaraj Ramaswamy, S Sountharrajan, N Yuuvaraj, and Amir H Gandomi. 2021. Analysis of false data detection rate in generative adversarial networks using recurrent neural network. In *Generative Adversarial Networks for Image-to-Image Translation*, pages 289–312. Elsevier.

Eun Byul Lee, Go Eun Heo, Chang Min Choi, and Min Song. 2022. Mlm-based typographical error correction of unstructured medical texts for named entity recognition. *BMC bioinformatics*, 23(1):1–16.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.

| Model | | R1F | BERT | BLEURT | AggScore | AggC | Error Flag | Error Sentence |
|---|---|---|---|---|---|---|---|---|
| MediFact_CORR | Corr+ \n Indexing | 0.342 | 0.355 | 0.419 | 0.372 | 0.508 | 0.737 | 0.600 |
| | QA-Model+ \n Indexing | **0.454** | **0.444** | **0.439** | **0.446** | **0.535** | **0.737** | **0.600** |
| | QA-Model | 0.409 | 0.401 | 0.418 | 0.409 | 0.353 | 0.507 | 0.398 |

Table 2: Performance comparison of different models on error correction tasks, including error flags accuracy and error sentence detection accuracy. The table showcases improvements achieved by employing a QA model and adopting a comprehensive approach to error flag annotation and error sentence detection (results before the competition).

Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Narjes Nikzad, Meysam Chenaghlu, and Jianfeng Gao. 2021. Deep learning–based text classification: a comprehensive review. *ACM computing surveys (CSUR)*, 54(3):1–40.

Gene Qian and Núria Morral. 2022. Role of non-coding rnas on liver metabolism and nafld pathogenesis. *Human Molecular Genetics*, 31(R1):R4–R21.

Luyl-Da Quach, Anh Nguyen Quynh, Nguyen Quoc Khang, and An Nguyen Thi Thu. 2023. Using the term frequency-inverse document frequency for the problem of identifying shrimp diseases with state description text. *International Journal of Advanced Computer Science and Applications*, 14(5).

Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*.

Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930–1940.

Haoqing Wang, Huiyu Mai, Zhi-hong Deng, Chao Yang, Luxia Zhang, and Huai-yu Wang. 2021. Distributed representations of diseases based on co-occurrence relationship. *Expert Systems with Applications*, 183:115418.

Hongyan Wang, WeiZhen Wu, Zhi Dou, Liangliang He, and Liqiang Yang. 2023. Performance and exploration of chatgpt in medical examination, records and education in chinese: Pave the way for medical ai. *International Journal of Medical Informatics*, 177:105173.

Zhengxuan Wu, Atticus Geiger, Joshua Rozner, Elisa Kreiss, Hanson Lu, Thomas Icard, Christopher Potts, and Noah Goodman. 2022. Causal distillation for language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4288–4295.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277.

# KnowLab_AIMed at MEDIQA-CORR 2024: Chain-of-Though (CoT) prompting strategies for medical error detection and correction

**Zhaolong Wu[1]\*, Abul Hasan[2]\*,**

**Jinge Wu[2], Yunsoo Kim[2], Jason P.Y. Cheung[1]†, Teng Zhang[1]†, Honghan Wu[2]†**

[1]Department of Orthopaedics and Traumatology, University of Hong Kong
[2]Institute of Health Informatics, University College London

{wuzl01}@connect.hku.hk,
{cheungjp, tgzhang}@hku.hk,
{a.kalam, jinge.wu.20, yunsoo.kim.23, honghan.wu}@ucl.ac.uk

## Abstract

This paper describes our submission to the MEDIQA-CORR 2024 shared task for automatically detecting and correcting medical errors in clinical notes. We report results for three methods of few-shot In-Context Learning (ICL) augmented with Chain-of-Thought (CoT) and reason prompts using a large language model (LLM). In the first method, we manually analyse a subset of train and validation dataset to infer three CoT prompts by examining error types in the clinical notes. In the second method, we utilise the training dataset to prompt the LLM to deduce reasons about their correctness or incorrectness. The constructed CoTs and reasons are then augmented with ICL examples to solve the tasks of error detection, span identification, and error correction. Finally, we combine the two methods using a rule-based ensemble method. Across the three sub-tasks, our ensemble method achieves a ranking of 3rd for both sub-task 1 and 2, while securing 7th place in sub-task 3 among all submissions.

## 1 Introduction

The rise of Large Language Models (LLMs) such as GPT4 (Achiam et al., 2023), Med-PaLM (Singhal et al., 2023), and LLaMA (Touvron et al., 2023a,b) have inspired investigations into their potential use in automatically analysing Electronic Health Records (EHRs). However, the usefulness of LLMs in clinical settings remains challenging due to the fact that these models are trained on large-scale corpora which may contain inaccuracies, common mistakes, and misinformation (Thirunavukarasu et al., 2023; Ji et al., 2023). To motivate research on the problem of identifying and correcting common sense medical errors in clinical

notes using LLMs, the MEDIQA-CORR (Medical Error Detection Correction) shared tasks are proposed. Herein, we describe our submissions to the shared tasks presenting two methodologies and an ensemble approach using GPT4, all utilising In-Context Learning (ICL) (Brown et al., 2020) in conjunction with Chain-of-thought (CoT) (Wei et al., 2022; Wang et al., 2022b) and reason prompts. The ensemble method achieves accuracies of 69.40% and 61.94% for sub-task 1 and sub-task 2, respectively, while obtaining a BLUERT score of 0.6541 for sub-task 3.

## 2 Shared Tasks and Dataset

### 2.1 Shared Tasks

The MEDIQA-CORR 2024(Ben Abacha et al., 2024a) proposes three sub-tasks:

1. **Binary Classification (sub-task 1)**: To detect whether a clinical note contains a medical error.

2. **Span Identification (sub-task 2)**: To identify the text span (i.e. Error Sentence ID) associated with the error, if a medical error exists in the clinical note.

3. **Natural Language Generation (sub-task 3)**: To generate a corrected text span, if a medical error exists in the clinical note.

### 2.2 Dataset

The training dataset is derived from a single source called as MS Training Set, where as the validation and test datasets are derived from two different sources termed as MS and UW Validation/Test set (Ben Abacha et al., 2024b). The MS Training Set is comprised of 2,189 clinical notes. The MS Validation Set includes 574 clinical notes, while the UW Validation Set includes 160 clinical notes. The

---

Test dataset has in total 926 clinical notes derived from two sources.

# 3 Methods

## 3.1 ICL-RAG- augmented with CoT prompting(ICL-RAG-CoT)

The Chain-of-Thought (CoT) prompting method, which includes a sequence of reasoning steps, has demonstrated enhancements in the problem-solving capabilities of LLMs over standard prompting techniques, particularly in solving mathematical tasks (Wei et al., 2022; Kojima et al., 2022; Yao et al., 2024). Recent studies, such as the one conducted by (Kim et al., 2023), have introduced datasets that incorporate CoT instructions aimed at addressing various Natural Language Processing (NLP) tasks. These tasks include question answering and natural language inference and have been tailored for smaller-scale language models like Flan-T5 (Longpre et al., 2023). Motivated by these developments, we conduct a manual analysis of a subset derived from both the MS Training set and UW Validation set to investigate the prevalent error types within clinical notes. Our examination reveals three broad categories of errors evident in the clinical notes; they are : (1) Diagnosis, (2) Intervention, and (3) Management. Using these categories we construct three separate prompts, shown in Figure 1, that are augmented with ICL examples.

To address the three sub-tasks, our initial approach, referred to as ICL-RAG augmented with CoT prompting (ICL-RAG-CoT), adopts a two-stage prompting methodology with GPT4. For the binary classification and span identification tasks (i.e. sub-task 1 and sub-task 2), we guide GPT4 systematically through a sequence of prompts, each tailored to detect and identify medical errors. The first prompt in the sequence is a standard prompting which tasks the model to detect errors in a clinical note, supplemented with in-context examples. If no medical error is detected, we proceed to prompt GPT4 iteratively by augmenting our CoTs in Figure 1 with ICL examples until an error is identified. Once all CoTs are exhausted, the clinical note is considered error-free. In the second stage, for the NLG task, we prompt GPT4 independently by specifying the predicted incorrect sentence number (i.e., Sentence ID) obtained from the first stage. A prompt template is provided in Appendix A; see Figure 4. In order to generate In-context examples

for prompting LLMs, our methodology incorporates the Retrieval-Augmented Generation (RAG) approach, as proposed by Lewis et al. (2020); Jin et al. (2024). Utilising the e5-large-unsupervised model (Wang et al., 2022a), we transform the MS-Training dataset into a vectorized database. This process involves applying cosine similarity to find the $k$-most similar training instances for each validation and test input. In our experiments we select $k$=4.

## 3.2 ICL-RAG- augmented with reason (ICL-RAG-Reason)

In our second method, referred to as ICL-augmented with reason (ICL-RAG-Reason), we aim to address three sub-tasks simultaneously using a single prompt containing ICL examples and their corresponding reasons for correctness or incorrectness. However, this method requires to prompt the LLM to pre-process the training data separately. Consequently, the ICL-RAG-Reason method begins by prompting GPT4 to generate a brief reason for the correctness or incorrectness of a clinical note from the MS Training set; see Figure 2 for an example. If a note contains an error, we prompt the LLM by concatenating it with the corrected sentence to explain why the clinical note is deemed incorrect. In the case of a correct training example, we prompt the GPT4 to provide us with the clinical characteristics that validate the note's correctness. Thus, we automatically construct reasoning instructions for each MS Training notes. We employ a similar RAG method to ICL-RAG-CoT; however, we utilize OpenAI embeddings [1] to embed all clinical notes across the three datasets. For every input validation and test note, we sample 4 (4-shot) training notes from a pool of its semantically most similar $k$ notes, comprising two correct and two incorrect notes. We augment selected training notes with their *Reasons* for being correct or incorrect and create the final prompt; ; see Figure 5 in Appendix A for an example of prompt template. The ICL-RAG-Reason method samples ICL examples three times to ensure that the model is shown different reasoning paths. This sampling strategy provides us with three different solutions which is resolve by majority voting to ensure consistency and then take the corrected sentence by randomly selecting one from two correct answers.

---

[1] https://platform.openai.com/docs/guides/embeddings

Figure 1: Three types of Chain-of-Thought (CoT) prompts utilised in the ICL-RAG-CoT method: (1), (2), and (3) direct the GPT4 model to focus on intervention, diagnostic, and management errors, respectively.



Figure 2: Reason generation template utlised in the ICL-RAG-Reason method

### 3.3 Ensemble

We integrate the ICL-RAG-CoT and ICL-RAG-Reason methods using a rule-based approach, henceforth termed as the Ensemble method. This approach initially considers predictions generated by the ICL-RAG-CoT method for sub-task 1 and sub-task 2 as correct, while predictions for sub-task 3 from ICL-RAG-Reason are also deemed correct. It then resolves conflicts by identifying clinical notes from the MS and UW Validation and Test sets that are predicted as incorrect by both methods but have differing Error Sentence IDs. Finally, the Ensemble method prompts GPT4 (see see Figure 6 in Appendix A for an example), providing it with

ICL examples, each containing an error, to generate a corrected sentence by specifying the Eorror Sentence ID predicted by the ICL-RAG-CoT.

### 3.4 Evaluation

We evaluate the performances of our methods with the official evaluation scripts on MS and UW Validation Set [2]. Sub-task 1 and 2 are evaluated by using Accuracy. The Natural Language Generation task (i.e. sub-task 3) is evaluated with with ROUGE (Lin, 2004), BERTScore (Zhang et al., 2019), and BLEURT (Sellam, Thibault and Das, Dipanjan and Parikh, Ankur, 2020). We report performances as

---

[2] https://github.com/abachaa/MEDIQA-CORR-2024

Table 1: Main results. Here Acc, AG, R1, and AGC denote Accuracy, Aggregate, ROUGE-1, and AggregateC scores, respectively.

| Method | Sub-task 1 Acc | Sub-task 2 Acc | Sub-task 3 | | | | |
|---|---|---|---|---|---|---|---|
| | | | AG | R1 | BERT | BLEURT | AGC |
| **MS Validation** | | | | | | | |
| ICL-RAG-CoT | **0.6620** | **0.6236** | **0.6350** | **0.6028** | **0.6658** | **0.6363** | **0.5067** |
| ICL-RAG-Reason | 0.6010 | 0.5644 | 0.6165 | 0.5739 | 0.6577 | 0.6178 | 0.4298 |
| Ensemble | **0.6620** | **0.6236** | 0.6184 | 0.5777 | 0.6560 | 0.6215 | 0.5048 |
| **UW Validation** | | | | | | | |
| ICL-RAG-CoT | **0.7437** | **0.6500** | 0.6525 | 0.6701 | 0.6519 | 0.6355 | 0.6091 |
| ICL-RAG-Reason | 0.6875 | 0.5625 | 0.6340 | 0.6180 | 0.6343 | 0.6499 | 0.5350 |
| Ensemble | **0.7437** | **0.6500** | **0.6740** | **0.6762** | **0.6729** | **0.6728** | **0.6174** |
| **Test** | | | | | | | |
| ICL-RAG-CoT | **0.6940** | **0.6194** | 0.6255 | 0.6130 | 0.6399 | 0.6235 | 0.5346 |
| ICL-RAG-Reason | 0.6540 | 0.5837 | 0.6509 | 0.6343 | 0.6703 | 0.6482 | 0.5119 |
| Ensemble | **0.6940** | **0.6194** | **0.6581** | **0.6434** | **0.6767** | **0.6541** | **0.5730** |

the arithmetic mean of ROUGE-1 F1, BERTScore, BLEURT-20. Furthermore, Aggregate scores and AggregateComposite scores, the overall measures across the mentioned metrics, are provided.

## 4 Results

We attain accuracies of 66.20%, 74.37%, and 69.40% on the MS Validation, UW Validation, and Test datasets, respectively, for the binary classification task of error detection (i.e. sub-task 1) using the ICL-RAG-CoT method; see Table 1. For the span identification task, i.e. sub-task 2, the same method achieves accuracies of 62.36%, 65.00%, and 61.94%, respectively. It is noteworthy that the Ensemble method achieves similar accuracies. In the sub-task 3, which involves Natural Language Generation (NLG), the ICL-RAG-CoT method performs less effectively compared to the ICL-RAG-Reason method. It reaches a BLEURT score of 0.6363 on the MS Validation Set. However, our Ensemble approach surpasses the other two methods, achieving BLEURT scores of 0.6729 and 0.6541 for the UW Validation and Test sets, respectively. We observe similar perfomances across other NLG metrics; see Table 1. This is because the reasoning generation method. i.e. ICL-RAG-Reason achieves better performances than the ICL-RAG-CoT method particularly in the NLG task.

## 5 Discussion

Our CoT prompting strategy works well in conjunction with the RAG system. As depicted in Figure 3, across various few-shot settings (e.g., 2, 3, 4, and



Figure 3: Comparison of few-shot examples with or without CoT using ICL-RAG-CoT method on the Binary Classification Task (i.e. sub-task 1) on the MS Validation Set

5-shot settings), the ICL-RAG-CoT method consistently outperforms scenarios where CoT is not employed alongside RAG in the binary classification task. We observe that both the 3-shot and 5-shot settings yield lower performance compared to the 2-shot and 4-shot settings. This disparity suggests that class imbalance in few-shot settings could potentially deteriorate performance. This motivates our selection of 4-shot setting consistently across all our experiments. One of the limitations of our study is that we do not rigorously evaluate the NLG Task, i.e. sub-task 3. Consequently, our overall ranking falls towards the lower end of the top 10 (ranked 7 over-all). While our Ensemble prompting strategy demonstrates a good performance by leveraging reasoning gathered independently from GPT4, there remains scope for improvement. For

instance, further enhancement could be achieved by evaluating the generation of LLMs against clinical and/or biomedical knowledge bases to verify their output.

# 6 Conclusion

We present our submission to the MEDIQA-CORR shared task for medical error detection and correction. Our study evaluates the effectiveness of the GPT4 model through various prompting strategies employing CoT prompting and Reasoning methods. Specifically, our CoT prompting strategies achieve high accuracies in error detection and identification tasks. Additionally, our Ensemble method, which combines outputs from both methods, demonstrates a better performance on the NLG task than the CoT prompting alone. In the future, we aim to explore our approach for other downstream tasks in the clinical domain using open-source LLMs.

# 7 Ethical Statement

Our research employs large language model (LLM) to improve the accuracy of medical records. However, before deploying and utilising the methods proposed with LLM, it is necessary to adhere to ethical and moral principles. The storage and use of patient data must strictly comply with data protection and privacy laws, such as Health Insurance Portability and Accountability Act (HIPAA) and General Data Protection Regulation (GDPR), to ensure that data access is strictly controlled and process transparency is maintained.

# 8 Acknowledgement

# References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems*, 33:1877–1901.

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 Computing Surveys*, 55(12):1–38.

Mingyu Jin, Qinkai Yu, Chong Zhang, Dong Shu, Suiyuan Zhu, Mengnan Du, Yongfeng Zhang, and Yanda Meng. 2024. Health-llm: Personalized retrieval-augmented disease prediction model. *arXiv preprint arXiv:2402.00746*.

Seungone Kim, Se Joo, Doyoung Kim, Joel Jang, Seonghyeon Ye, Jamin Shin, and Minjoon Seo. 2023. The CoT Collection: Improving Zero-shot and Few-shot Learning of Language Models via Chain-of-Thought Fine-Tuning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12685–12708.

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.

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.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The Flan Collection: Designing Data and Methods for Effective Instruction Tuning. In *International Conference on Machine Learning*, pages 22631–22648. PMLR.

Sellam, Thibault and Das, Dipanjan and Parikh, Ankur. 2020. BLEURT: Learning Robust Metrics for Text

Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892. 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, et al. 2023. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.

Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930–1940.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. LLaMA:OpenandEfficient Foundation Language Models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: OpenFoundation and Fine-Tuned Chat Models. *arXiv preprint arXiv:2307.09288*.

Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022a. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022b. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *The Eleventh International Conference on Learning Representations*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, 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.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTScore: Evaluating Text Generation with BERT. In *International Conference on Learning Representations*.

# A  Prompt Templates



**Few-shot Prompt**

Detect whether the text below contains a medical error? If an error is found, set the Error Flag to 1 and output Error Sentence ID; otherwise, set Error Flag to 0 and Error Sentence ID to -1.
Output must follow output format!
The output should not exceed 50 words!
Output format:
Error Flag: ⟨number⟩
Error Sentence ID: ⟨number⟩
Input:
⟨Same as (a)⟩[. . .]
**Please read the example below:**
**Example 1:**
0 A 6-year-old girl is brought to the physician for intermittent fevers and painful swelling of the left ankle for 2 weeks.
1 She has no history of trauma to the ankle.
2 She has a history of sickle cell disease.
3 Current medications include hydroxyurea and acetaminophen for pain.
4 Her temperature is 38.4 C (101.2 F) and pulse is 112/min.
5 Examination shows a tender, swollen, and erythematous left ankle with point tenderness over the medial malleolus.
6 A bone biopsy culture confirms the diagnosis as it grew Streptococcus pneumoniae.
Error Flag: 1
Error Sentence ID: 6
Error Sentence: A bone biopsy culture confirms the diagnosis as it grew Streptococcus pneumoniae.
Corrected Sentence: A bone biopsy culture confirms the diagnosis as it grew Salmonella enterica.
**Example 2:**
⟨another example ⟩[. . .]
**[CoT Part]**

**Answer**

**Error Flag**: 1
**Error Sentence ID**: 5

Figure 4: A template used in ICL-RAG-CoT for the few-shot prompting to solve sub-task 1 and 2.

Figure 5: A template used in ICL-RAG-Reason for the few-shot prompting to solve all sub-tasks simultaneously.

Figure 6: A template used in Ensemble method for the few-shot prompting to solve the sub-task 3.

# PromptMind Team at EHRSQL-2024: Improving Reliability of SQL Generation using Ensemble LLMs

**Satya K Gundabathula**
satyakesav123@gmail.com

**Sriram R Kolar**
sriramrakshithkolar@gmail.com

## Abstract

This paper presents our approach to the EHRSQL-2024 shared task, which aims to develop a reliable Text-to-SQL system for electronic health records. We propose two approaches that leverage large language models (LLMs) for prompting and fine-tuning to generate EHRSQL queries. In both techniques, we concentrate on bridging the gap between the real-world knowledge on which LLMs are trained and the domain-specific knowledge required for the task. The paper provides the results of each approach individually, demonstrating that they achieve high execution accuracy. Additionally, we show that an ensemble approach further enhances generation reliability by reducing errors. This approach secured us 2nd place in the shared task competition. The methodologies outlined in this paper are designed to be transferable to domain-specific Text-to-SQL problems that emphasize both accuracy and reliability.

## 1 Introduction

Text-to-SQL technology translates natural language questions into executable SQL queries that can answer the questions using a provided database. A robust Text-to-SQL system could significantly increase productivity for anyone using databases by providing an easy-to-use natural language interface and reducing the need for expertise in different SQL dialects. These systems are particularly more valuable in domains where SQL knowledge is not essential, such as healthcare, where healthcare professionals like doctors, nurses, and hospital administrators spend a significant amount of time interacting with patient health records stored in databases.

In the era of Large Language Models (LLMs), the field of Text-to-SQL is gaining prominence as these models demonstrate impressive text generation capabilities without the need for fine-tuning.

Introduced in 2017, WikiSQL (Zhong et al., 2017) remains one of the largest datasets for Text-to-SQL and primarily caters to relatively simple queries. Subsequently, the SPIDER (Yu et al., 2018) and MULTI-SPIDER (Dou et al., 2023) datasets were developed. These datasets posed challenges with complex queries that required an understanding of the database schema and support for various languages. BIRD-Bench was introduced to bridge the gap between research and real-world applications by providing large and imperfect databases (Li et al., 2024). These datasets are good representations of typical Text-to-SQL tasks. However, the healthcare domain differs from these generic datasets for the following reasons:

- The questions asked by users maybe highly specialized and specific to the medical field.

- To answer such questions, systems must also possess an understanding of clinical terminology.

- Reliability is of paramount importance as errors can have serious consequences.

These differences present unique challenges for developing a reliable Text-to-SQL system for the healthcare domain. EHRSQL is the first dataset that closely captures the needs of hospital staff and serves appropriately for building and testing Text-to-SQL systems in the healthcare domain (Lee et al., 2022).

Our solution aims to create a Text to SQL system that emphasizes both reliability and accuracy. To achieve this, we divide the task into two phases:

- SQL Generation

- SQL Validation

In the first stage, we focus on SQL generation employing different techniques that include

360

prompting and fine-tuning of LLMs. In both approaches, we use the same prompting strategy to provide the LLM with database information and question-related context. Specifically, we use table schemas combined with sample column values as the database context, and similar questions from the training data as the task context. To identify similar questions from the training data, we employ an embedding-based similarity technique. Then, our goal is to maximize the LLM's ability to generate highly accurate SQL statements utilizing this approach.

There are several reasons why LLMs may fail to generate correct SQL for a given question. Some common reasons include:

- Misinterpretation of question's intent
- Incorrect assumptions or hallucinations about the database's tables or columns
- Inaccuracies or hallucinations in the generated SQL query

Unlike many text generation tasks, Text-to-SQL tasks have a limited number of correct answers but potentially infinite incorrect ones. Inspired by this, we develop a second stage that evaluates the accuracy of the generated SQL. To evaluate the same, we propose an approach for Text-to-SQL that combines the results of multiple robust LLMs. Stronger LLMs often produce consistent outputs despite variations in temperature or other parameters, while smaller LLMs show lower consistency and accuracy. By leveraging the strengths of several robust LLMs, our approach minimizes the number of incorrect SQL queries and enhances the overall robustness and reliability of the Text-to-SQL system.

In the remainder of this paper, we discuss related work, introduce the EHRSQL-2024 task and dataset, and present our two-stage approach. We then provide the results of our experiments and conclude with a summary of our findings.

## 2 Related Work

Prior to the advent of LLMs, the primary focus of research in natural language processing involved refining specialized models using innovative strategies (Wang et al., 2020). Additionally, substantial efforts were devoted to developing sophisticated pre-training methodologies, such as those proposed by STAR (Cai et al., 2022), and exploring decoding strategies, as exemplified by PICARD (Scholak et al., 2021). However, these approaches typically require substantial computational resources and novel techniques.

Large Language Models (LLMs) have been trained extensively on textual data, which has equipped them with vast knowledge. As a result, they exhibit exceptional probabilistic reasoning abilities and can excel at various tasks even without explicit training. Zero-shot prompting techniques, when used with LLMs, have not only narrowed the performance gap on Text-to-SQL but have also surpassed specialized pre-trained or fine-tuned models. Several prompt techniques have been developed based on this zero-shot approach for Text-to-SQL tasks, leading to remarkable achievements on datasets such as SPIDER (Dong et al., 2023), (Liu et al., 2023). Zero-shot generation capabilities can be further enhanced through techniques like in-context learning (ICL) and few-shot prompting.

DIN-SQL (Pourreza and Rafiei, 2023) adopts an in-context learning approach to break down complex SQL generation into manageable subtasks, leading to improved performance on intricate queries. Another technique, retrieval-augmented generation, provides relevant and helpful examples as a few-shot to guide SQL generation (Guo et al., 2024). These approaches have proven effective on general Text-to-SQL tasks but they have not yet been studied rigorously on domain-specific Text-to-SQL problems. Retrieval Augmented Fine-tuning (RAFT) introduces a novel fine-tuning technique that improves the in-domain performance of RAG while integrating domain-specific knowledge (Lewis et al., 2020).

Through our work, we delve into the application of these techniques for the EHRSQL-2024 task.

## 3 Shared Task and Dataset

The EHRSQL-2024 shared task (Lee et al., 2024) is aimed at creating a reliable SQL for answering questions posed in natural language on the MIMIC-IV demo database. The MIMIC-IV database consists of anonymized electronic health records of patients admitted to the Beth Israel Deaconess Medical Center. These records primarily cover two modules: hospital records and ICU records. The publicly available demo version of the database contains a subset of patient records for 100 individuals. It consists of 17 tables from both modules, encompassing a total of 109 columns.

|        | Total Samples | % Unanswerable |
|--------|---------------|----------------|
| **Train** | 5124 | 8.78 % |
| **Valid** | 1163 | 19.95 % |
| **Test**  | 1167 | 19.97 % |

Table 1: EHRSQL-2024 Dataset Statistics

## 3.1 Task Definition

The task aims to accurately generate SQL queries for answerable questions and predict null ($\phi$) for unanswerable ones. Each correct answer earns a score of 1, while incorrect answers receive a score of $-c$, where $c$ is the associated cost. The overall score $RS$ for a cost $c$ and prompt parameter $\theta$ can be expressed as below.

$$RS_\theta(C) = \Sigma_{i=1}^N \mathbb{1}(E(LLM_\theta(Q_i)) = E(GT_i))$$
$$-C * \mathbb{1}(E(LLM_\theta(Q_i)) \neq E(GT_i))$$
$$(1)$$

where $LLM$ represents the model that generates SQL based on a given question $Q_i$. $GT_i$ denotes the ground truth SQL query for the question, and $E$ signifies the executed value of the SQL query when run on a specific database. $\mathbb{1}$ is the indicator function.

The objective of this task is to find the optimal value of $\theta$ at a cost $c$ with respect to the function $RS_\theta(C = c)$.

## 3.2 Dataset

The dataset contains a combination of answerable and unanswerable questions across three subsets: train, valid, and test. Table 1 provides an overview of the composition of each subset.

## 4 Approach

The reliable Text-to-SQL solution is decomposed into two stages as follows.

## 4.1 SQL Generation

To begin, we concentrate solely on boosting the number of accurately produced SQLs without being concerned with reducing the number of incorrect responses. As a result, the objective function becomes:

$$RS_\theta(C = 0) = \Sigma_{i=1}^N \mathbb{1}(E(LLM_\theta(Q_i)) = E(GT_i))$$
$$(2)$$

Maximizing the success and minimizing hallucinations of the LLMs generation task require the provision of the correct context. To achieve this, the following information is essential regarding the task at hand:

- **Database Schema** Comprising tables, columns, and their interrelationships, the database schema serves as a blueprint for the data stored in the database. This information guides the LLM in selecting the appropriate tables and columns.

- **Database Column Values** The actual values stored in the table columns offer additional information. This helps the LLM comprehend and perform operations such as data validation, manipulation, and filtering

- **Training Data** Providing questions (with corresponding SQL answers) similar to the current question aids the LLM in comprehending query formats, syntax, semantics, ambiguity resolution, and bridging the real-world knowledge gap with EHRSQL.

To produce SQL queries for each question, we employ an **in-context learning** approach. Here, the LLM is provided with similar question-SQL pairs, along with the relevant database content. To retrieve similar questions from the training data, we calculate cosine similarities between the evaluation question embedding and the training question embeddings.

We utilize the AnglE model based on BERT, which aims to minimize the angle difference in a complex space (Li and Li, 2023). This approach helps overcome the negative impact caused by the saturation zone of the cosine function. The AnglE embedding model ranks among the top 10 in the Massive Text Embedding Benchmark (MTEB), encompassing eight embedding tasks and 58 datasets (Muennighoff et al., 2022). While AnglE effectively captures the semantic similarity between the intent of questions, it faces challenges in capturing the similarity between clinical terminology, which is also crucial for this task.

To bridge this gap, we combine AnglE embeddings with PubMedBERT embeddings (Gu et al., 2020), trained on the PubMed literature. This allows us to enhance the system's ability to capture clinical terminology. Since embedding similarity scores are not directly comparable across different models due to varying dimensionality, we perform z-normalization to ensure comparability. Algorithm 1 provides an overview of how we retrieve

the top N similar questions for a given question using two different embedding models.

To generate the SQL, we employed ICL and fine-tuning approaches with a consistent prompt template. A shorter version of the final prompt template is provided below for reference. For ICL, we utilized pre-trained models, such as GPT-4 (OpenAI et al., 2024) and Claude-3 Opus (Anthropic and others, 2024), with their default settings for temperature, top_p, and top_k parameters. To fine-tune GPT-3.5, we leveraged the retrieval augmented fine-tuning (RAFT) technique. For each training question, we generated similar questions using the multi-embedding retrieval approach while maintaining the prompt template. Given the limited size of the training set, we conducted fine-tuning with default parameters for only one epoch to prevent overfitting.

---

**Prompt Template**

This is a task converting a natural language question to an SQLite query for a database. You will be provided with the schema of the SQLite database followed by a few examples. You need to generate the SQLite query for a given question and you may return "null" if the question cannot be answered.

```
[Database Tables]
CREATE TABLE patients
(
  row_id int not null primary key, -- 42
  subject_id int not null unique, -- 201
  gender varchar(5) not null, -- 'm'
  dob timestamp(0) not null,
  dod timestamp(0)
);
...
...
[Examples]
[Q]  : How many patients are there in
       total?
[SQL]: SELECT count(subject_id) FROM
       patients
[Q]  : What is the gender of patient
       1002?
[SQL]: SELECT gender FROM patients
       WHERE subject_id = 1002
[Q]  : What is the date of birth of
       patient 1002?
SQL:
```

---

Figure 1 illustrates the complete process of generating SQL using an LLM post training for a given question.

## 4.2 SQL Validation

LLMs have a tendency to generate inaccurate and imaginary responses, regardless of the quality of the context they are provided. Therefore, we implement a second stage using an ensemble approach to eliminate errors generated during the initial generation stage. To verify whether the SQL generated by a two-model or three-model ensemble is correct, each query result is obtained by evaluating it against the database. Subsequently, the results are compared, and a match among all the results qualifies the query as correct.

## 5 Results

In this section, we present the comparison of the reliability scores of the individual models followed by ensembles.

### 5.1 Individual Models

Table 2 presents the reliability scores along with the percentage of unanswered questions for each model i.e. GPT-4, Claude-3 Opus and fine-tuned GPT-3.5.

Overall, Claude-3 Opus answered the most number of questions correctly while also answering them wrong more than others which led to the lowest RS10. GPT-4 appears to be more conservative in generating SQLs and has generated the most nulls. As refraining from generating for unanswerable questions is more important in this task, this led to achieving the best score on RS10 for GPT-4. Although the GPT-3.5 model is significantly less performant than GPT-4, the fine-tuned version brought the generation capability close to the GPT-4 model.

### 5.2 Ensemble

To select the ensemble model that achieves the best performance, we comprehensively evaluated all possible combinations of 2-model and 3-model ensembles. Table 3 provides a detailed comparison of the reliability scores achieved by these various model ensembles.

Among the 2-model ensembles, the combination of fine-tuned GPT-3.5 and Claude-3 Opus achieved the highest RS10 score, outperforming other models. Notably, the ensemble approach involving the fine-tuned GPT-3.5 model exhibited a significant

**Algorithm 1:** Multi-Embedding Retrieval

---

**Data:** train_questions, test_questions, N
**Result:** Similar train questions for test questions
`// Same size as test_questions`
`// Each element contains top N similar train questions and scores`
result ← [];
**foreach** *embed_model* ∈ *M* **do**
  train_embeddings ← create_embeddings(embed_model, train_questions);
  test_embeddings ← create_embeddings(embed_model, test_questions);
  questions, scores ← compute_similarity(test_embeddings, train_embeddings, top_n=N);
  $\mu$ ← compute_mean(scores);
  $\sigma$ ← compute_std(scores);
  z_scores ← (scores $-\mu$) /$\sigma$;
  `// Sort and merge top N current questions with result`
  `// If questions overlap, update with max score`
  result ← sort_and_merge(result, questions, z_scores);

---

Figure 1: SQL Generation Process



| Model | RS0 | RS10 | Unanswered % |
|-------|-----|------|--------------|
| GPT-4 | 88.51 | 40.53 | 25.71 |
| FT GPT-3.5 | 88.08 | 22.96 | 23.14 |
| Opus | 88.94 | 18.68 | 22.28 |

Table 2: Reliability Scores of Individual Models

reduction in errors compared to pre-trained models. This finding suggests that the fine-tuned model produces distinct errors from the pre-trained models, thus maximizing the validation benefits of ensemble approaches. The 3-model ensemble, however, achieved the best RS10 score among all approaches. To illustrate the effectiveness of Ensemble mod-els, Figure 2 demonstrates the reliability scores of top-performing models from the individual, 2-model ensemble, and 3-model ensemble categories. When comparing against a stand-alone model, both 2-model and 3-model ensembles substantially minimize errors and obtain roughly equivalent but large RS10 scores. These results clearly demonstrate that ensemble approaches are effective validation mechanisms for creating reliable and accurate SQL generation systems.

| Ensemble | RS0 | RS10 |
|---|---|---|
| GPT-4 + Opus | 84.57 | 65.72 |
| FT GPT-3.5 + GPT-4 | 84.83 | 71.97 |
| FT GPT-3.5 + Opus | 85.08 | 73.09 |
| All | 82.6 | 74.89 |

Table 3: Reliability Scores of Ensemble Models

Figure 2: Individual vs Ensemble Models



| Prompt Type | Executable % | RS0 | RS10 |
|---|---|---|---|
| No Few-shot | 83.84 | 32.84 | -440.06 |
| One Embedding Few-Shot | 95.89 | 66.98 | 7.65 |
| Two Embeddings Few-Shot | 98.34 | 69.13 | 11.52 |
| Two Embeddings Few-Shot + Column Values | 95.71 | 69.3 | 15.99 |

Table 4: Reliability Scores of GPT-3.5 with Different Prompt

all metrics by a good margin. While adding column values to the few-shot prompt decreased executable queries potentially leading to an increase in RS10, it also showed an improvement in RS0, indicating its usefulness as a signal. Through these experiments, we arrived at the final prompt, which enabled us to develop a highly reliable Text-to-SQL system.

## 7 Conclusion

Our work primarily aims to enhance the reliability of SQL generation, which is of paramount importance for the EHRSQL-2024 shared task. Although in-context learning with advanced LLMs such as GPT-4, Claude-3 Opus, or fine-tuning GPT-3.5 yields excellent RS0, errors still seem inevitable. The model's ability to solve a specific task is heavily influenced by the training data. Repeatedly generating using the same prompt (or) the same model to validate often fails to minimize errors since hallucinations mainly originate from the training data. Fine-tuning GPT-3.5 resulted in different error tendencies compared to pre-trained models, even when using the same prompt. Therefore, ensemble LLMs, particularly those with a fine-tuned model, offer a superior approach for SQL validation, improving robustness and reliability. This approach has also secured us 2nd place in the competition.

## 8 Limitations and Risks

Our approach, while successful in this context, requires careful planning for real-world deployment due to certain limitations. Fine-tuning GPT-3.5 is computationally expensive and necessitates high-quality training data. Ensemble methods, though powerful for validation, introduce trade-offs in terms of cost and complexity. Crucially, it's vital to evaluate potential biases inherited from the

## 6 Ablation Study

To assess the significance of each parameter in the final prompt employed for ICL and fine-tuning, we conduct an ablation study. In these experiments, we focus solely on pre-trained models because fine-tuning experiments are more expensive and time-consuming. To accelerate the process and maintain costs, we leverage GPT-3.5, a compact and less potent yet faster variant of the GPT family. Through these experiments, we extrapolate the efficacy of each parameter for prompting using more robust and advanced models such as GPT-4 and Claude-3 Opus. Table 4 provides the reliability scores obtained by progressively constructing a prompt with varying levels of complexity.

Incorporating few-shot examples in the prompt has substantially improved both the executable queries and reliability scores. This demonstrates the critical role of ICL with few-shot in Text-to-SQL tasks, particularly in the context of EHRSQL. The one-embedding few-shot experiment employs non-medical AnglE embeddings (Li and Li, 2023), while the two-embeddings few-shot additionally leverages PubMedBERT (Gu et al., 2020). It is evident that adding medical embeddings enhances

LLM's training data to ensure fair and reliable performance in practical applications.

# References

Anthropic and others. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf. Online; accessed March 2024.

Zefeng Cai, Xiangyu Li, Binyuan Hui, Min Yang, Bowen Li, Binhua Li, Zheng Cao, Weijie Li, Fei Huang, Luo Si, and Yongbin Li. 2022. STAR: SQL guided pre-training for context-dependent text-to-SQL parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1235–1247, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, lu Chen, Jinshu Lin, and Dongfang Lou. 2023. C3: Zero-shot text-to-sql with chatgpt. *Preprint*, arXiv:2307.07306.

Longxu Dou, Yan Gao, Mingyang Pan, Dingzirui Wang, Wanxiang Che, Dechen Zhan, and Jian-Guang Lou. 2023. Multispider: Towards benchmarking multilingual text-to-sql semantic parsing. In *AAAI Conference on Artificial Intelligence*.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.

Chunxi Guo, Zhiliang Tian, Jintao Tang, Shasha Li, Zhihua Wen, Kaixuan Wang, and Ting Wang. 2024. Retrieval-augmented gpt-3.5-based text-to-sql framework with sample-aware prompting and dynamic revision chain. In *Neural Information Processing*, pages 341–356, Singapore. Springer Nature Singapore.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Patrick S. H. 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.

Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. 2024. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36.

Xianming Li and Jing Li. 2023. Angle-optimized text embeddings. *arXiv preprint arXiv:2309.12871*.

Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S. Yu. 2023. A comprehensive evaluation of chatgpt's zero-shot text-to-sql capability. *Preprint*, arXiv:2303.13547.

Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.

OpenAI, Josh Achiam, et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *arXiv preprint arXiv:2304.11015*.

Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. PICARD: Parsing incrementally for constrained auto-regressive decoding from language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9895–9901. Association for Computational Linguistics.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2020. RAT-SQL: Relation-aware schema encoding and linking for text-to-SQL parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7567–7578, Online. Association for Computational Linguistics.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Brussels, Belgium. Association for Computational Linguistics.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *CoRR*, abs/1709.00103.

# PromptMind Team at MEDIQA-CORR 2024: Improving Clinical Text Correction with Error Categorization and LLM Ensembles

**Satya K Gundabathula**
satyakesav123@gmail.com

**Sriram R Kolar**
sriramrakshithkolar@gmail.com

## Abstract

This paper describes our approach to the MEDIQA-CORR shared task, which involves error detection and correction in clinical notes curated by medical professionals. This task involves handling three subtasks: detecting the presence of errors, identifying the specific sentence containing the error, and correcting it. Through our work, we aim to assess the capabilities of Large Language Models (LLMs) trained on a vast corpora of internet data that contain both factual and unreliable information. We propose to comprehensively address all subtasks together, and suggest employing a unique prompt-based in-context learning strategy. We will evaluate its efficacy in this specialized task demanding a combination of general reasoning and medical knowledge. In medical systems where prediction errors can have grave consequences, we propose leveraging self-consistency and ensemble methods to enhance error correction and error detection performance.

## 1 Introduction

With rapid advancements in Natural Language Processing (NLP), we are witnessing a surge of its applications across various domains, including healthcare. Incorporating NLP in clinical settings brings about a multitude of advantages. It enhances clinical decision-making through advanced support, making health information more accessible, streamlines documentation, and accelerates research initiatives (Hossain et al., 2023). These developments contribute to improved patient care, reduced healthcare costs, and alleviated physician burnout.

NLP for healthcare applications pose inherent challenges due to the need for medical expertise. However, advancements in LLMs trained on internet data including medical information have significantly enhanced their knowledge and reasoning capabilities, enabling them to tackle more complex problems in the healthcare domain involving text

processing and generation. Few recent applications in healthcare include information extraction, question answering, summarization, and translation, all while comprehending intricate medical knowledge (Nazi and Peng, 2023). Despite these advancements, safety and accuracy remain major concerns as training data may contain unreliable and misleading information that could have adverse effects if not handled appropriately. Nevertheless, the effective utilization of these LLMs has the potential to revolutionize healthcare and bring immense benefits to society (Clusmann et al., 2023).

In the healthcare industry, there is a need for automated systems capable of efficiently analyzing and interpreting clinical texts that improves patient's safety, quality of care and costs. Processing the clinical texts presents a unique and significant challenge due to the complexities introduced by medical jargon, abbreviations, syntactic variations, and context-specific nuances. The MEDIQA-CORR shared task (Abacha et al., 2024), part of the ClinicalNLP 2024 workshop, seeks to address this issue of identifying and correcting (common sense) medical errors found in clinical notes.

The MEDIQA-CORR shared task involves three subtasks: detecting errors in clinical notes, identifying specific error sentences, and correcting those sentences. Our approach involves tackling all three subtasks simultaneously using a single prompt for LLMs invoking chain-of-thought. By doing this, we seek to assess the complex reasoning capabilities of LLMs in the clinical domain.

First, we analyze the dataset and curate a list of the most common types of errors in clinical notes. We then utilize this information to create task specific instructions for the LLM. We leverage contemporary LLMs using in-context learning (ICL) with a few-shot approach. Since publicly accessible LLMs are instruction-tuned models, considering the approach of directing them towards assessing the clinical note based on specific error

types maximizes the objective while only utilizing a few training examples.

By employing single-prompt ICL approaches to solve end-to-end tasks, we pave the way forward for building more complex healthcare applications using simple yet intuitive strategies leveraging the capabilities of advanced LLMs.

In addition to the task prompt, an LLM's performance on a particular task is mainly influenced by two factors: training techniques and underlying training data. Consistency, which measures how frequently an LLM produces the same output given the same instructions, can be viewed as another dimension that can affect the quality of an LLM. Leveraging self-consistency can significantly improve the accuracy of an LLM, particularly for complex reasoning tasks (Wang et al., 2023). In healthcare datasets, where hallucinations in LLMs can occur more frequently due to training on factually unverified data, this could lead to serious problems. While self-consistency is one approach to obtaining more accurate results, LLM ensemble, which has not yet been fully explored, presents a promising opportunity. We validate the results of each LLM by using the output of other LLMs that are trained on different corpora. In our approach, we investigate both self-consistency and ensemble concepts.

The remainder of the paper includes related work, the MEDIQA-CORR task and dataset, our approach, results and conclusion.

## 2 Shared Task and Dataset

The shared task focuses on leveraging LLMs for the following three subtasks:

- **Binary Classification:** Detect if the text from a clinical note includes a medical error.

- **Span Identification:** Identify the text span (in the sentence) associated with the error, if an error exists.

- **Natural Language Generation:** Generate the corrected text, if an error exists.

### 2.1 Subtasks and Metrics

### 2.1.1 Error Detection

For each clinical text, the prediction is assigned a value of 1 when an error is detected, and 0 if no error is detected. Given the binary nature of this classification task, accuracy serves as the primary metric for performance evaluation.

### 2.1.2 Error Span Identification

Each clinical text comprises sentences associated with unique IDs. The subtask involves predicting the error ID, which is an integer between 0 and the highest sentence ID. If no error is detected, the prediction should be -1. The primary evaluation metric is accuracy, calculated based on all samples, including those with and without errors.

### 2.1.3 Correct Sentence Generation

If a model identifies an error sentence in the previous subtasks, it should also generate a corrected sentence as prediction for this subtask. Here, the full corrected sentence is evaluated against the ground truth sentence for measuring the performance. Clinical note generation tasks are challenging to evaluate automatically due to the large number of possible correct answers. Metric ensembles (Abacha et al., 2023) have been shown to outperform individual state-of-the-art measures, such as ROGUE for such tasks. The evaluation metric for this subtask is computed as an unweighted average of the following three scores:

- **ROUGE-1F** measures the similarity between system-generated and human-written texts by measuring the overlap of unigrams (Lin, 2004). It uses the F-1 score to assess the quality of the generated sentence.

- **BERTScore** leverages contextual word embeddings obtained from BERT models to assess the similarity between a candidate sentence and a reference sentence (Zhang* et al., 2020). In this context, it signifies the F-1 score of the semantic similarity performed using the DeBERTa XL model (He et al., 2021).

- **BLEURT-20** is a learned metric trained on human ratings that aims to better correlate with human judgments for measuring quality compared to traditional BLEU (Sellam et al., 2020).

### 2.2 Dataset

The train data consists of clinical texts from MS data while the valid and test data contains MS and UW collections. Each entry in the datasets includes a text, its ID and sentences as inputs. Table 1 shows the composition of the dataset.

## 3 Approach

We propose to tackle all subtasks concurrently within a single prompt for the following reasons:

| Dataset Type | # Samples | % of Error Samples |
|---|---|---|
| MS Training | 2189 | 55.69 |
| MS Validation | 574 | 55.57 |
| UW Validation | 160 | 50.00 |
| MS + UW Test | 925 | 51.35 |

Table 1: MEDIQA-CORR Dataset

- **Comprehensive Evaluation:** To enable performance evaluation on a complex task, rather than assessing them on isolated, simpler tasks. This provides a more holistic view of the LLMs' capabilities.

- **Efficiency Optimization:** To minimize inference costs and developmental efforts by eliminating the need for multiple models (or) sequential processing. It streamlines the process, making it more efficient and cost-effective.

- **Ease of Adoption:** To alleviate the adoption burden in practical applications and facilitate seamless upgrades to more advanced LLMs amidst the rapid technological advancements.

Publicly accessible LLMs are models that are fine-tuned to follow instructions, with the aim of performing user-defined instructions as accurately as possible. The success of the task then depends on the quality of the instructions provided and the model's ability to follow them effectively. Our approach focuses on refining the instructions for the LLM to facilitate comprehensive learning of all subtasks using ICL. We initiate this process by analyzing error types in the dataset followed by curating the prompt and inferencing with different LLMs.

### 3.1 Error Analysis

In the MEDIQA-CORR task, the definition of an error is loosely defined and can be interpreted differently by humans or systems without examining the dataset. To address this, we perform error type classification in clinical texts by extracting error sentences and corresponding corrected sentences from the training data. We create an LLM prompting task to broadly categorize the entities modified from the error sentence to the corrected sentence within the clinical domain. We utilize GPT-3.5 for categorizing the errors and cluster these generated free-form categories into a manually defined set. This categorization results in the identification of

various error types as depicted in figure 1. Finally, we use the well-defined error categories for the rest of the task while handling "Others" category discreetly.



Figure 1: Error Types Extracted from Training Data

### 3.2 Prompt Curation

When prompted to identify errors in clinical texts without being specific, LLMs may introduce biases from their training data and flag non-critical errors adhering to their own standards of composing clinical notes. To mitigate this, we propose converting an abstract definition of a clinical note error into a concrete and approximate one by analyzing and categorizing the errors. Consequently, we expand the original task from identifying the error to include error classification, which facilitates chain-of-thought for the LLM. We conduct ablation studies to demonstrate the effectiveness of these techniques and present in section 4.

To illustrate the task more thoroughly, we incorporate reasoning within the task prompt for the LLM. Through this, we aim to provide more generalizable patterns for detecting and identifying errors. Additionally, it adds explainability to the LLM systems, which is crucial for real-world applications. Finally, we adopt a few-shot approach utilizing random samples from the training data to teach the LLM how to detect, identify, classify errors, provide reasoning and demonstrate the corrected text. We utilize the same samples for few-shot throughout the task as we

need to manually generate the additional fields such as error category and reason. The designed task prompt is provided as below. Note that this prompt is a tailored version for demonstration purposes. Finally, errors that fall into the "Others" category are processed as "No Error" as they are unimportant for this task.

---

**Prompt Template**

In this task, you will be given a clinical report presented as sentences while each sentence is associated with a sentence number. Now, you need to go through the report line by line and identify if there is any error in the sentence. The error can fall into one of the following main categories:
1. Medications
2. Medical Conditions, Virus or Bacteria
3. Reports, Diagnosis and Monitoring
4. Clinical Procedures and Treatments
5. Clinical Plans and Recommendations
6. Medical Devices
7. Others including clarity/improper usage of terminology

The error can be identified by looking at the entities present in each sentence and check if these entities fall into one of the aforementioned categories and validate if the entire report is correct with this entity. If there is an error, you need to output details as shown in the examples. Use all your medical knowledge and make the right judgements. Here are a few examples for your understanding:

/* Five random samples from training data with manually curated error category and reason */
**Example Clinical Report:**
0 Mr. <Name> is admitted ..
1 He has a surgical ..
2 He is also being managed ..
**Output:**
{
"Error Sentence ID": 1,
"Error Category": "Medical Devices",
"Reason": "The device should be .."
"Corrected Sentence": "He has a surgical .."
}
...

---

...
**Test Clinical Report:**
0 A 45 year old woman ..
1 She is experiencing ..
2 She had prior examamination ..
**Output:**

---

## 3.3 Model Selection

Using the designed prompt, we utilize the following LLMs for performing the task:

- **GPT-3.5:** A model from the OpenAI's generative pre-trained transformer (GPT) family that can understand as well as generate natural language or code (Ye et al., 2023).

- **GPT-4:** Latest model from the GPT family with broader general knowledge and advanced reasoning capabilities (OpenAI et al., 2024).

- **Claude-3 Opus:** Anthropic's largest model, released in Feb 2024, which sets new industry benchmarks across a wide range of cognitive tasks and outperforms its peers on most of the common evaluation benchmarks for AI systems (Anthropic and others, 2024).

Due to its affordability, speed, and reliability, GPT-3.5 is an excellent choice for experimentation. As a result, we employed GPT-3.5 for error analysis and prompt design, reserving the more advanced GPT-4 and Claude-3 Opus models for the final test runs.

## 3.4 Enhancing Robustness

Due to the potential limitations such as hallucinations and inconsistent results, which can affect the quality of the LLMs, we investigate two concepts to improve performance on the subtasks: self-consistency and ensemble. The Claude-3 Opus model has slower speed, higher cost, and stricter token limits compared to GPT-4. Therefore, we utilize GPT-4 for self-consistency by generating four outputs per test sample and aggregate them, while only generating one output per test sample for Claude-3 Opus.

To enhance the quality of predictions, we combine the results from both models to predict the outputs for all three subtasks. Figure 2 provides a visual representation of the overall process. The results aggregator module validates and combines the outputs i.e. predicted error flag, predicted error sentence ID and corrected sentence, from GPT-4 and Claude-3 Opus models to generate the final error flag, error sentence ID and corrected sentence.

Figure 2: LLM Ensemble Approach for MEDIQA-CORR task

| Model | Prompt | Task-1 Accuracy | Task-2 Accuracy |
|-------|--------|-----------------|-----------------|
| GPT-3.5 | No error categories | 48.75% | 22.5% |
| GPT-3.5 | Error categories | 58.44% | 38.55% |
| GPT-4 | Error categories | 63.07% | 58.17% |

Table 2: Performance improvement with error categorization in prompt using GPT-3.5 and assessing GPT-4 performance

| Model | Task-1 Accuracy | Task-2 Accuracy | Task-3 Aggregate Score |
|-------|-----------------|-----------------|------------------------|
| GPT-4 | 62.05% | 56.43% | 0.6172 |
| Claude-3 Opus | 62.59% | 58.48% | 0.6669 |

Table 3: Comparison of GPT-4 and Claude-3 Results

| Model | Task-1 Accuracy | Task-2 Accuracy | Task-3 Aggregate Score |
|-------|-----------------|-----------------|------------------------|
| GPT-4 with consistency (Majority=3/4) | 62.91% | 59.45% | 0.6390 |
| GPT-4 + Claude-3 Opus (Majority=4/4) | 62.16% | 60.86% | **0.7865** |
| GPT-4 + Claude-3 Opus (Majority=3/4) | **63.78%** | **62.48%** | 0.7492 |

Table 4: Self-consistent GPT-4 and its ensemble with Claude-3 Opus Results

## 4 Ablation Study

We begin by presenting the results obtained from incorporating error categorization into the final prompt, which demonstrates an improvement in performance on both error detection and text span identification tasks. In order to make the comparison, we utilize GPT-3.5 with prompts including and excluding error categorization. Additionally, we assess the performance of GPT-4 to ascertain the extent to which it surpasses GPT-3.5 for the finalized prompt. The results obtained using the combined MS and UW validation sets (during the development phase) are presented in Table 2.

The results indicate that by integrating error categorization which initiates an intermediate chain-of-thought, results in a significant performance boost of nearly 10% and 16% for Task-1 and Task-2, respectively. Additionally, GPT-4 outperforms GPT-3.5, confirming its enhanced reasoning capabilities. These advancements make GPT-4 a preferred choice for final test runs.

## 5 Results

We present the results of individual models first followed by incorporating self-consistency and ensembles on the test data. The performance of GPT-4 and Claude-3 Opus using the final prompt on all

three subtasks is presented in Table 3.

Claude-3 Opus surpasses GPT-4 in error detection and significantly in error sentence identification. However, GPT-4 tends to be more verbose during error correction, leading to lower scores in metrics such as ROUGE, BERT, and BLEURT. Although Claude-3 Opus exhibits superior performance, its daily token limit, slower inference, and shorter test phase duration hinder its usability for self-consistency. Therefore, we employ GPT-4 to demonstrate the performance enhancement using self-consistency in large language models (LLMs). Additionally, we ensemble the self-consistent GPT-4 with Claude-3 Opus to showcase further improvement. Table 4 presents the results for all subtasks using the aforementioned methods:

The majority ratio x/y for GPT-4 results is a measure of how often the model produces the same result for a given input. For example, a majority ratio of 3/4 means that at least three out of the four

results should be the same to qualify the predicted error sentence ID as correct. In the ensemble approach, a prediction is considered correct if the self-consistent GPT-4 result matches the Claude-3 Opus result. Otherwise, the prediction is considered "no error". To select the best corrected sentence from the ensemble, the ROGUE score is used to estimate the distance between each corrected sentence and the error sentence. The sentence with the highest ROGUE score is used for the evaluation because LLMs tend to generate verbose corrected sentences which decreases the aggregate scores. In our experiments, using a majority ratio of 3/4 for GPT-4 in the ensemble resulted in the best Task-1 and Task-2 performances. Using a majority ratio of 4/4 resulted in the best Task-3 aggregate score. Our best subtask-3 score is ranked 2nd among all participants in the competition, and our best subtask-2 performance is among the top 3 according to the reported scores (Ben Abacha et al., 2024).

## 6 Related Work

In recent years, there has been a surge of research exploring the potential of prompt engineering techniques with large language models (LLMs) in healthcare. These techniques have shown promising results in various healthcare tasks, often achieving state-of-the-art performances (Zhou et al., 2023), (He et al., 2023). One notable study, Med-Prompt, highlighted several research directions demonstrating the power of prompt exploration for LLMs (Nori et al., 2023). LLMs exhibited impressive knowledge and reasoning abilities, tackling various tasks effectively. These advancements showcase the potential of LLMs in healthcare, offering new opportunities for leveraging language models to address healthcare challenges.

Evaluating common sense reasoning is essential for computer systems, as it impacts language comprehension, communication reliability, and general task performance. SemEval-2020 ComVE aims to address general common sense questions and seeks logical justification for correct responses, assessing reasoning abilities. Pretrained language models naturally acquire common sense through training on vast word tokens obtained from the web (Wang et al., 2020). MEDIQA-CORR, specifically tailored to identify and correct errors in clinical notes, offers a valuable resource for evaluating pretrained LLMs in medical common sense reasoning. Inspired by prompt-based explorations, our research also focuses on utilizing pretrained LLMs to assess reasoning capabilities in medical common sense scenarios.

## 7 Conclusion

Our research demonstrates that incorporating error categorization into the prompt enhances the performance of large language models (LLMs) in detecting, identifying, and classifying clinical note errors. By initiating an intermediate chain-of-thought, this approach facilitates better reasoning and aids the LLM in providing more accurate and explainable results. Furthermore, our findings suggest that self-consistency and ensembles can further enhance the robustness and performance of LLMs on these tasks. These advancements pave the way for the development of more reliable and interpretable AI systems for clinical documentation analysis, ultimately contributing to improved healthcare outcomes.

## 8 Limitations and Risks

While promising, our approach has limitations. LLMs trained on general data may lack specific medical knowledge, potentially leading to misinterpretations and inaccurate corrections. Despite mitigation efforts, the risk of hallucinations and inconsistencies in LLM outputs remains a concern. Additionally, the effectiveness of our approach relies heavily on prompt engineering, which requires expertise and may not be easily generalizable.

The black box nature of LLMs also presents challenges in terms of explainability and building trust in medical contexts. To mitigate these limitations, continuous improvements of training data, more robust evaluation metrics, and human oversight are crucial. Further research is needed to explore the full potential and limitations of LLMs in healthcare, ensuring their safe and responsible application for improved patient care.

## 9 Ethical Considerations

The use of LLMs for medical error detection and correction raises significant ethical concerns. Potential biases in training data and algorithms must be carefully mitigated to prevent propagating existing healthcare disparities and ensure fairness. Transparency in how prompts are designed and how the LLM reaches its decisions is vital for building trust and ensuring accountability. Additionally, robust data security and de-identification practices

are paramount for protecting sensitive patient information.

It is essential to remember that LLMs should serve as tools to augment the expertise of healthcare professionals, not replace them. Clear lines of responsibility, ongoing human oversight, and continuous research and collaboration are necessary to address these ethical challenges. This will ensure the responsible use of LLMs and their positive contribution to improved healthcare outcomes.

# References

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Asma Ben Abacha, Wen wai Yim, George Michalopoulos, and Thomas Lin. 2023. An investigation of evaluation metrics for automated medical note generation. *ArXiv*, abs/2305.17364.

Anthropic and others. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf. Online; accessed March 2024.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Jan Clusmann, Fiona R Kolbinger, et al. 2023. The future landscape of large language models in medicine.

Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. 2023. A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics. *Preprint*, arXiv:2310.05694.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*.

Elias Hossain, Rajib Rana, Niall Higgins, Jeffrey Soar, Prabal Datta Barua, Anthony R. Pisani, and Kathryn Turner. 2023. Natural language processing in electronic health records in relation to healthcare decision-making: A systematic review. *Computers in Biology and Medicine*, 155:106649.

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.

Zabir Al Nazi and Wei Peng. 2023. Large language models in healthcare and medical domain: A review. *Preprint*, arXiv:2401.06775.

Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoifung Poon, Tao Qin, Naoto Usuyama, Chris White, and Eric Horvitz. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *Preprint*, arXiv:2311.16452.

OpenAI, Josh Achiam, et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

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.

Cunxiang Wang, Shuailong Liang, Yili Jin, Yilong Wang, Xiaodan Zhu, and Yue Zhang. 2020. SemEval-2020 task 4: Commonsense validation and explanation. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 307–321, Barcelona (online). International Committee for Computational Linguistics.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.

Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models. *Preprint*, arXiv:2303.10420.

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

Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu Zou, Jinfa Huang, Jinge Wu, Yiru Li, Sam S. Chen, Peilin Zhou, Yining Hua Junling Liu, Chengfeng Mao, Xian Wu, Yefeng Zheng, Lei Clifton, Zheng Li, Jiebo Luo, and David A. Clifton. 2023. A survey of large language models in medicine: Progress, application, and challenge. *arXiv preprint arXiv:2311.05112*.

# Maven at MEDIQA-CORR 2024: Leveraging RAG and Medical LLM for Error Detection and Correction in Medical Notes

**Suramya Jadhav**\* , **Abhay Shanbhag**\* , **Sumedh Joshi**\* , **Atharva Date**\* , **Sheetal Sonawane**

SCTR's Pune Institute of Computer Technology

{2018suramyajadhav,abhayshanbhag0110,sumedhjoshi463,atharva2718}@gmail.com,
sssonawane@pict.edu

## Abstract

Addressing the critical challenge of identifying and rectifying medical errors in clinical notes, we present a novel approach tailored for the MEDIQA-CORR task @ NAACL-ClinicalNLP 2024, which comprises three subtasks: binary classification, span identification, and natural language generation for error detection and correction. Binary classification involves detecting whether the text contains a medical error; span identification entails identifying the text span associated with any detected error; and natural language generation focuses on providing a free text correction if a medical error exists. Our proposed architecture leverages Named Entity Recognition (NER) for identifying disease-related terms, Retrieval-Augmented Generation (RAG) for contextual understanding from external datasets, and a quantized and fine-tuned Palmyra model for error correction. Our model achieved a global rank of **5** with an aggregate score of **0.73298**, calculated as the mean of ROUGE-1-F, BERTScore, and BLEURT scores.

## 1 Introduction

Clinical notes typically include details about the patient's medical history, symptoms, physical examinations, diagnostic tests, treatments administered, and any other relevant information related to the patient's health status and care plan.

Accurate documentation is crucial for patient care, as errors in clinical notes can lead to misdiagnosis, improper treatment, and potential harm to patients. By automating the process of error detection and correction, healthcare providers can ensure the integrity and reliability of patient records, ultimately improving the quality of care delivered. Research indicates that a substantial proportion of adverse events in healthcare settings are due to errors in documentation, highlighting the need for effective error detection and correction mechanisms.

In this task of Medical Error Detection Correction Ben Abacha et al., 2024. We seek to address the problem

of identifying and correcting medical errors in clinical notes. This task had 3 subtasks. In subtask 1 (Binary Classification) researchers had to detect whether the clinical notes included a medical error or not. Subtask 2 named Span Identification was to identify the text span associated with the error if a medical error exists. Subtask 3 (Natural Language Generation) was specifically to provide error-free text after making corrections if a medical error exists.

In our approach, we initially conducted Named Entity Recognition (NER) using GEMINI to identify words representing diseases or pathogens or suggestions in the text. After masking these identified words, we implemented the Retrieval-Augmented Generation (RAG) model on textbooks and external datasets. If the RAG score fell below a certain threshold, we passed the input to our model, which was made by using 4-bit quantization on Palmyra 20b and then fine-tuned the quantized Palmyra model using the QLoRA technique on MEDQA data Jin et al., 2020. If the word provided by Palmyra or the RAG model matched the word detected by NER, no error was detected. Otherwise, if a different word was obtained, it was replaced with the masked word identified by NER. Finally, the error sentence is mapped with the sentence ID to get the output in the desired format. This approach helped us in getting a Global Rank 5 with an Aggregate Score of 0.73298. The Aggregate score is calculated as the mean of ROUGE-1-F, BERTScore, and BLEURT . Our model achieved R1F, BERTSCORE, and BLEURT scores of 0.70306, 0.74372 and 0.75217 respectively.

The rest of the paper is organized as follows: Review of related work and background information in Sections 2 and 3 respectively, to provide context for our study. Following this, we elucidate the system architecture in Section 4 and describe the experimental setup in Section 5. Subsequently, we present our findings in Section 6, discuss limitations encountered in Section 7, and propose avenues for future research in Section 8. Finally, we have concluded our discussion in Section 9.

## 2 Background

The med dataset provided by organizers had 2 types of clinical notes - MS and UW. Upon meticulous examination of the datasets , it became clear that the medical dataset which was divided into MS and UW clinical notes presented some unique difficulties. The MS sub-

---

\*first author, equal contribution

| | MS | UW |
|---|---|---|
| **Text ID** | ms-val-108 | uw-val-51 |
| **Text** | A 3175-g (7-lb) female newborn is delivered at term. Initial examination shows a flat perineal Colonic atresia is confirmed when dark green discharge is coming out of the vulva. | Mr. <NAME/> has been noted to have documentation of thrombocytopenia on <DATE/> in the Medicine note. Plt 101 on admission. Thrombocytopenia was present on admission (POA). |
| **Sentences** | 0 A 3175-g<br>1 (7-lb) female newborn is delivered at term.<br>2 Initial examination shows a . . . .<br>3 Colonic atresia is confirmed . . . | 0 Mr. <NAME/> has been noted to have documentation of thrombocytopenia on <DATE/> in the Medicine note.<br>1 Plt 101 on admission.<br>2 Thrombocytopenia was present on admission (POA). |
| **Error Flag** | 1 | 0 |
| **Error Sentence ID** | 3 | -1 |
| **Error Sentence** | Colonic atresia is confirmed. . . | NA |
| **Corrected Sentence** | Imperforate anus is confirmed when dark green discharge is coming out of the vulva. | NA |
| **Corrected Text** | A 3175-g (7-lb) female . . . opening. Imperforate anus is confirmed when dark green discharge is coming out of the vulva. | NA |

Table 1: Dataset Glimpse

set, which came from Microsoft, had incredibly small flaws. So much so that a great deal of faults appeared to be subtle, making it difficult for the physicians on our team to recognize them. Yet, it was clear from closely examining the training set's corrected text that the corrections frequently represented ideal completions.

The UW subset, which came from University of Washington, on the other hand, showed a distinct scene. This subset of clinical notes seemed to more closely resemble real-world situations, which made errors easier to identify in them.

MS dataset was split into train (2189) and val (574), and UW into val dataset (160). The testing data was a mixture of MS and UW formats.

The dataset is in CSV format and consists of labeled text data. Each row represents a unique input text and includes columns named Text ID, Text, Sentences, Error Flag, Error Sentence ID, Error Sentence, Corrected Sentence, and Corrected Text. The Text column contains the complete text, while the Sentences column divides the text into individual sentences with corresponding IDs starting from 0. The Error Flag column indicates whether there is an error in the text, with 0 representing no error and 1 representing an error. If there is an error, the Error Sentence ID column specifies the ID of the sentence containing the error, and the Error Sentence column provides the erroneous part of the text containing the error. The Corrected Sentence column contains the error-corrected version of the sentence, and the Corrected Text column includes the complete text with corrected sentences. When there is no error, Error Flag is 0, Error Sentence ID is -1, and the Error Sentence, Corrected Sentence, and Corrected Text columns contain "NA" values. This structured format facilitates

error detection and correction tasks within the dataset. Table 1 offers a glimpse into MS and UW datasets.

The MEDQA dataset is a collection of question-answer pairs related to the medical field specifically derived from professional medical board exams, like the United States Medical Licensing Examination (USMLE). It covers a wide range of medical topics and is available in three languages: English, Simplified Chinese, and Traditional Chinese.

**Question-Answer Pairs**: The dataset consists of multiple-choice questions along with their corresponding answers. The number of questions varies depending on the language:
English: 12,723 questions
Simplified Chinese: 34,251 questions
Traditional Chinese: 14,123 questions

**Medical Textbooks**: The dataset also provides access to a large corpus of medical textbook content to aid models in comprehending the medical context for answering the questions.

For this task we used the just the English QA corpus. Here's an example of a question-answer pair in MEDQA dataset.

**Question** A 55-year-old female patient presents with a chief complaint of progressive shortness of breath over the past 6 months. She denies chest pain, cough, fever, or chills. On physical exam, her vital signs are normal. Her lungs are clear to auscultation bilaterally.What is the most likely diagnosis for this patient's shortness of breath?
**Options**
A. Heart failure
B. Asthma
C. Chronic obstructive pulmonary disease (COPD)

D. Pneumonia

**Answer-idx :** C

# 3 Related Work

Zhu et al., 2024 unveils REALM, a Retrieval-Augmented Generation framework, addressing limitations in existing clinical predictive models by enhancing multimodal Electronic Health Records (EHR) representations. Integrating clinical notes and time-series EHR data, REALM leverages Large Language Models (LLM) and GRU models for encoding, while incorporating external knowledge from a labeled knowledge graph (PrimeKG). By aligning with clinical standards, the framework eliminates hallucinations and ensures consistency, culminating in an adaptive multimodal fusion network. Extensive experiments on MIMICIII tasks demonstrate REALM's superior performance, highlighting its effectiveness in refining multimodal EHR data utilization and enhancing nuanced medical context for informed clinical predictions.

Elgedawy et al., 2024 presented a conversational interface powered by large language models (LLMs) for efficiently accessing information within clinical notes. Utilizing Langchain framework and transformer-based models, users can interactively query and retrieve relevant details from unstructured clinical data. Evaluation experiments, including advanced language models and semantic embedding techniques, demonstrate promising results, with Wizard Vicuna showing the highest accuracy despite computational demands. Model optimization techniques, such as weight quantization, significantly improve inference latency. However, challenges like model hallucinations and limited evaluation across diverse medical cases remain, indicating avenues for future research in enhancing clinical decision-making through AI-driven approaches.

Singhal et al., 2023 outlines Med-PaLM 2, a significant advancement in medical question answering, achieving an impressive accuracy of 86.5 % on the MedQA dataset. Compared to its predecessor, Med-PaLM, which scored 67.2% on the same dataset, Med-PaLM 2 represents a substantial improvement. By leveraging enhancements in base large language models (LLMs), domain-specific fine-tuning, and novel prompting strategies, Med-PaLM 2 demonstrates promising progress towards attaining physician-level performance in medical question answering across various datasets, including MedQA, PubmedQA Jin et al., 2019, MMLU, and MedMCQA Pal et al., 2022.

Jin et al., 2020 elucidates MEDQA, the inaugural free-form multiple-choice OpenQA dataset for medical problem-solving, sourced from professional medical board exams in English, simplified Chinese, and traditional Chinese. With question counts of 12,723, 34,251, and 14,123 across the three languages respectively, MEDQA provides a robust benchmark. Despite employing both rule-based and neural methods, even the most advanced model achieves only 36.7%, 42.0%,

and 70.1% test accuracy on English, traditional Chinese, and simplified Chinese questions. MEDQA poses significant challenges to current OpenQA systems, encouraging the NLP community to develop more robust models for medical applications.

Chen et al., 2023 introduces MEDITRON, an open-source suite of Large Language Models (LLMs) tailored for the medical domain, ranging from 7B to 70B parameters. Leveraging Nvidia's Megatron-LM Shoeybi et al., 2020 distributed trainer and a carefully curated medical corpus, including PubMed articles and international medical guidelines, MEDITRON outperforms state-of-the-art baselines across four major medical benchmarks. The study underscores the impact of increasing model parameters on medical LLM performance, highlighting MEDITRON's competitive edge against closed-source counterparts like GPT-3.5 and Med-PaLM. Notably, MEDITRON achieves performance levels within 5% of GPT-4 and 10% of Med-PaLM-2, thus potentially democratizing access to extensive medical knowledge.

The recent development of LLMs Boiko et al., 2023,Tamkin et al., 2021 has generated a great deal of enthusiasm due to their exceptional performance in natural language generation and understanding, as well as their adaptability in handling a variety of tasks. To improve the performance of Large Language Models (LLMs), particularly for disease identification and classification tasks. Oniani et al., 2024 proposed an ensemble prompting method called Models-Vote Prompting (MVP). The way MVP operates is that multiple LLMs are given the same task, and their results are combined via a majority voting procedure. The utility of MVP is demonstrated by experiments showing better results on one-shot unusual disease diagnostic tasks compared to distinct models in the ensemble. Additionally, the researchers provide a novel rare disease dataset, which is made available to researchers under the terms of the MIMIC-IV Data Use Agreement (DUA). For doing research and evaluating in the field, this set of data is a helpful resource.

The Retrieval Augmented Generation (RAG) Lewis et al., 2020 method is a natural language processing model that combines retrieval and generation components to handle knowledge-intensive tasks. In this paper Jin et al., 2024 used LLMs along with RAG to evaluate health reports with a novel feature extraction method. They used RAG to retrieve knowledge from the professional knowledge base. Researchers employ an automated feature engineering approach to train a classification model XGBoost for final disease prediction. The accuracy of GPT-4 combined with information retrieval by RAG for disease diagnosis is 0.68, and the F1 score is 0.71, while their framework achieved an accuracy of 0.833 and an F1 score of 0.762, respectively.

Dettmers et al., 2023 formerly employed QLoRA, an effective finetuning technique that maintained full 16-bit fine-tuning task performance while reducing memory usage to the point where a single 48GB GPU could finetune a 65B parameter model. The Guanaco model

family, described in the research as the top model family, achieves 99.3% of ChatGPT's performance level on the Vicuna test, beating out all other publicly available models in under 24 hours of fine-tuning on just one GPU. Results from this approach consistently demonstrate that, on educational standards with widely recognized evaluation settings, 4-bit QLORA with NF4 data type matches 16-bit complete finetuning and 16-bit LoRA finetuning performance. Additionally, they have demonstrated that NF4 (4-bit NormalFloat) outperforms FP4 (4-bit Float) and even indicated that performance is not diminished by double quantization.

A significant advancement in the field has been made recently with the development of HEAL Yuan et al., 2024, a Large Language Model (LLM) designed specifically for automated scribing and medical conversations. Based on the widely taught 13B LLaMA2 architecture, HEAL provides a novel approach to the unique issues associated with medical communication. An evaluation of HEAL on tasks like PubMedQA yields an excellent accuracy of 78.4%, proving its superiority over current LLMs like GPT-4 and PMC-LLaMA Wu et al., 2023. Furthermore, when it comes to producing medical notes, HEAL performs similarly to GPT-4, demonstrating its effectiveness in clinical documentation activities. Notably, HEAL outperforms human scribes and other similar models in terms of accuracy and completeness, and it outperforms GPT-4 and Med-PaLM 2 in terms of reliably identifying medical ideas.

## 4 System Description

The subsequent sections provide a list of the submodules used. We will describe why and how each model was utilized, and assess its relevance to our problem statement.

### 4.1 RAG using GEMINI

Large language models (LLMs) function best when Retrieval-Augmented Generation (RAG) Lewis et al., 2020 extends their capabilities to internal knowledge bases or specialized domains without requiring retraining. By guaranteeing that LLM output is accurate, pertinent, and usable in a variety of circumstances, this technique improves LLM output. Giving end users out-of-date or generic information when they're looking for specific answers is a prevalent problem with LLMs. This problem is solved by RAG, which instructs LLMs to obtain relevant information from reliable, pre-selected knowledge sources, improving accuracy and dependability.

Domain-specific or pertinent data is loaded, split into appropriately sized chunks to preserve context, and finally embedded using embedding models. The resultant embeddings are kept in a vector database so that documents with similar semantic content may be quickly retrieved. Data is then extracted from these embeddings according to how closely the query supplied by the user matches the documents. We use RAG with Gemini as



Figure 1: Proposed Model - Quantised Palmyra with RAG

the foundational LLM because of Gemini's extensive knowledge base as well as its large context window which allows chunks with higher semantic lengths to be supplied by the retriever.

LangChain simplifies the implementation of RAG by providing tools to load relevant datasets, such as the MedQA dataset, through its Data Loaders. It facilitates the chunking of data and the creation of embeddings using predefined functions and embedding models. The user's query is incorporated into a template and given as input into the LLM, while a Retriever component assists in finding similar documents based on query similarity. Utilizing MedQA data enhanced Gemini's answering ability, resulting in improved accuracy and relevance in responses. This integrated approach underscores the effectiveness of RAG in augmenting LLM performance specifically in the domain of Medical science.

### 4.2 Palmyra Quantised version

In our experiments, we employed a big decoder-only transformer model, known as Palmyra-20b. The Pile dataset Gao et al., 2020, which was tokenized with the GPT2 Radford et al., 2019 BPE tokenizer, served as the pre-training dataset for Palmyra-20b. It is a GPT-based model with 48 attention heads, a hidden size of 6144, 44 transformer layers, and a sequence length of 2048. The distributed Adam optimizer was used to train the model, which has two parallelism configura-

tions: pipeline parallelism of 1 and tensor parallelism of 4. Given the constraints of limited computational resources, we implemented 4-bit quantization on the model to mitigate computational demands while preserving efficiency. Quantization Gholami et al., 2021 is a technique that involves the process of converting the weights of the model from a higher precision to a lower precision. In our approach, we used 4-bit quantization to reduce the precision of weights and activations of Palmyra-20b to only 4-bit integer format. By quantization, we were able to significantly decrease memory and computational requirements without compromising model performance substantially to give accurate predictions by analyzing the symptoms provided to the model.

## 4.3 QLoRA on palmyra

Since fine-tuning LLMs like Palmyra20B is highly computationally expensive, we used PEFT ( Parameter Efficient Fine Tuning ) to make sure the training could be carried out on consumer-grade GPUs. In particular, we used QLORA Dettmers et al., 2023 ( Quantized Low-Rank Adaptation ) which quantizes a pre-trained model to 4-bit weights and adds an Adaptor - a low-rank tensor of trainable weights that can then be used to fine-tune the model through back-propagation. QLORA achieves far more efficient fine-tuning through the use of 4-bit Normal Float datatype which has been empirically proven to yield superior results to 4bit Floats. QLORA also employs double-quantization where not only are the weights but the quantization constants themselves are also quantized saving further memory. Finally, this approach uses Paged Optimisers allowing NVIDIA to manage memory effectively and ensuring that QLORA gives optimal results in parallel processing.

## 4.4 Proposed Model - Quantised Palmyra with RAG

In this, we incorporated 3 modules for the Error detection and correction task. The first one was the RAG module as explained in the previous section the second was the quantized and finetuned Palmyra med 20B and the third NER module.

We initially conducted Named Entity Recognition (NER) using GEMINI to identify words representing diseases or vaccines in the text. After masking these identified words, we implemented the Retrieval-Augmented Generation (RAG) model on an external dataset Jin et al., 2020. If the RAG score fell below a certain threshold, we passed the input to our model i.e. Palmyra(quantized and finetuned version). If the word provided by our Palmyra or RAG model matched the word detected by NER, no error was detected. Otherwise, if a different word is obtained from the model, then it is replaced with the masked word identified by NER. Finally, the error sentence is mapped with the sentence ID to get the output in the desired format. Our proposed model is illustrated in Figure 1, which provides a visual representation of the key components and relationships

within our framework. For determining the error flag NER plays a pivotal role. It is so because irrespective of what flow the text takes (i.e. RAG or Palmyra), the output will always be compared with the NER's output for the error flag.

## 4.5 Metrics Used

To evaluate our model and assess its accuracy in light of the corrected sentence, we have adopted the following metrics for evaluation

### 4.5.1 R1-F

The ROUGE-1 F1-score is a metric commonly utilized in natural language processing tasks, particularly in the evaluation of automatic text summarization systems. ROUGE, which stands for Recall-Oriented Understudy for Gisting Evaluation, focuses on measuring the quality of summaries generated by algorithms in comparison to human-generated reference summaries.

Specifically, the ROUGE-1 F1-score assesses the overlap of unigrams (individual words) between the generated summary and the reference summary. It is computed by taking into account both precision and recall of unigrams. Precision measures the proportion of correctly included unigrams in the generated summary relative to all unigrams present, while recall measures the proportion of correctly included unigrams relative to all unigrams in the reference summary.

$$F1 = 2 \times \frac{precision \times recall}{precision + recall} \quad (1)$$

Here, precision is the number of samples correctly predicted out of the number of samples predicted in that category. Recall is the number of samples predicted correctly out of the number of samples present for that class.

### 4.5.2 BERT SCORE

BERTScore is a collection of three metrics - BERT-Precision, BERT-Recall, and BERT-F1. As the names imply, BERT-Precision measures how well the candidate texts avoid introducing irrelevant content. BERT-Recall measures how well the candidate texts avoid omitting relevant content. BERT-F1 is a combination of both Precision and Recall to measure how well the candidate texts capture and retain relevant information from the reference texts.

$$P_{BERT} = \frac{1}{|\hat{x}|} \sum_{\hat{x}_j \in \hat{x}} \max_{x_i \in x} (x_i^T \cdot \hat{x}_j) \quad (2)$$

$$R_{BERT} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\hat{x}_j \in \hat{x}} (x_i^T \cdot \hat{x}_j) \quad (3)$$

$$F1 = 2 \times \frac{P_{bert} \times R_{bert}}{P_{bert} + R_{bert}} \quad (4)$$

| Model | Score | | | |
|---|---|---|---|---|
| | **R1F Score** | **BERT Score** | **BLEURT Score** | **Aggregate Score** |
| Quantised Palmyra | 0.46277 | 0.48681 | 0.49753 | 0.482371 |
| Quantised+QLoRa | 0.54802 | 0.57079 | 0.55477 | 0.55786 |
| Pure RAG | 0.66376 | 0.64557 | 0.60720 | 0.63884 |
| **Quantised+QLoRa+RAG** | **0.70306** | **0.74372** | **0.75217** | **0.73298** |

Table 2: Scores for various model

### 4.5.3 BLEURT

BLEURT (Bilingual Evaluation Understudy with Representations from Transformers) Sellam et al., 2020 is a novel, machine learning-based automatic metric for Natural Language Generation BLEURT that can capture non-trivial semantic similarities between sentences. It takes a pair of sentences as input, a reference, and a candidate, and it returns a score that indicates to what extent the candidate is fluent and conveys the meaning of the reference.

### 4.5.4 Aggregate Score

The aggregate score is calculated as the Mean of ROUGE-1-F, BERTScore, and BLEURT

$$Aggregate = \frac{R1F + BERTScore + BLEURT}{3}$$

(5)

| Parameter | Value |
|---|---|
| per_device_train_batch_size | 4 |
| gradient_accumulation_steps | 4 |
| optim | paged_adamw_32bit |
| logging_steps | 1 |
| learning_rate | 1e-4 |
| fp16 | True |
| max_grad_norm | 0.3 |
| num_train_epochs | 2 |
| evaluation_strategy | steps |
| eval_steps | 0.2 |
| warmup_ratio | 0.05 |
| save_strategy | epoch |
| group_by_length | True |
| save_safetensors | True |
| lr_scheduler_type | Cosine |
| Seed | 42 |

Table 3: Hyperparameters for Fine Tuning

## 5 Experimental Setup

We primarily used Google Colab notebooks for our workflow as well as for less computationally demanding tasks such as NER, EDA, text masking, RAG, etc.

Colab notebooks provide free access to a single T4 GPU (12GB RAM, 8GB VRAM, 64GB disk space). However, running quantized LLMs or fine-tuning had much higher computational requirements, and we therefore used Kaggle notebooks, which provide limited access to 2x T4 GPUs (15 GB of VRAM each). Please refer to Table 3 for a comprehensive overview of the parameters employed during the fine-tuning process. Since dataset preparation requires disk storage and frequent reads and writes, we use Jupyter Kernels for the same. We used the BitsAndBytes library for 4-bit quantization as well as the PEFT, Accelerate, and Datasets libraries by Huggingface for fine-tuning.

For performing NER on text, we used the GEMINI API from Google AI Studio. It had a maximum query limit of 60 queries per second. Since we were using GEMINI for NER as well as for RAG, this became our bottleneck, which sometimes led the session to crash. To address this, we imposed a timeout after every few API calls as well as made frequent local saves to the inferred results.

We implemented RAG using the Langchain framework, using GEMINI as our LLM. For implementing retrieval in our knowledge base, we used GEMINI embeddings to populate our vector store, which was a locally created ChromaDB instance.

## 6 Result

In our study, we employed a series of approaches aimed at enhancing the accuracy of our model. Initially, we implemented the quantized Palmyra approach, in which we tested the model that we built after the 4-bit quantization of Palmyra-20b. This gave a modest aggregate score of 0.482371. However, recognizing the room for improvement, we continued to refine our methodology. Building upon the quantized palmyra framework, we introduced the quantized+ QLoRa approach. In quantised palmyra, we fine-tuned using QLora on MEDQA data, which demonstrated a notable improvement, yielding an aggregate score of 0.55786. Encouraged by this progress, we further augmented our model with the Pure RAG technique, resulting in a substantial enhancement in aggregate score to 0.63884. Finally, through the integration of all three approaches—quantized, QLoRa, and RAG—into our model, we achieved the highest aggregate score of **0.73298**. The detailed scores for each approach are described in Table 2.

## 7 Limitations

The model struggles to give the correct output if the error is not related to a disease or pathogen. NER plays a crucial role in detecting pathogens or diseases from the text and therefore proves to be a bottleneck for accuracy since if NER fails to accurately determine the

disease, pathogen, or suggestion, the result will not be accurate regardless of the robustness of the model. The RAG approach fails for symptoms that are phrased very differently from those in the principal texts.

## 8 Future Work

**Using Larger and More Powerful LLMs**: Larger LLMs like Meditron-70b and Palmyra-med-40b can be used for achieving better accuracy in error detection and correction in clinical notes given sufficient computational power. The greater number of weights in these larger models allows them to capture more intricate patterns and nuances in the data during training.

**FineTuning on a larger dataset**, which will contain richer and more diverse medical information, can improve the model performance. Integrating multimodal information, such as images or structured data from electronic health records, alongside text data could provide richer context and improve error detection and correction accuracy.

**Enhancing Model Robustness:** The model can be made more robust against failures by having an end-to-end architecture where individual modules like NER, error detection, etc. are not carried out independently.

## 9 Conclusion

To conclude with this work for the MEDIQA-CORR task at NAACL, In ClinicalNLP 2024, we investigated four approaches for detecting and correcting errors in clinical notes. Our experiments demonstrated that the combined approach of Quantised Palmyra with RAG achieved the best performance, with an aggregate score of 0.73298. However, a key limitation identified is the reliance on named entity recognition (NER). Errors in NER can impact the overall performance of the system. Looking towards the future, research efforts should focus on mitigating the dependence on NER. Additionally, exploring alternative techniques and leveraging a larger, more comprehensive dataset holds promise for further improving the accuracy of error detection and correction in clinical notes. This will ultimately lead to a more robust and reliable system for enhancing the quality of clinical documentation.

## References

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Daniil A. Boiko, Robert MacKnight, and Gabe Gomes. 2023. Emergent autonomous scientific research capabilities of large language models.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023. Meditron-70b: Scaling medical pretraining for large language models.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms.

Ran Elgedawy, Sudarshan Srinivasan, and Ioana Danciu. 2024. Dynamic qa of clinical documents with large language models.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling.

Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W. Mahoney, and Kurt Keutzer. 2021. A survey of quantization methods for efficient neural network inference.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams.

Mingyu Jin, Qinkai Yu, Dong Shu, Chong Zhang, Lizhou Fan, Wenyue Hua, Suiyuan Zhu, Yanda Meng, Zhenting Wang, Mengnan Du, and Yongfeng Zhang. 2024. Health-llm: Personalized retrieval-augmented disease prediction system.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W. Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering.

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.

David Oniani, Jordan Hilsman, Hang Dong, Fengyi Gao, Shiven Verma, and Yanshan Wang. 2024. Large language models vote: Prompting for rare disease identification.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa : A large-scale multi-subject multi-choice dataset for medical domain question answering.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. Bleurt: Learning robust metrics for text generation.

Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2020. Megatron-lm: Training multi-billion parameter language models using model parallelism.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. 2023. Towards expert-level medical question answering with large language models.

Alex Tamkin, Miles Brundage, Jack Clark, and Deep Ganguli. 2021. Understanding the capabilities, limitations, and societal impact of large language models.

Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Towards building open-source language models for medicine.

Dong Yuan, Eti Rastogi, Gautam Naik, Sree Prasanna Rajagopal, Sagar Goyal, Fen Zhao, Bharath Chintagunta, and Jeff Ward. 2024. A continued pretrained llm approach for automatic medical note generation.

Yinghao Zhu, Changyu Ren, Shiyun Xie, Shukai Liu, Hangyuan Ji, Zixiang Wang, Tao Sun, Long He, Zhoujun Li, Xi Zhu, and Chengwei Pan. 2024. Realm: Rag-driven enhancement of multimodal electronic health records analysis via large language models.

# LAILab at Chemotimelines 2024: Finetuning sequence-to-sequence language models for temporal relation extraction towards cancer patient undergoing chemotherapy treatment

**Shohreh Haddadan[1], Tuan-Dung Le[1,2], Thanh Duong[1,2], Thanh Q. Thieu[1,2],**

[1]Moffitt Cancer Center and Research Institute, USA
[2]University of South Florida, USA
{shohreh.haddadan, tuandung.le , thanh.duong, thanh.thieu}@moffitt.org

## Abstract

In this paper, we report our effort to tackle the challenge of extracting chemotimelines from EHR notes across a dataset of three cancer types. We focus on the two subtasks: 1) detection and classification of temporal relations given the annotated chemotherapy events and time expressions and 2) directly extracting patient chemotherapy timelines from EHR notes. We address both subtasks using Large Language Models. Our best-performing methods in both subtasks use Flan-T5, an instruction-tuned language model. Our proposed system achieves the highest average score in both subtasks. Our results underscore the effectiveness of finetuning general-domain large language models in domain-specific and unseen tasks.

## 1 Introduction

Patient health records contain a wealth of information that can offer valuable insights to healthcare professionals and researchers, aiding in the enhancement of diagnosis, treatment, and disease prevention. Cancer patients often undergo lengthy treatment regimens, resulting in extensive electronic health record (EHR) documentation over time. The sheer volume of data available to healthcare providers is substantial, making manual curation impractical and cost-prohibitive.

A crucial aspect of cancer patient records is their chemotherapy treatment status documentation. Automatically extracting information regarding the timelines of chemotherapy events offers several advantages, including the ability to evaluate treatment efficacy across various cancer types. This automated extraction process also facilitates the creation of concise summaries for future attending physicians.

Two main tasks have been defined and addressed in association with temporal relation extraction from clinical notes: DocTimeRel and TLINK detection and classification. The first task is to iden-tify and classify the relation between events in an EHR note and the creation time of the document. TLINK detection and classification identify relations between event mentions and time expressions in EHR notes.

In this paper, we deal with the latter, the temporal relation extraction on a dataset of three cancer types. The shared task defines two subtasks. Subtask one aims at discerning a temporal relationship between a pair, consisting of a chemotherapy event and a time expression, subsequently classifying this relationship into one of the following categories: CONTAINS, BEGINS-ON, or ENDS-ON. In the second subtask, the only given input is the patient notes. The desired output for both subtasks is patient-level chemotherapy timelines. For detailed information on the definition of the subtasks, baseline methods, dataset, and evaluation criteria, see (Yao et al., 2024).

We approach both subtasks using large language models (LLMs). For the first subtask, we reformulate the relation classification problem into a text generation task and benefit from instruction-tuned language models to predict the relation. In the second subtask, we experiment with a sequence-to-sequence fine-tuning method with relations transformed into target sequences using a triplet linearization algorithm and also a pipeline method consisting of a rule-based event and time expression module and our best-performing model on the first subtask to identify and classify the pairwise relations.

We achieved the highest average scores on the test results as announced by the organizers (Yao et al., 2024).

In the following, we describe how we have utilized LLMs in detection, classification, and the end-to-end approach to chemotherapy timeline extraction from clinical notes.

Figure 1: Low-rank adaptation instruction fine-tuning for Subtask 1

## 2 Methodology

### 2.1 Subtask1

With the chemotherapy events and time expressions in each patient's note already provided by the organizer, the first subtask aims to identify temporal relations between them and subsequently generate patient-level timelines.

Prior to training a model, we need to prepare the dataset to train the model for the temporal relation classification task. The annotated relations with their respective pair of events and time expressions in the gold standard training/development dataset are added as positive instances. We create negative instances labeled as NO-RELATION by pairing events and time expressions within a patient note with no temporal relations in the gold-standard dataset.

However, incorporating every potential negative instance would lead to a significant imbalance in the training dataset as well as additional computational costs for training and inferencing the model. To mitigate this, we exclude instances where the positional distance between the event mentions and time expressions in the EHR note exceeds a maximum number of characters. Table 1 reports the maximum distance and number of NO-RELATION label instances added to the dataset. With this empirical observation, we set the maximum distance to 250 characters. We also create a heuristic rule that any pairs with a distance greater than the threshold are automatically predicted as NO-RELATION during inference on the test set. Applying this rule to the test set reduces the number of possible pairs from 12762 to 3042, thus enabling a more computationally efficient inference process.

During preprocessing, we first employ the

"mimic" model from the Stanza library, developed by (Zhang et al., 2021), for sentence segmentation. Then, we construct the context for the input sequences using two different approaches: concatenated context and bounded context. If the event and the time expression in the pair occur within the same sentence, both methods consider the sentence as the context. Otherwise, if the event and time expression are located in different sentences, the two sentences are concatenated to form the concatenated context. In the bounded context method, any sentence occurring between these two sentences is also included in the context. In addition, we add markers denoted by <e> followed by </e>, and <t> followed by </t> to respectively delineate events and time expressions.

We reformulate the temporal relation classification task as a generation task by finetuning a large language model to directly generate a label from the predefined set of relation types: CONTAINS, BEGINS-ON, ENDS-ON, NO-RELATION. We prepend the instruction describing the task to each input context. This method conditions the model to generate the relation type immediately following an anchor prompt "Relation:". Figure 1 illustrates our approach to tackling the first subtask, including our model's input and expected output. In the instruction, we use the definition of events, time expressions, and temporal relations provided in the data descriptions of the shared task. Our instruction format leverages the prompt structure used in relation extraction tasks, as described by (Lai et al., 2023). We also experiment with finetuning the model without adding the task instruction to the input contexts.

During our preliminary experiments, we fine-

| Split | Cancer type | Gold relation pairs# | Max character distance | No relation pairs # |
|-------|-------------|---------------------|------------------------|---------------------|
| Train | brca | 455 | 99 | 381 |
|       | mela | 48 | 218 | 35 |
|       | ovca | 494 | 143 | 336 |
| Dev   | brca | 113 | 213 | 132 |
|       | mela | 201 | 144 | 191 |
|       | ovca | 226 | 173 | 220 |

Table 1: Maximum (character) distance between event and time mentions of relation pairs in the gold standard dataset used as a threshold to reduce the number of NO-RELATION pairs. The number of gold relations is provided for comparison.

tune three instruction-tuned models: Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), Flan-T5-xxl 11B (Chung et al., 2022) and Llama-2-13B-chat (Touvron et al., 2023).

Flan-T5-xxl consistently achieved superior performance on the development set compared to the other two models. Thus, we use Flan-T5-xxl for further experimentation in this subtask.

## 2.2 Subtask2

As Yao et al. (2024) describe, in the second subtask, the input to the system is only the patient's EHR notes. Therefore, an end-to-end system that integrates identifying chemotherapy events and time expressions to extract the final chemotherapy patient-level timeline is required. We consider two different approaches to address this subtask.

In the first approach, we train a sequence-to-sequence model with input snippets from the EHR notes. The output is sequences containing the temporal relations (each a triplet of <event, relation type, time expression>) found in that snippet. The training objective is to simultaneously identify the events and time expressions in the context of each sentence in the EHR note and to detect and classify the relation as CONTAINS, BEGINS-ON, or ENDS-ON. We use the annotated data for the first subtask to train the models and evaluate our models on the development set provided for the second subtask.

We consider the context of each sentence to be its neighboring sentences (one preceding and one succeeding) joined by the separator token of the

corresponding tokenizer as defined in equation (1).

$$Context(s_i) = s_{i-1} + [SEP] + s_i + [SEP] + s_{i+1} \quad (1)$$

Huguet Cabot and Navigli (2021) have neatly introduced a triplet linearization algorithm for generating target sequences incorporating one or more relations between entities. We adopt this algorithm to transform the annotated temporal relations to target sequences.

Our approach differs from Huguet Cabot and Navigli (2021)'s approach in several ways. Firstly, their approach is to identify more than 200 relation types; thus, contrary to our setting, they are not limited to a restricted set. We add the relation types (CONTAINS, BEGINS-ON, ENDS-ON) to the special tokens of the tokenizer so they are not split during the tokenization process and the model learns them as defined in the target sequences. Since the events in our problem settings are domain-specific, we observed that the approach used in Huguet Cabot and Navigli (2021) identifies any event (not only the chemotherapy events) as a chemotherapy event after training. To prevent the generation of false positive events, we include additional chemotherapy events annotated in the gold standard data set, which are not in any relation with a time expression, to the training set. Similarly, to create negative instances, we add the input sequences that have no annotation of chemotherapy events, time expressions, or relations to the training data. Figure 2 shows different input sequence and their corresponding target sequence.[1]

We then experiment with finetuning various versions of two pre-trained models with the encoder-decoder structure, which have proven to perform well for sequence-to-sequence tasks, namely BART and Flan-T5. The reasoning behind choosing BART is that it is trained for sequence-to-sequence tasks and has proven to perform well on sequence-generation tasks. We chose Flan-T5 to test the effectiveness of this instruction-tuned model on an unseen task. In this subtask, we do not add instructions to input sequences while finetuning Flan-T5. We conduct experiments on various available model sizes for BART and Flan-T5.

In the second approach, we use a pipeline method that consists of two steps: the first step extracts the chemotherapy events and time expres-

---

[1]To abide by the terms of the data agreement, we refrain from quoting exact snippets from the EHR notes. The examples are altered and, therefore, might not be medically accurate.

| | |
|---|---|
| **Input Sequence₁** | They underwent surgery. On day of admission, they had their first dose of Taxol. Their blood glucose was 456. |
| **Target Sequence₁** | \<triplet> day of admission \<subj> Taxol \<obj> CONTAINS |
| **Input Sequence₂** | Vital signs are stable. Culture results significant cancer, currently getting chemotherapy. They present for evaluation today. |
| **Target Sequence₂** | \<triplet> currently \<subj> chemotherapy \<obj> CONTAINS |
| **Input Sequence₃** | Patient seen on 04/12. They received first dose of aflibercept today and second dose 05/16 prior to admission for high dose. Transferred to unit. |
| **Target Sequence₃** | \<triplet> aflibercept \<subj> today \<obj> BEGINS-ON \<triplet> 05/16 \<subj> aflibercept \<obj> CONTAINS |
| **Input Sequence₄** | No known medicinal allergies. They were initiated on the ipilimumab arm. They continue on the recommended regimen. |
| **Target Sequence₄** | \<triplet> \<subj> ipilimumab \<obj> |
| **Input Sequence₅** | Malignant melanoma of other specified site. Patient here for cycle #2 of IL-2, on study with aflibercept. follow electrolytes and renal function |
| **Target Sequence₅** | \<triplet> \<subj> IL-2 \<obj> \<triplet> \<subj> aflibercept \<obj> |
| **Input Sequence₆** | Patient was seen and examined. History and physical exam were reviewed. I agree with physical findings. |
| **Target Sequence₆** | \<triplet> \<subj> \<obj> |

Figure 2: The input sequences are the contexts, including a sentence and its preceding and succeeding sentence in the EHR note joined by the separator token of the corresponding tokenizer. Target sequences are the linearized triplets taken from the gold standard annotations. Following the encoding in Huguet Cabot and Navigli (2021), \<triplet> marks the start of a new temporal relation with a new head entity, followed by the tokens representing the head entity in the input text; \<subj> marks the end of the head entity and the start of the tail entity's tokens; \<obj> marks the end of the tail entity and the start of the relation type between the head and tail entity. The head/tail entities can be either a chemotherapy event or a time expression depending on their relative position in the text.

sions, and in the second step, we utilize our best-performing model on the first subtask to detect and classify the relations between pair of events and time expressions. We extract the time expressions using the Python wrapper for Stanford CoreNLP's SUTime Java library developed by Manning et al. (2014)[2]. We utilize a rule-based system with a predefined dictionary for the event extraction task. We compile a list of chemotherapy events from three different sources: 1) the baseline system[3]. 2) all chemotherapy events extracted from the training set, and 3) all the cancer drugs mentioned on the Cancer Research UK website[4].

## 2.3 Finetuning process

Our approach uses the Huggingface[5] implementation of the Seq2SeqTrainer to finetune trained models.

In the first subtask, we set the maximum length of the input as 450 tokens and the maximum target length as 10 tokens to fit the instruction. We finetune Flan-T5-xxl model using LoRA (Hu et al., 2021) for 10 epochs, employing early stopping with a patience of 3 epochs.

In the second subtask, we set the maximum

length of the input as 256 tokens and the maximum target length as 32 tokens. We then pad input and target sequences to maximum length with the pad token of the tokenizer specific to the model. We run the finetuning for 10 epochs in the BART setting and 5 epochs in Flan-T5 setting. The parameter efficient module (LoRA) was enabled while finetuning Flan-T5 models for this subtask. For more details on the models, see Appendix A.1.

For both finetuning experiments, we used the implementation of LoRA in Huggingface library. Parameters for LoRA are set to $\alpha = 32$, dropout $= 0.05$, and $r = 16$ and are added to $[q, k, v, o]$ layers in both tasks. Appendix A.2 briefly describes LoRA.

## 2.4 Preparing data for evaluation

Most of the time expressions in the EHR notes are relative and must be normalized using the document time (DOCTIME). The document time in the first subtask can be extracted from the gold standard annotated data or the headers of each patient EHR note. In the case of the second subtask, only the latter is feasible due to the absence of gold standard annotations. The headers of the patient records are provided in a standard format, so the document time can be precisely extracted using regular expressions.

To normalize relative time expressions such as "two weeks ago", "today", "currently", we use the timenorm library (Xu et al., 2019). We discard the extracted relations for which timenorm fails to

---

[2]https://pypi.org/project/sutime
[3]https://github.com/HealthNLPorg/chemoTimelinesBaselineSystem/tree/main/timelines/instance-generator/src/user/resources/org/apache/ctakes/dictionary/lookup/fast/bsv
[4]https://www.cancerresearchuk.org/about-cancer/treatment/drugs
[5]https://huggingface.co/

| Cont. | Inst. | brca | | | mela | | | ovca | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | F1 | RF1 | TF1 | F1 | RF1 | TF1 | F1 | RF1 | TF1 |
| Bound | No | 0.893 | **0.992** | 0.941 | 0.922 | 0.938 | **0.887** | 0.879 | 0.968 | 0.852 |
| Bound | Yes | **0.922** | 0.980 | **0.962** | **0.960** | **0.977** | **0.887** | **0.916** | **0.987** | 0.793 |
| Concat | No | 0.913 | 0.980 | 0.937 | 0.898 | 0.916 | **0.887** | 0.890 | 0.968 | **0.871** |
| Concat | Yes | 0.919 | 0.967 | 0.918 | 0.934 | 0.954 | **0.887** | 0.893 | 0.960 | 0.810 |

Table 2: Results for the first subtask on the development set. The terms F1 and RF1 represent the F1-score and relaxed F1-score of our classification model, respectively. TF1 is the official F1-score for the final timelines calculated using the evaluation system.

normalize the time expression. Examples of such time expressions include "at this time", "January 10 or 11", "05/2012" , "day one" , "16-09", etc.

For both subtasks, we use the baseline system provided by Yao et al. (2024) for de-duplication and creation of final timelines.[6]

As an extra step in the pipeline approach to the second subtask, we manually omitted some events and time expressions from the results of the rule-based systems. Examples of such omissions are "continues" event (which appears in the train set) and time expressions "1842", "1255" and "1000".

## 3 Results

We use the evaluation system provided by (Yao et al., 2024) for both subtasks on the development set.[7] The evaluation script receives the gold standard timelines, and the system prediction for all patients in each cancer type as input.

All the experiments on the development set have been executed before the test set results were announced.

### 3.1 Subtask1

In addition to reporting the timeline score on the development set using the organizer's evaluation system, we also evaluate our model's performance on the pairwise temporal classification task (Table 2). We implement two metrics: micro F1 and relaxed micro F1. CONTAINS and BEGINS-ON, CONTAINS and ENDS-ON are interchangeable in the relaxed F1-score computation.

Finetuned Flan-T5-xxl with instruction and bounded context achieved the highest scores on almost all metrics. Finetuning bounded context shows a marginal improvement in relaxed micro

F1 compared to the concatenated context. This suggests that incorporating sentences between sentences containing event and time expression might be beneficial for classifying NON-RELATION pairs.

Our classification model scores do not correlate well with timeline scores. For instance, in ovarian cancer results, fine-tuned Flan-T5-xxl bounded context and instruction achieves the highest F1-score on the classification task but the lowest timeline score. We suspect that this difference originates from official results being based on average macro F1 score across all patients. Further reasons might be related to the errors of the post-processing steps in creating the final patient timelines, such as the normalization of time expressions and the de-duplication process. We select three submissions with the highest average F1-score of F1-scores, relaxed F-scores and final timeline F-scores for all cancer types as presented in Table 2.

Our submission outperformed the baseline for breast cancer, ovarian cancer, and the average score. It achieved the same score as the baseline system for melanoma cancer (Table 3).

### 3.2 Subtask2

The end2end approach with Flan-T5-xxl + LoRA achieves the highest results across all other methods and the baseline system results for melanoma and ovarian cancer as shown in Table 4. For breast cancer, on the other hand even though this method performs best among other implemented methods, it does not surpass the baseline system results on the development set.

Considering the relaxed setting, Flan-T5-xxl + LoRA has achieved the highest precision rate across all cancer types. However, the methods that extract event types using a rule-based or dictionary-based system (baseline system and the pipeline approach) have gained higher recall scores in the

| | Method | brca | mela | ovca | Average score |
|---|---|---|---|---|---|
| **Subtask 1** | Baseline system | 0.93 | **0.87** | 0.88 | 0.89 |
| | Flan-T5-xxl + bound context + instruction | **0.96** | **0.87** | 0.88 | **0.90** |
| | Flan-T5-xxl + bound context | 0.95 | 0.85 | **0.89** | **0.90** |
| | Flan-T5-xxl + concat context | 0.95 | 0.84 | **0.89** | **0.90** |
| | Highest score on the leader board | 0.96 | 0.87 | 0.89 | 0.90 |
| **Subtask 2** | Baseline system | 0.59 | 0.43 | 0.71 | 0.58 |
| | End2end BART-large | 0.52 | 0.57 | 0.59 | 0.56 |
| | End2end Flan-T5-xxl + LoRA | 0.62 | **0.74** | **0.74** | **0.70** |
| | Pipeline system | 0.53 | 0.38 | 0.49 | 0.47 |
| | Highest score on the leader board | 0.68 | 0.74 | 0.74 | 0.70 |

Table 3: Evaluation published by the organizers for our submission on the held-out test set

same setting. The low score in the strict evaluation setting for Flan-T5 is due to its failure to identify ENDS-ON relations in any cancer type, possibly because of the label's low frequency in the training set. See Appendix B for detailed results on precision and recall.

We chose to submit the results of the end2end method with both BART and Flan-T5 with the highest scores, which are the largest models we experimented with, namely BART-large and Flan-T5-xxl and the pipeline approach as our third submission.

We first performed a sentence tokenization step on the test data and extracted the contexts as input sequences. We used the models to infer target sequences. The results of these three approaches on the test data provided by the organizers are presented in Table 3. The end2end approach using the pre-trained Flan-T5-xxl model with LoRA exceeds in all evaluations except for the breast results. Albeit, the breast cohort results surpass the test set's baseline score contrary to the experiments on the development set. The average score on this approach gains the highest score among other submissions, as reported by the organizers.

| Method | brca | mela | ovca |
|---|---|---|---|
| Baseline system | **0.857** | 0.456 | 0.329 |
| Pipeline Approach | 0.529 | 0.511 | 0.470 |
| End2End BART-large | 0.700 | 0.618 | 0.496 |
| End2End Flan-T5-xxl | 0.749 | **0.720** | **0.647** |

Table 4: Evaluation for the second subtask on the development set.

## 4 Error Analysis

Since the gold standard timeline and annotations on the test set have not been released to enable future editions of the task, we will provide the error analysis on the results of the development set.

### 4.1 Subtask1

Our best model, Flan-T5-xxl finetune bounded context with instruction, achieved a low error rate of 20 incorrect predictions out of 1,083 tested pairs. Possible error sources include misspellings, potentially missing or spurious annotations, or unclear or complex contexts. Complex contexts occur when notes include tables that have lost their structures in the plain text files. We list some examples of mispredictions below.

- Misspelling: "Condition **<t>yesterdat</t>** appeared improved with treatment, and **<e>chemo</e>** cycle discontinued.", Label: ENDS-ON, Predict: NO-RELATION.

- Missing annotation: "Patient with metastatic melanoma enrolled in protocol and received first dose of **<e>aflibercept</e>** 9/4 and second dose **<t>09/18</t>** prior to admission for high dose IL2 (first cycle)Thus far has received 9/12 planned doses.", Label: NO-RELATION, Predict: CONTAINS.

- Spurious annotation: "Patient enrolled in protocol and received first dose of alibercept 9/4 and second dose **<t>09/18</t>** prior to admission for high dose **<e>IL2</e>** (first cycle)", Label: CONTAINS, Predict: NO-RELATION.

- Unclear context: "Cycle #2 was initiated on **<t>September 10 , 2011</t>**; however, the patient had a severe reaction during the **<e>paclitaxel</e>** infusion.", Label: CONTAINS, Predict: NO-RELATION.

### 4.2 Subtask2

One source of error in subtask2 is the emergence of medical events or drugs as output events that

are not particularly chemotherapy events such as "radiation", "iv decadron", "bolus", "anti-vegf antibody", "augmentin" and so on. The following example shows one of the incident where our best performing model incorrectly identifies "radiation" as a chemotherapy event in temporal relation CONTAINS with "June 1st". [8]

- ... she is undergoing consultation with Dr. X for possible **radiation** therapy on June 1st.

We noticed that balancing the negative instances with the positive examples of temporal relations worsens this problem. Thus, we keep all negative instances in the training set to improve the identification of chemotherapy events. These negative instances include the ones containing events/time expressions but no relations, for example, target sequences 4 and 5 in Figure 2. And also the ones with no events and no time expressions, for example, input sequence 6 in Figure 2. We suspect this problem is caused by the bias of the pre-trained model in identifying all entities beyond the chemotherapy events. This approach improved the results; however, it's not completely resolved. It can further be alleviated either manually or by applying a post-processing filter created by experts to only keep the temporal relations with chemotherapy drugs/treatments.

In examining the distribution of relation types across various cancer types within the development set for the second subtask, we observed an imbalance in the dataset. Specifically, the ENDS-ON relation type was found to occur with frequencies of approximately 30%, 2%, and 14% concerning all chemotherapies within the final gold timelines for breast cancer, melanoma, and ovarian cancer, respectively. Given our approach's reduced accuracy in identifying the ENDS-ON relation type, this discrepancy explains the lower accuracy observed compared to the baseline system specifically concerning breast cancer within both the development and potentially the test set (Assuming the distribution of relation types on the test set is close to the distribution on the development set).

Another source of the model's confusion is the chemotherapy events that were not annotated in the training dataset. The first example was identified as a "chemotherapy" event in CONTAINS temporal relation with time expression "2003" and the second as "docetaxel", BEGINS-ON, "oct 3rd" by our

end2end model, however, we do not find the equivalent of this chemotherapy event instance in the annotated development set. In order to resolve this particular error, we would need further information about the annotation rules.

- History of Present Illness: Patient was diagnosed with disease in 2003 and treated with surgery, chemotherapy, and radiation per the patient.

- Patient says they are now taking docetaxel with 1st dose Oct 3rd and second due in early november.

We can also associate a fraction of errors to the normalization errors originating from the *timnorm* library, for example, in cases where time expressions containing two-digit years are inaccurately resolved to the 1900s.

## 5 Related work

Numerous studies focus on annotation (O'Gorman et al., 2016; Wang et al., 2022; Alsayyahi and Batista-Navarro, 2023), detection and classification (Lim et al., 2023; Huang et al., 2023) of temporal relations in the general domain.

In the medical domain, temporal relation extraction also received attention for its benefits in longitudinal studies of medication, treatments, and diseases, as well as in summarizing clinical notes for physicians' further reference. THYME annotation guidelines and corpus (Styler IV et al., 2014) and its extension (Wright-Bettner et al., 2020) is a considerable effort in the specification of process of temporal relation annotation process in clinical narratives based on ISO-TimeML (Pustejovsky et al., 2010).

Prior to the introduction of transformer-based language models a few studies approached various tasks of temporal relation extraction problem with feature-based supervised machine learning algorithms and sequential neural networks (Xu et al., 2013; Lee et al., 2016; Alfattni et al., 2020, 2021). Moreover, Lin et al. (2018) utilized unlabeled data by self-training neural networks in clinical temporal relation extraction.

After the rise of transformer-based models, temporal relation extraction from clinical notes also benefited from this significant development in NLP methods using models such as BERT (Lin et al., 2020; Zhou et al., 2021), BioBERT and BART (Wright-Bettner et al., 2020) for clinical text representation. Lin et al. (2021) continue training BERT using a masking method called entity-centric masking strategy, where they use the MIMIC III dataset

---

[8]Examples in this section have been altered to abide by the data use agreement.

as their training data. Their results on temporal relation extraction shows improvements on baselines using the model pretrained using this approach.

Most end-to-end systems for temporal relation extraction in the clinical domain have been tackled using a pipeline approach consisting of modules for event and time expression extraction and pairwise temporal relation detection and classification. Dligach et al. (2022) on the other hand, explore the use of sequence-to-sequence models in extracting temporal relations from text. They experiment with various input/output representations and adopt those representations, which enable the reconstruction of the snippets with several relations and repetitive event names in a text snippet. They report this approach's results utilizing different sequence-to-sequence LLMs such as BART and T5. Miller et al. (2023) approach temporal relation extraction problem as an end-to-end task without given events and time expressions using a combination of domain-specific pre-trained language model PubmedBERT (Gu et al., 2021) and a multi-headed attention classifier on THYME2 dataset (Wright-Bettner et al., 2020).

Bethard et al. (2016, 2017) organized previous shared tasks to incentivize the research on temporal relation extraction from clinical notes.

## 6 Conclusions

This paper presents our effort in participating in the Chemotimeline shared task. We apply an instruction finetuning method for temporal relation detection and classification and a sequence-to-sequence approach for extracting timelines directly from EHR notes to solve the first and second subtasks. Our approach, leveraging the power of general-domain Large Language Models and further fine-tuning them with parameter-efficient methods, secured the highest average scores across the different cancer types for both subtasks. The results of our approach using Flan-T5-xxl + LoRA underscore the potential of instruction finetuning in enhancing the capabilities of LLMs for unseen natural language understanding and generation tasks, even on domain-specific data. In future work, we aim at augmenting the data for low-frequency relation types and also harnessing the power of provided unlabeled data to continue pre-training Large Language models and to investigate the effect on the results of extracting temporal relations from cancer patient EHR notes.

## Limitations

There are several limitations to our experiments.

Firstly, our experiments were bounded by computational resource limitations. Specifically, our experiments employed the Flan-T5 model with parameter-efficient techniques due to constraints in available computational power on shared hardware and time. This limitation prevents us from comparing our methodology with implementations of Flan-T5 without LoRA approach. Secondly, we do not test our experiments on other datasets since annotated data in the medical domain on such a specific task is extremely scarce. Thus, we cannot claim that our results will be as high on different datasets. Moreover, since our method is fine-tuned on the provided data, practical use and release of the model are legally bound by the data agreement usage. Finally, our method uses a deep learning approach and, therefore, is limited by the explainability and interpretability constraints of such techniques.

## References

Ghada Alfattni, Niels Peek, and Goran Nenadic. 2020. Extraction of temporal relations from clinical free text: A systematic review of current approaches. *Journal of Biomedical Informatics*, 108:103488.

Ghada Alfattni, Niels Peek, and Goran Nenadic. 2021. Attention-based bidirectional long short-term memory networks for extracting temporal relationships from clinical discharge summaries. *Journal of Biomedical Informatics*, 123:103915.

Sarah Alsayyahi and Riza Batista-Navarro. 2023. TIMELINE: Exhaustive annotation of temporal relations supporting the automatic ordering of events in news articles. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16336–16348, Singapore. Association for Computational Linguistics.

Steven Bethard, Guergana Savova, Wei-Te Chen, Leon Derczynski, James Pustejovsky, and Marc Verhagen. 2016. SemEval-2016 task 12: Clinical TempEval. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1052–1062, San Diego, California. Association for Computational Linguistics.

Steven Bethard, Guergana Savova, Martha Palmer, and James Pustejovsky. 2017. SemEval-2017 task 12: Clinical TempEval. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 565–572, Vancouver, Canada. Association for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *arXiv preprint*.

Dmitriy Dligach, Steven Bethard, Timothy Miller, and Guergana Savova. 2022. Exploring text representations for generative temporal relation extraction. In *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 109–113, Seattle, WA. Association for Computational Linguistics.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685.

Quzhe Huang, Yutong Hu, Shengqi Zhu, Yansong Feng, Chang Liu, and Dongyan Zhao. 2023. More than classification: A unified framework for event temporal relation extraction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9631–9646, Toronto, Canada. Association for Computational Linguistics.

Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. REBEL: Relation extraction by end-to-end language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370–2381, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Viet Lai, Nghia Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Nguyen. 2023. ChatGPT beyond English: Towards a comprehensive evaluation of large language models in multilingual learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13171–13189, Singapore. Association for Computational Linguistics.

Hee-Jin Lee, Hua Xu, Jingqi Wang, Yaoyun Zhang, Sungrim Moon, Jun Xu, and Yonghui Wu. 2016. UTHealth at SemEval-2016 task 12: an end-to-end system for temporal information extraction from clinical notes. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1292–1297, San Diego, California. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Stanley Lim, Da Yin, and Nanyun Peng. 2023. LEAF: Linguistically enhanced event temporal relation framework. In *Proceedings of the 2nd Workshop on Pattern-based Approaches to NLP in the Age of Deep Learning*, pages 6–19, Singapore. Association for Computational Linguistics.

Chen Lin, Timothy Miller, Dmitriy Dligach, Hadi Amiri, Steven Bethard, and Guergana Savova. 2018. Self-training improves recurrent neural networks performance for temporal relation extraction. In *Proceedings of the Ninth International Workshop on Health Text Mining and Information Analysis*, pages 165–176, Brussels, Belgium. Association for Computational Linguistics.

Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2021. EntityBERT: Entity-centric masking strategy for model pretraining for the clinical domain. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 191–201, Online. Association for Computational Linguistics.

Chen Lin, Timothy Miller, Dmitriy Dligach, Farig Sadeque, Steven Bethard, and Guergana Savova. 2020. A BERT-based one-pass multi-task model for clinical temporal relation extraction. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*, pages 70–75, Online. 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.

Timothy Miller, Steven Bethard, Dmitriy Dligach, and Guergana Savova. 2023. End-to-end clinical temporal information extraction with multi-head attention. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2023, page 313. NIH Public Access.

Tim O'Gorman, Kristin Wright-Bettner, and Martha Palmer. 2016. Richer event description: Integrating event coreference with temporal, causal and bridging

annotation. In *Proceedings of the 2nd Workshop on Computing News Storylines (CNS 2016)*, pages 47–56, Austin, Texas. Association for Computational Linguistics.

James Pustejovsky, Kiyong Lee, Harry Bunt, and Laurent Romary. 2010. ISO-TimeML: An international standard for semantic annotation. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).

William F. Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, and James Pustejovsky. 2014. Temporal annotation in the clinical domain. *Transactions of the Association for Computational Linguistics*, 2:143–154.

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*.

Xiaozhi Wang, Yulin Chen, Ning Ding, Hao Peng, Zimu Wang, Yankai Lin, Xu Han, Lei Hou, Juanzi Li, Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. MAVEN-ERE: A unified large-scale dataset for event coreference, temporal, causal, and subevent relation extraction. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 926–941, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Kristin Wright-Bettner, Chen Lin, Timothy Miller, Steven Bethard, Dmitriy Dligach, Martha Palmer, James H. Martin, and Guergana Savova. 2020. Defining and learning refined temporal relations in the clinical narrative. In *Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis*, pages 104–114, Online. Association for Computational Linguistics.

Dongfang Xu, Egoitz Laparra, and Steven Bethard. 2019. Pre-trained contextualized character embeddings lead to major improvements in time normalization: a detailed analysis. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 68–74, Minneapolis, Minnesota. Association for Computational Linguistics.

Yan Xu, Yining Wang, Tianren Liu, Junichi Tsujii, and Eric I-Chao Chang. 2013. An end-to-end system to identify temporal relation in discharge summaries: 2012 i2b2 challenge. *Journal of the American Medical Informatics Association*, 20(5):849–858.

*Jiarui Yao, *Harry Hochheiser, WonJin Yoon, Eli Goldner, and Guergana Savova, editors. 2024. *Overview of the 2024 Shared Task on Chemotherapy Treatment Timeline Extraction, Proceedings of the 6th Clinical Natural Language Processing Workshop, , NAACL June 2024*. Mexico City, Mexico.

Yuhao Zhang, Yuhui Zhang, Peng Qi, Christopher D Manning, and Curtis P Langlotz. 2021. Biomedical and clinical english model packages for the stanza python nlp library. *Journal of the American Medical Informatics Association*, 28(9):1892–1899.

Yichao Zhou, Yu Yan, Rujun Han, J Harry Caufield, Kai-Wei Chang, Yizhou Sun, Peipei Ping, and Wei Wang. 2021. Clinical temporal relation extraction with probabilistic soft logic regularization and global inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14647–14655.

## A Preliminaries

### A.1 Pretrained Models

In our experiments, we have employed two main pre-trained model weights available for public use: BART (Lewis et al., 2020), Flan-T5 Chung et al. (2022) which we briefly introduce in this section. The details of how we finetune them for our specific task are described in sections 2.1 and 2.2.

- **BART** (Lewis et al., 2020) BART is a sequence-to-sequence model based on an encoder-decoder architecture, which is composed of a BERT-based bidirectional encoder and an auto-regressive GPT-based left to right decoder. BART is trained by the task of reconstructing a corrupted input sentence into its original text and it has proven to perform well for text-to-text generation tasks such as Summarization.

- **Flan-T5** (Chung et al., 2022) Instruction-tuning is a technique to explicitly guide Large Language models to perform specific tasks. Flan-T5 is a sequence-to-sequence Large Language model that has been fine-tuned using this technique on a mixture of tasks. Flan-T5 has shown performance improvement on unseen tasks.

### A.2 LoRA

Large language models inherent to their title have billions of parameters. Finetuning large language models for a specific task or domain is expensive and infeasible in terms of time and computational resource limitations. Hu et al. (2021) introduced Low-Rank Adaptation of Large Language (LoRA) models method to make the finetuning process of these models more efficient and conclusively more accessible by freezing the pre-trained weights of the model and injection of trainable rank decomposition matrices into different layers of the transformer architecture. This method drastically reduces the number of training parameters. It has been shown to perform comparably well to full-parameter finetuning methods and, in some cases, outperforms several baselines with comparable or fewer trainable parameters.

## B Detailed Results for the second subtask

We report the detailed results of strict and relaxed settings for all our experiments in the second subtask using the evaluation system in this section.

Table 5 contains the results of our experiments for the second subtask. We have experimented with the end2end approach described in section 2.2 using BART and Flan-T5 models with various sizes. Not surprisingly bigger models have performed better across all cancer types for both strict and relaxed evaluation settings. The pipeline approach achieves high recall scores for melanoma and ovarian cancer since it extracts events in a rule-based manner. However, the precision score is low in the pipeline approach, since it identifies drugs and treatments other than chemotherapy-specific ones.

|  |  | Micro | | | Macro Type A | | | Macro Type B | | | Official Score |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | Pr | Re | F1 | Pr | Re | F1 | Pr | Re | F1 | F1 |
| | **Strict** | | | | | | | | | | |
| | Baseline system | 0.622 | 0.718 | 0.667 | 0.811 | 0.871 | 0.835 | 0.663 | 0.823 | 0.727 | |
| | Bart-base | 0.333 | 0.256 | 0.290 | 0.646 | 0.640 | 0.642 | 0.222 | 0.207 | 0.213 | |
| | Bart-large | 0.625 | 0.385 | 0.476 | 0.826 | 0.785 | 0.797 | 0.537 | 0.727 | 0.455 | |
| | Flan-T5-large + LoRA | 0.409 | 0.231 | 0.295 | 0.667 | 0.645 | 0.653 | 0.278 | 0.220 | 0.242 | |
| | Flan-T5-xl + LoRA | 0.5 | 0.282 | 0.360 | 0.799 | 0.746 | 0.763 | 0.464 | 0.322 | 0.369 | |
| | Flan-T5-xxl + LoRA | 0.667 | 0.308 | 0.421 | 0.851 | 0.753 | 0.781 | 0.602 | 0.342 | 0.417 | |
| | Pipeline approach | 0.337 | 0.718 | 0.459 | 0.341 | 0.451 | 0.379 | 0.409 | 0.702 | 0.510 | |
| brca | **Relaxed** | | | | | | | | | | |
| | Baseline system | 0.795 | 0.816 | 0.805 | 0.866 | **0.894** | 0.876 | 0.809 | **0.855** | 0.837 | **0.857** |
| | Bart-base | 0.429 | 0.353 | 0.387 | 0.688 | 0.668 | 0.676 | 0.334 | 0.281 | 0.302 | 0.489 |
| | Bart-large | 0.818 | 0.5 | 0.621 | 0.888 | 0.813 | 0.837 | 0.701 | 0.501 | 0.564 | 0.700 |
| | Flan-T5-large + LoRA | 0.692 | 0.529 | 0.600 | 0.859 | 0.769 | 0.801 | 0.791 | 0.552 | 0.635 | 0.718 |
| | Flan-T5-xl + LoRA | 0.696 | 0.457 | 0.552 | 0.905 | 0.827 | 0.853 | 0.748 | 0.540 | 0.607 | 0.730 |
| | Flan-T5-xxl + LoRA | **0.833** | 0.441 | 0.577 | **0.944** | 0.830 | 0.863 | **0.851** | 0.547 | 0.634 | **0.749** |
| | Pipeline approach | 0.405 | **0.833** | 0.545 | 0.381 | 0.507 | 0.425 | 0.515 | 0.852 | 0.633 | 0.529 |
| | **Strict** | | | | | | | | | | |
| | Baseline system | 0.667 | 0.667 | 0.667 | 0.571 | 0.571 | 0.571 | 0.357 | 0.357 | 0.357 | |
| | Bart-base | 0.483 | 0.333 | 0.395 | 0.217 | 0.119 | 0.154 | 0.326 | 0.179 | 0.231 | |
| | Bart-large | 0.585 | 0.533 | 0.558 | 0.660 | 0.833 | 0.658 | 0.490 | 0.75 | 0.487 | |
| | Flan-T5-large + LoRA | 0.72 | 0.4 | 0.514 | 0.725 | 0.682 | 0.656 | 0.588 | 0.524 | 0.484 | |
| | Flan-T5-xl + LoRA | 0.629 | 0.489 | 0.550 | 0.702 | 0.817 | 0.714 | 0.553 | 0.726 | 0.571 | |
| | Flan-T5-xxl + LoRA | 0.686 | 0.533 | 0.6 | 0.726 | 0.833 | 0.733 | 0.590 | 0.75 | 0.6 | |
| | Pipeline approach | 0.347 | 0.911 | 0.503 | 0.499 | 0.865 | 0.574 | 0.249 | 0.798 | 0.362 | |
| mela | **Relaxed** | | | | | | | | | | |
| | Baseline system | 0.630 | 0.630 | 0.630 | 0.570 | 0.56 | 0.565 | 0.354 | 0.34 | 0.347 | 0.456 |
| | Bart-base | 0.44 | 0.458 | 0.449 | 0.204 | 0.167 | 0.183 | 0.305 | 0.25 | 0.275 | 0.229 |
| | Bart-large | 0.586 | 0.739 | 0.654 | 0.663 | 0.905 | 0.694 | 0.495 | 0.857 | 0.542 | 0.618 |
| | Flan-T5-large + LoRA | 0.72 | 0.565 | 0.634 | 0.708 | 0.690 | 0.664 | 0.561 | 0.536 | 0.496 | 0.580 |
| | Flan-T5-xl + Lora | 0.667 | 0.75 | 0.706 | 0.698 | 0.910 | 0.748 | 0.548 | 0.864 | 0.622 | 0.685 |
| | Flan-T5-xxl + Lora | **0.731** | 0.827 | 0.775 | **0.728** | 0.936 | 0.776 | **0.592** | 0.905 | 0.665 | **0.720** |
| | Pipeline approach | 0.3375 | **1.0** | 0.505 | 0.517 | **1.0** | 0.608 | 0.275 | **1.0** | 0.413 | 0.511 |
| | **Strict** | | | | | | | | | | |
| | Baseline system | 0.4 | 0.306 | 0.347 | 0.224 | 0.358 | 0.239 | 0.224 | 0.358 | 0.239 | |
| | Bart-base | 0.350 | 0.4 | 0.374 | 0.391 | 0.486 | 0.378 | 0.391 | 0.486 | 0.378 | |
| | Bart-large | 0.340 | 0.423 | 0.377 | 0.351 | 0.357 | 0.341 | 0.351 | 0.357 | 0.341 | |
| | Flan-T5-large + LoRA | 0.494 | 0.494 | 0.494 | 0.471 | 0.426 | 0.437 | 0.471 | 0.426 | 0.437 | |
| | Flan-T5-xl + LoRA | 0.557 | 0.518 | 0.537 | 0.488 | 0.559 | 0.483 | 0.488 | 0.558 | 0.483 | |
| | Flan-T5-xxl + LoRA | 0.564 | 0.411 | 0.476 | 0.581 | 0.545 | 0.504 | 0.581 | 0.544 | 0.504 | |
| | Pipeline approach | 0.265 | 0.659 | 0.378 | 0.297 | 0.692 | 0.389 | 0.297 | 0.692 | 0.389 | |
| ovca | **Relaxed** | | | | | | | | | | |
| | Baseline system | 0.558 | 0.426 | 0.483 | 0.280 | 0.465 | 0.329 | 0.280 | 0.465 | 0.329 | 0.329 |
| | Bart-base | 0.434 | 0.554 | 0.486 | 0.440 | 0.574 | 0.457 | 0.440 | 0.574 | 0.457 | 0.457 |
| | Bart-large | 0.506 | 0.620 | 0.557 | 0.498 | 0.590 | 0.496 | 0.498 | 0.590 | 0.496 | 0.496 |
| | Flan-T5-large + LoRA | 0.633 | 0.769 | 0.694 | 0.581 | 0.646 | 0.592 | 0.581 | 0.646 | 0.592 | 0.592 |
| | Flan-T5-xl + LoRA | 0.677 | 0.646 | 0.661 | 0.658 | 0.677 | 0.642 | 0.658 | 0.677 | 0.642 | 0.642 |
| | Flan-T5-xxl + LoRA | **0.686** | 0.515 | 0.588 | **0.726** | 0.592 | 0.647 | **0.756** | 0.592 | 0.647 | **0.647** |
| | Pipeline approach | 0.318 | **0.742** | 0.445 | 0.365 | **0.812** | 0.470 | 0.365 | **0.812** | 0.470 | 0.470 |

Table 5: System results for the second subtask on the development set

# Lexicans at Chemotimelines 2024: Chemotimeline Chronicles - Leveraging Large Language Models (LLMs) for Temporal Relations Extraction in Oncological Electronic Health Records

**Vishakha Sharma[1], Andres Fernandez[2], Andrei Ioanovici[2], David Talby[2], Frederik Buijs[3]**
[1]Roche Diagnostics, California, USA
[2]John Snow Labs, Delaware, USA
[3]F. Hoffmann-La Roche Ltd, Basel, Switzerland

## Abstract

Automatic generation of chemotherapy treatment timelines from electronic health records (EHRs) notes not only streamlines clinical workflows but also promotes better coordination and improvements in cancer treatment and quality of care. This paper describes the submission to the Chemotimelines 2024 shared task that aims to automatically build a chemotherapy treatment timeline for each patient using their complete set of EHR notes, spanning various sources such as primary care provider, oncology, discharge summaries, emergency department, pathology, radiology, and more. We report results from two large language models (LLMs), namely Llama 2 and Mistral 7B, applied to the shared task data using zero-shot prompting.

## 1 Introduction

Electronic Health Records (EHRs) are a rich repository of patient information, encompassing a wide array of formats and sources including physician notes, laboratory results, radiology images, and pathology reports. Due to the heterogeneous and unstructured nature of clinical data, it is cumbersome to visualize patient journeys or extract meaningful information from Electronic Health Records (EHRs) to help guide clinical decision making (Anand and Sadhna, 2023; Najafabadipour et al., 2020). EHR data is often dispersed, recorded in free text with substantial variability in terminology, and embedded in narrative formats that are not easy to process or normalize across healthcare settings and systems. In addition, privacy concerns further limit the use of clinical data across hospitals and geographical borders further compounding complexity (Reisman, 2017; Kehl et al., 2020; Levine et al., 2019; Banerjee et al., 2019) and difficulty to leverage EHR data for insights generation.

Large Language Models (LLMs), with their advanced natural language (Guevara et al., 2024; Chen et al., 2023a,c; Hochheiser et al., 2023; Bitterman et al., 2023) understanding capabilities, offer a transformative solution to these challenges. They can be trained to interpret complex language found in EHRs, extracting relevant clinical events and concepts, and mapping these onto a coherent information or treatment timelines which can be difficult to realize manually by humans. LLMs are appropriate for handling the variability and ambiguity that arise in medical documentation, enabling them to identify and organize critical information such as chemotherapy treatments, such as drug names, dosages, administration dates, and associated clinical outcomes (Jahan et al., 2024).

Moreover, by leveraging the latest advancements in transfer learning and domain-specific fine-tuning, LLMs can be programmed in such a way to understand the specific lexicon and data structures unique to domains as complex as oncology and chemotherapy treatment regimes (Chen et al., 2023b).

All in all, this can help with the generation of comprehensive, accurate, and personalized chemotherapy treatment timelines that are an essential component for advancing precision oncology, and also supporting the development and assessment of patient-centric therapeutic strategies (Levine et al., 2019; Banerjee et al., 2019).

To better understand the impact of various factors on tumor behavior and responsiveness, particularly in the context of precision oncology, the Chemotimelines 2024 shared tasks has been proposed (Yao et al., 2024). In this work, we describe our submission to Subtask 1, which aims to build timelines of chemotherapy treatments for individual patients using their Electronic Health Records (EHR) notes. We achieved a 5th place ranking in Subtask 1, with an averaged accuracy across breast, ovarian, and melanoma indications.

The contributions of our paper can be outlined as follows:

1. We developed a Large Language Model (LLM)-based system customized for description and exploration, providing substantial value in tasks related to natural language understanding.

2. We employed multiple LLMs and prompts across diverse development and training datasets, our approach aimed to improve performance and enhance generalization.

3. We introduced a framework that presents a modular strategy for zero-shot relation extraction, leveraging well-established LLMs.

## 2 Related Work

In recent years, there has been a growing body of research demonstrating the effectiveness of Large Language Models (LLMs) in comprehending medical text data and extracting valuable insights from Electronic Health Records (EHRs) across various clinical domains (Beam et al., 2019; Van Veen et al., 2024; Wong et al., 2023; Eriksen and Ryg, 2023). Prior investigations have shown the application of Natural Language Processing (NLP) in healthcare, encompassing tasks such as clinical text classification, medical entity recognition, and patient risk prediction. Efforts to construct clinical timelines from EHR data have predominantly focused on structured data such as procedure codes, diagnosis codes, and laboratory results (Rajkomar et al., 2018; Mullenbach et al., 2018).

Within oncology, NLP methodologies have been employed to analyze cancer-related textual data, including pathology reports, clinical notes, and research articles (Bodenreider, 2004; Meystre et al., 2008). Researchers have investigated the utility of NLP in extracting treatment regimens, identifying adverse drug events, and predicting treatment outcomes among cancer patients (Savova et al., 2010; Xu et al., 2019). Techniques for temporal event extraction and sequence modeling have been explored extensively to develop patient timelines for disease progression tracking and treatment monitoring (Ebadi et al., 2021). Temporal reasoning techniques have found applications in healthcare to analyze the temporal associations between clinical events, treatments, and patient outcomes (Sun et al., 2013). Studies have explored temporal logic, temporal abstraction, and probabilistic models to represent and analyze temporal data in healthcare contexts (Orphanou et al., 2014).

Transformer based large language models have demonstrated remarkable performance improvements across various NLP benchmarks (Devlin et al., 2018; Chiu and Nichols, 2016). Furthermore, healthcare-specific models (Lee et al., 2020) have exhibited state-of-the-art accuracy in biomedical entity recognition (Kocaman and Talby, 2020) and relation extraction (Kocaman and Talby, 2021).

The current state of the art lies in several notable Large Language Models (LLMs), each featuring distinct model architectures and sizes (Pan et al., 2024). Prominent examples include Llama 2 (Touvron et al., 2023), Mistral 7B (Jiang et al., 2023), Zephyr (Tunstall et al., 2023), MEDITRON-70B (Chen et al., 2023d), and Mixtral of Experts (Jiang et al., 2024). LLMs possess the capability to analyze extensive textual data, and the task of summarizing crucial information from electronic health records (Van Veen et al., 2024) can significantly impact how clinicians manage their time, enabling them to dedicate more time to interacting with patients (Khairat et al., 2018) and improve quality of care.

## 3 Task and Dataset Details

### 3.1 Task Details

Chemotimelines 2024 at NAACL-ClinicalNLP Workshop is a shared task (Yao et al., 2024) that focuses on building a timeline of chemotherapy treatment for each patient given all the available Electronic Health Records (EHRs) notes of that patient. The shared task has 2 subtasks. Subtask 1 involves using provided gold annotations of chemotherapy events (EVENTs) and time expressions (TIMEX3s) along with Electronic Health Record (EHRs) notes to predict temporal relations between them and generate patient-level timelines. This task requires deduplicating and resolving conflicts in pairwise temporal relations, with the option to derive timelines without relying on pairwise relations. Additionally, attributes such as modality and relation to document creation time are included. Subtask 2 entails building an end-to-end system for chemotherapy timeline extraction using only patient EHR notes. Both subtasks are evaluated against gold patient-level timelines. We submitted the results of Subtask 1. The submission scripts for evaluation can be found here[1].

---

[1] https://github.com/HealthNLPorg/chemoTimelinesEval

| Indications (Train, Dev and Test) | # Patients | # Reports | # Entities | # Relations |
|---|---|---|---|---|
| Breast (Train) | 23 | 236 | 1599 | 455 |
| Melanoma (Train) | 3 | 32 | 225 | 48 |
| Ovarian (Train) | 14 | 273 | 1765 | 494 |
| Breast (Dev) | 10 | 61 | 425 | 113 |
| Melanoma (Dev) | 2 | 99 | 1050 | 20 |
| Ovarian (Dev) | 8 | 138 | 1102 | 226 |
| Breast (Test) | 25 | 379 | 3678 | 0 |
| Melanoma (Test) | 4 | 77 | 591 | 0 |
| Ovarian (Test) | 6 | 143 | 1045 | 0 |
| **Total** | 95 | 1438 | 11480 | 1356 |

Table 1: Summary of Dataset Statistics: Number of Patients, Reports, Entities, and Relations across Training (Train), Development (Dev) and Testing (Test) Sets for different indications (Breast, Melanoma and Ovarian).

| Indications (Train and Dev) | BEGINS-ON | CONTAINS | ENDS-ON |
|---|---|---|---|
| Breast (Train) | 131 | 298 | 26 |
| Melanoma (Train) | 10 | 37 | 1 |
| Ovarian (Train) | 101 | 327 | 66 |
| Breast (Dev) | 27 | 57 | 29 |
| Melanoma (Dev) | 42 | 157 | 2 |
| Ovarian (Dev) | 34 | 140 | 52 |

Table 2: Summary of Dataset Statistics: Indications (breast, melanoma, and ovarian) across training and development sets, including the three types of temporal relations.

## 3.2 Dataset

The dataset comprises 95 patients with 1438 reports. Table 1 summarizes dataset statistics, including indications (breast, melanoma and ovarian) for training, development, and testing sets, along with the number of patients, reports, entities, and relations. The annotated dataset has been using THYME ontology (Styler IV et al., 2014) and temporal relation annotations (Wright-Bettner et al., 2020) with three different temporal relations used for TLINKs (temporal links): BEGINS-ON, CONTAINS and ENDS-ON. Table 2 presents summarized statistics for indications (breast, melanoma, and ovarian) across training and development sets, including the three types of temporal relations.

## 4 Approach

We aimed to significantly contribute to the development of advanced cutting-edge methodologies and techniques for automatically constructing chemotherapy treatment timelines from Electronic Health Records (EHRs) clinical notes of individual patients. We leveraged current state of the art Large Language models (LLMs) for this shared task. We tested various LLMs with different sizes and archi-

tectures to determine which model works best for relation extraction (See Figure 1).

### 4.1 Natural Language Processing (NLP) Pipeline with Language Representations

#### 4.1.1 Document Chunking

We divided the documents into paragraphs or groups of paragraphs (sections) to facilitate manageable processing units.

**Sequence Length** The experiments involved evaluating various sequence lengths, which determine the number of words or tokens processed by the model at once. Assessing lengths of 1024, 512, and 256 tokens provides insights into how input length impacts the system's accuracy in extracting relations.

**Paragraph and Sentence Detection Paragraph Detection** NLP plays a crucial role in enhancing contextual understanding within Electronic Health Records (EHRs) by segmenting the text into meaningful units. By identifying paragraphs, NLP models can discern distinct sections of the EHRs, such as patient history, symptoms, diagnoses, and treatment plans. This segmentation enables the model

Figure 1: NLP Pipeline for Subtask 1

to focus on specific aspects of the patient's medical information, facilitating more accurate analysis and interpretation. We have incorporated section chunking and paragraph detection techniques into our system. This involves identifying individual sentences within the text data. By isolating paragraphs, the system can focus on extracting relations specifically from relevant pairs of entities within each paragraph, which enhances precision. We have incorporated section chunking and paragraph detection techniques into our system. This involves identifying individual sentences and paragraphs within the text data. By isolating paragraphs, the system can focus on extracting relations specifically from relevant pairs of entities within each paragraph (Kocaman and Talby, 2020), which enhances precision.

- In terms of extracting relations from various document paragraphs, our sequences already extend beyond a single paragraph, as our sequence length is configured to accommodate 256 tokens. Nonetheless, such occurrences are rare within this dataset. If necessary, we can concatenate adjacent or contiguous paragraphs or clusters of paragraphs to enable the extraction of relations spanning multiple paragraphs.

- To address chunking concerns, we implemented an overlap parameter for enhanced performance. This parameter prevents the inadvertent separation of essential information

by preserving sentence integrity, even without overlap. It facilitates the reconciliation of fragmented data, mitigating the risk of context loss and preserving predictive accuracy. The risk of reduced recall arises from potential pairs not being prompted for relation classification. Encouragingly, the model's metrics exhibit no specific recall-related issues, signaling positive performance in this regard.

### 4.1.2 Zero-Shot Prompting for Related Pairs

We developed structured prompts to guide the system in identifying and extracting relations between pairs of entities. These prompts serve as cues for the system to recognize and analyze relevant information in the text pertaining to the specified entities (See Figure 2).

This process involved prompt engineering techniques aimed at refining the instructions within the relation extraction pipeline, optimizing them to extract more precise and relevant information during subsequent stages. The zero-shot prompt gave us a reasonably high precision by leveraging the prompt templates that guided the LLMs to generate responses that closely match the desired output without requiring explicit training data for each class or category. Prompt 1 was used for the submission and for evaluation. We tried Prompt 2 but encountered challenges in labeling the relations from distinct lists. (See Figure 2)

The input to the LLMs involves combining the prompt with the paragraph or groups of paragraphs (sections).

397

### 4.1.3 Tokenization and Embedding

Each paragraph is tokenized, and the tokens are encoded using the tokenizer specific to the chosen Large Language Model (LLM).

### 4.1.4 Embedding Decoding

The encoded tokens were fed to the LLM, resulting in serialized outputs.

### 4.1.5 Semantic Object Construction

Using the outputs from the LLMs and predefined validation classes for each of the prompts, we construct semantic-rich objects that encapsulate the information extracted from the text.

**Directed Acyclic Graph (DAG)** We constructed a simplified DAG to outline the logical framework guiding the construction of the output taxonomy, enabling a structured representation of the reasoning process. (See Figure 1)

After establishing the relations with the output from the LLMs, we leveraged Pydantic (Colvin and contributors, 2024), a Python library for data validation and settings management. Pydantic (Colvin and contributors, 2024) facilitates data parsing and validation, ensuring that the data adheres to the expected types specified using Python's standard type hints. A Directed Acyclic Graph (DAG) can impact model accuracy positively by ensuring that validation functions are executed in a specific, predictable order. This helped in maintaining data integrity and correctness, thereby reducing the likelihood of errors or inconsistencies in the model's predictions. Additionally, DAGs prevent cyclic dependencies, which led to more stable and reliable model behavior.

**Date Normalization** We normalized the data using both the natural language representation for the temporal entity and the document time as a reference. We then transformed the temporal entity to an absolute datetime.

- The date normalization process is integrated into the validation procedure through a dedicated class field validator. It involves multiple steps to handle various date formats and potential failure scenarios, including cases where external services like Duckling (Rasa, 2024) may not parse the input successfully.

- Initially, the validator attempts to parse the raw date string using the parse_timex func-

tion, which sends the string to Duckling (Rasa, 2024) and, if unsuccessful, to the SparkNLP date normalizer annotator (John Snow Labs, 2024). These tools excel at interpreting natural language and complex date expressions, providing robust initial parsing. If successful, the parsed value undergoes further processing with dateutil to ensure compatibility with Python's datetime object format.

- In case of failure with Duckling (Rasa, 2024) and SparkNLP parsing (John Snow Labs, 2024), the validator employs fallback strategies. It checks for ISO week format dates and year-month-only strings, attempting to convert them into complete dates. If these strategies fail, the validator employs a battery of parsers (e.g., dateutil_parser.parse, pd.to_datetime, arrow.get) in a loop until successful parsing occurs.

- Throughout the process, detailed logging captures various states and errors, aiding in debugging and understanding parsing issues. Finally, if a valid date is obtained through any of these strategies, it is stored as the normalized value in the model, which may represent a full date or just the year and month, depending on the input string and specified context.

### 4.1.6 Serialization for Submission

Finally, we aggregated and serialized these semantic objects into the submission format specified by the competition guidelines.

### 4.2 Baseline Models

We fine-tined pre-trained Llama 2 (Touvron et al., 2023) and Mistral 7B (Jiang et al., 2023) for our submission to this shared task.

**Llama 2** Llama 2 (Touvron et al., 2023) is a collection of large language models (LLMs) ranging from 7 billion to 70 billion parameters. They are fine-tuned LLMs optimized for dialogue applications.

**Mistal 7B** Mistral 7B (Jiang et al., 2023) is a language model consisting of 7 billion parameters designed to deliver superior performance and efficiency. Mistral 7B demonstrates superior performance compared to the best open 13B model (Llama 2) (Touvron et al., 2023) across all assessed benchmarks and outperforms the leading released

Figure 2: Prompt template: Prompt 1 for the relation label from the pairs (left) and Prompt 2 for the relation label from the separate lists of drugs and dates (right).

34B model (Llama 1) in tasks such as reasoning, mathematics, and code generation.

### 4.2.1 Validation and Quality of LLM Response

We rewrote the THYME ontology on top of a typed validation framework based on (pydantic (Colvin and contributors, 2024)) library. We binded every result from a prompted task to one of these objects: Thyme; in Subtask 1 the validation class for the LLM response is the graph representation defined in `TypedTimedEvents:List[Tuple[Event, Timex, TLinkType]]`. We forced the output from the LLM to conform to this type, and if not we kept refining the prompt. After obtaining accurately processed outputs from the LLM, the next step involved aggregation. This entails concatenating the parsed subgraphs from each chunk of the LLM output into a deduplicated timeline at the patient level.

During the inference phase, we focused on the post processing techniques, such as parsing and refining, applied to the output generated by the Large Language Models (LLMs). These techniques aim to enhance the quality and accuracy of the extracted information, ensuring its suitability for downstream analysis and applications.

### 4.3 Evaluation Metrics

Models were evaluated with the official evaluation script[2] on the test set. The following metrics were used: Precision, Recall and F-score (Hossin and Sulaiman, 2015). We reported performance as the arithmetic mean of F-score.

### 4.4 Human Evaluation

In our study, two medical professionals conducted a comparative analysis of chemotherapy timelines generated by LLMs, specifically using the Llama 2 model for our initial submission, against a ground truth established by the dataset (train and dev set) provided by the challenge. The dataset combines training, development, and testing sets, encompassing a total of ninety five (n=95) patients. Train and dev set contain sixty patients (n=60) patients and the test set contains thirty five (n=35) patients. The gold standard for the test set of thirty five (n=35) patients was not released. Therefore, the two medical professionals randomly selected five patients (n = 5) from each indication (breast, melanoma and ovarian) and manually reviewed the predictions generated by the LLMs, performed a qualitative evaluation.

The LLMs demonstrated a tendency to misclas-

| Indications (Train and Dev) | Baseline | Predictions | Llama 2 | Mistral 7B |
|---|---|---|---|---|
| Breast (Train) | 0.427713 | 0.800827 | 0.695125 | 0.606543 |
| Breast (Dev) | **0.863988** | **0.888878** | 0.768916 | 0.723611 |
| Melanoma (Dev) | 0.455782 | 0.797009 | 0.633271 | 0.767574 |
| Melanoma (Train) | 0.765196 | 0.842803 | **0.882037** | **0.799432** |
| Ovarian (Dev) | 0.715926 | 0.607934 | 0.561085 | 0.625625 |
| Ovarian (Train) | 0.715137 | 0.816064 | 0.647571 | 0.595842 |

Table 3: Performance on relation extraction by approach.

| Parameters | Value | Description |
|---|---|---|
| Chunk Size | 256 | number of tokens or words processed at a time during training or inference |
| Temperature | 0.1 | controls the randomness of the generated output |
| Seed | 123 | predefined starting point for the random no. generator used during training |

Table 4: Hyperparameters used with LLama 2 (Touvron et al., 2023) and Mistral 7B (Jiang et al., 2023) for the Chemotimelines 2024 Subtask 1.

sify CONTAINS relation over the BEGINS-ON and ENDS-ON, resulting in low recall for BEGINS-ON and ENDS-ON, and low precision for CONTAINS. For instance, in one case, where the actual relationship indicated Taxotere ENDS-ON at a specific date, the model incorrectly predicted it as a CONTAINS relation.

Another noteworthy observation was the occasional complete oversight of a relation by the LLMs. Additionally, discrepancies arose when the year was occasionally misinterpreted as a future date due to errors in the dates mentioned in the reports.

## 5 Results and Discussion

As previously stated, our study utilizes the two large language models, Llama 2 (Touvron et al., 2023) and Mistral 7B (Jiang et al., 2023). Table 3 shows the performance metrics on the relation extraction NLP task for the training and development set across three indications (Breast, Melanoma and Ovarian). We evaluated the performance of both Llama 2 and Mistral 7B against the baseline. Notably, we attained the highest performance on the Melanoma training set with both Llama 2 and Mistral 7B.

We utilized the default parameters for both Llama 2 and Mistral 7B, except for the chunk size, which was set to 256, temperature set to 0.1, and seed set to 123 (Refer to Table 4). Chunk size refers to the number of tokens or words processed at a time during training or inference. Temperature regulates the randomness of the generated output, while the seed serves as the predefined starting

point for the random number generator used during model training.

Table 5 illustrates the results of three test data runs utilizing Llama 2 and Mistral 7B for Subtask 1. Our highest performing model was Llama 2, achieving an F1 average score of 0.71, while Mistral 7B attained an average F1 of 0.61. Llama 2 exhibited superior performance compared to Mistral 7B, resulting in a higher rank. Specifically, Llama 2 secured the 5th position in the average score for Subtask 1, the 4th position for the Melanoma indication, and the 7th position for Breast and Ovarian indications.

**Error Analysis** Error analysis in Large Language Models (LLMs) involves scrutinizing the model's prediction errors to discern their types, frequency, and underlying causes. This entails evaluating the model's performance on a test dataset and categorizing errors into various types, including false positives, false negatives, ambiguous cases, out-of-distribution errors, and conceptual errors. By analyzing these errors, insights can be gleaned regarding patterns and areas for improvement in the model. This analysis guides strategies for enhancing the model's performance through fine-tuning, refining training data, and optimizing input representations. Furthermore, error analysis is crucial for establishing confidence in the model's predictions and comprehending its limitations in real-world scenarios.

Figure 3 shows the error analysis presented compares Llama 2 and Mistral for the baseline established by the organizers, as well as prediction

| Runs | LLMs | Average Score | Breast | Melanoma | Ovarian |
|------|------|---------------|--------|----------|---------|
| Run 1 | Llama 2 | 0.71 | 0.68 | 0.83 | 0.61 |
| Run 2 | Llama 2 | 0.68 | 0.66 | 0.80 | 0.59 |
| Run 3 | Mistral 7B | 0.61 | 0.62 | 0.59 | 0.62 |

Table 5: Results of Runs on Test Data for Subtask 1.



Figure 3: Error Analysis for Subtask 1.

scores derived by directly utilizing the golden relations as the timeline. The results from Llama 2 and Mistral 7B are based on the question answering prompting approach used to generate our timelines.

The metrics reveal lower precision within the system, characterized by exceptionally high recall. Further investigation into the distribution of false positives across event types or relation categories may unveil discernible patterns. It appears that the Large Language Model (LLM) is indiscriminately predicting all instances as if they are related events to timelines.

## 6 Conclusion and Future Work

In this paper, we present our submission to the Chemotimelines 2024 shared tasks (Yao et al., 2024) to build chemotherapy treatment timelines using Electronic Health Records (EHRs) notes from various sources, such as primary care

providers, oncology departments, discharge summaries, emergency department, pathology, radiology, and more. We used zero shot prompted relation extraction (Wang et al., 2023; Jun and et al., 2022) driven by the THYME ontology (Styler IV et al., 2014) and temporal relation annotations (Wright-Bettner et al., 2020).

We evaluated pre-trained Large Language Models (LLMs) like Llama 2 (Touvron et al., 2023), Mistral 7B (Jiang et al., 2023), Zephyr (Tunstall et al., 2023), MEDITRON-70B (Chen et al., 2023d), and Mixtral of Experts (Jiang et al., 2024) with different sizes and architectures. We only reported results on Llama 2 (Touvron et al., 2023) and Mistral 7B (Jiang et al., 2023). We conducted a series of experiments with different setups to improve the system's performance. From our analysis, we conclude that our approach helped us determine which model works best for this shared task. We

conclude that LLMs provide a promising path forward for extracting timelines that contextualize cancer treatment, which were previously unavailable. We also show that our model provides high recall that is beneficial for instances where high sensitivity is required such as with output-sensitive predictions like cancer prediction models. However, due to the auto-generation approach and minimizing human intervention, the models we developed demonstrated relative low precision. Precision was evaluated with two physicians validating the accuracy of the generated chemotherapy timelines.

Our proposed methodology represents a significant advancement in the field, providing a flexible and efficient solution for relation extraction tasks in natural language processing. Large Language Models (LLMs) offer a promising approach for auto-generating chemotimelines from Electronic Health Records (EHRs) due to their advanced natural language understanding capabilities, contextual understanding, and semantic representation of medical information. LLMs can comprehend complex medical texts, capture contextual relationships between different medical events, and generate rich semantic representations of medical concepts and events mentioned in EHRs. As we see in our current study, our effort to attempt fully auto-generated chemotherapy timelines have shown great promise in terms of recall but have a negative impact on precision. In future studies we will explore further training rounds or human-in-the-loop models to explore the right balance between automation and human guided outputs. Nevertheless, our study demonstrates great promise in integrated LLM-generated chemotherapy timelines that have the potential to alleviate documentation and data harmonization burdens, potentially easing clinician workload and enhancing quality of patient care.

Exploring the potential of LLMs is an emerging area in research. We have experimented with two state-of-the-art LLMs (Llama 2 and Mistral 7B) for this task, comparing each with the gold standard for various cancer types. Our approach was to maintain a broad, domain-agnostic perspective, treating it as a high-level NLP relation detection task. We assumed that the underlying LLMs were general-purpose. In the future, we aim to explore domain-specific LLMs tailored for biomedical texts, such as JSL-MedMNX-7B (JSL-Med-Sft-Llama-3-8B, 2024), which could offer improved accuracy by better handling specialized language and data structures inherent in this domain.

Furthermore, we aim to validate the effectiveness of our LLM-based system across diverse healthcare datasets to enhance its performance. Additionally, we intend to conduct comprehensive analysis of the generated chemotherapy timelines to fine-tune them further and improve precision. This includes conducting in-depth error analyses to pinpoint the root causes of false positives. Our goal is to identify any consistent patterns, words, or phrases that the model may misinterpret, facilitating targeted improvements to enhance its accuracy.

## Limitations

While leveraging Large Language Models (LLMs) for creating chemotherapy timelines from clinical notes offers numerous benefits, it also presents several limitations: 1. The accuracy and reliability of generated timelines heavily depend on the quality and consistency of input clinical notes, potentially leading to inaccuracies or omissions. 2. LLMs may exhibit biases inherent in the training data, leading to disparities, inaccuracies or generalization in the generated timelines, especially when applied to diverse patient populations. 3. LLMs are complex models with billions of parameters, making it challenging to interpret their decision-making processes, limiting clinicians' ability to trust and validate the generated outputs. 4. Training and fine-tuning LLMs for healthcare applications, including generating chemotherapy timelines, require significant computational resources, expertise, and time. Due to time constraints, we investigated a narrow range of models and hyperparameter configurations. Given their demonstrated proficiency in natural language processing, these models serve as an ideal starting point for extracting pertinent clinical events and concepts essential for constructing treatment timelines. 5. Despite the automation capabilities of LLMs, human oversight and validation are still essential to ensure the accuracy and relevance of the generated chemotherapy timelines. Clinicians must review and validate the outputs to identify and correct any inaccuracies or inconsistencies. In our study, two medical professionals compared chemotherapy timelines generated by LLMs, particularly the Llama 2 (Touvron et al., 2023) model, with a ground truth dataset provided by the challenge.

## Ethics Statement

Leveraging Large Language Models (LLMs) for constructing timelines of chemotherapy treatments using Electronic Health Records (EHR) notes raises numerous ethical considerations. Foremost among these is the imperative to safeguard patient privacy and confidentiality, given the sensitive nature of personal health information stored in EHRs. By leveraging openly available LLMs, physicians can inadvertently expose patient data to private companies (Blease, 2024). Robust data security measures, and digital literacy training is essential to thwart unauthorized patient data exposure to LLMs or data breaches, thereby averting potential cyber threats. Additionally, obtaining informed consent from patients regarding the utilization of their health data is paramount to uphold patient autonomy and foster transparency. Ensuring the accuracy and integrity of the data is vital to mitigate risks of erroneous treatment timelines that could lead to patient harm. Moreover, LLMs may perpetuate biases inherent in the data, thereby introducing disparities or unfairness in the generated timelines (Singh et al., 2023). Prioritizing algorithmic transparency and accountability is imperative to identify and mitigate biases in LLM decision-making processes. Furthermore, granting patients control over their health data, including access and consent for research or analytical purposes, is fundamental in upholding patient autonomy and fostering trust in the healthcare system. The organizers of the Chemotimelines 2024 at NAACL-ClinicalNLP Workshop shared tasks (Yao et al., 2024) have provided a de-identified dataset.

In leveraging Large Language Models (LLMs) for Open Book Question Answering (QA), it's crucial to address the potential ethical concerns surrounding the minimization of generation divergence risk. This entails ensuring that the responses generated by LLMs align closely with the intended context and accurately reflect the information available in the open book. By minimizing generation divergence risk, we aim to uphold the integrity of the QA process, promote transparency, and mitigate the dissemination of misinformation or biased responses. Additionally, efforts should be made to continually evaluate and refine LLMs to enhance their reliability and trustworthiness in providing accurate and contextually appropriate answers.

It is noteworthy that LLMs often demonstrate a propensity to produce hallucinations when generating coherent answers, underscoring the necessity for human supervision in their utilization. Ensuring human supervision during the deployment of LLMs in healthcare contexts is crucial to validate the accuracy, appropriateness and potential harmfulness of the generated outputs and to mitigate potential risks or errors (Chen et al., 2023a). Moreover, it is crucial to recognize that the present system serves as an experimental tool intended to catalyze further research, including additional fine-tuning and model explainability studies. Such endeavors are indispensable before these systems can be safely incorporated into clinical settings, ensuring their reliability and efficacy in supporting clinical decision-making processes. Additionally, another critical aspect deserving careful consideration is the explainability and interpretability of Language Models (LLMs) when deployed in healthcare contexts.

## References

Gaurav Anand and Divya Sadhna. 2023. Electronic health record interoperability using fhir and blockchain: A bibliometric analysis and future perspective. *Perspectives in Clinical Research*.

Imon Banerjee, Selen Bozkurt, Jennifer Lee Caswell-Jin, Allison W Kurian, and Daniel L Rubin. 2019. Natural language processing approaches to detect the timeline of metastatic recurrence of breast cancer. *JCO clinical cancer informatics*, 3:1–12.

Andrew L Beam, Benjamin Kompa, Allen Schmaltz, Inbar Fried, Griffin Weber, Nathan Palmer, Xu Shi, Tianxi Cai, and Isaac S Kohane. 2019. Clinical concept embeddings learned from massive sources of multimodal medical data. In *Pacific Symposium on Biocomputing 2020*, pages 295–306. World Scientific.

Danielle S Bitterman, Eli Goldner, Sean Finan, David Harris, Eric B Durbin, Harry Hochheiser, Jeremy L Warner, Raymond H Mak, Timothy Miller, and Guergana K Savova. 2023. An end-to-end natural language processing system for automatically extracting radiation therapy events from clinical texts. *International Journal of Radiation Oncology* Biology* Physics*, 117(1):262–273.

Charlotte Blease. 2024. Open AI meets open notes: surveillance capitalism, patient privacy and online record access. *Journal of Medical Ethics*, 50(2):84–89.

Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270.

Shan Chen, Benjamin H Kann, Michael B Foote, Hugo JWL Aerts, Guergana K Savova, Raymond H Mak, and Danielle S Bitterman. 2023a. Use of artificial intelligence chatbots for cancer treatment information. *JAMA oncology*, 9(10):1459–1462.

Shan Chen, Benjamin H Kann, Michael B Foote, Hugo JWL Aerts, Guergana K Savova, Raymond H Mak, and Danielle S Bitterman. 2023b. The utility of chatgpt for cancer treatment information. *MedrXiv*, pages 2023–03.

Shan Chen, Yingya Li, Sheng Lu, Hoang Van, Hugo JWL Aerts, Guergana K Savova, and Danielle S Bitterman. 2023c. Evaluation of chatgpt family of models for biomedical reasoning and classification. *arXiv preprint arXiv:2304.02496*.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. 2023d. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*.

Jason P Chiu and Eric Nichols. 2016. Named entity recognition with bidirectional lstm-cnns. *Transactions of the Association for Computational Linguistics*, 4:357–370.

Samuel Colvin and contributors. 2024. Pydantic: Data validation and settings management library for python. https://github.com/samuelcolvin/pydantic.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Ashkan Ebadi, Pengcheng Xi, Stéphane Tremblay, Bruce Spencer, Raman Pall, and Alexander Wong. 2021. Understanding the temporal evolution of covid-19 research through machine learning and natural language processing. *Scientometrics*, 126:725–739.

Sören Möller Eriksen, Alexander V. and Jesper Ryg. 2023. Use of gpt-4 to diagnose complex clinical cases. *NEJM AI*.

Marco Guevara, Shan Chen, Spencer Thomas, Tafadzwa L Chaunzwa, Idalid Franco, Benjamin H Kann, Shalini Moningi, Jack M Qian, Madeleine Goldstein, Susan Harper, et al. 2024. Large language models to identify social determinants of health in electronic health records. *NPJ digital medicine*, 7(1):6.

Harry Hochheiser, Sean Finan, Zhou Yuan, Eric B Durbin, Jong Cheol Jeong, Isaac Hands, David Rust, Ramakanth Kavuluru, Xiao-Cheng Wu, Jeremy L Warner, et al. 2023. Deepphe-cr: Natural language processing software services for cancer registrar case abstraction. *JCO Clinical Cancer Informatics*, 7:e2300156.

Mohammad Hossin and Md Nasir Sulaiman. 2015. A review on evaluation metrics for data classification evaluations. *International journal of data mining & knowledge management process*, 5(2):1.

Israt Jahan, Md Tahmid Rahman Laskar, Chun Peng, and Jimmy Xiangji Huang. 2024. A comprehensive evaluation of large language models on benchmark biomedical text processing tasks. *Computers in Biology and Medicine*, page 108189.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

John Snow Labs. 2024. Spark nlp - date normalizer annotator documentation. https://nlp.johnsnowlabs.com/licensed/api/python/modules/sparknlp_jsl/annotator/normalizer/date_normalizer.html.

JSL-Med-Sft-Llama-3-8B. 2024. John Snow Labs jsl-med-sft-llama-3-8b. https://huggingface.co/johnsnowlabs/JSL-Med-Sft-Llama-3-8B.

Zhao Jun and et al. 2022. An exploration of prompt-based zero-shot relation extraction method. In *Proceedings of the 21st Chinese National Conference on Computational Linguistics*.

Kenneth L Kehl, Wenxin Xu, Eva Lepisto, Haitham Elmarakeby, Michael J Hassett, Eliezer M Van Allen, Bruce E Johnson, and Deborah Schrag. 2020. Natural language processing to ascertain cancer outcomes from medical oncologist notes. *JCO Clinical Cancer Informatics*, 4:680–690.

Saif Sherif Khairat, Aniesha Dukkipati, Heather Alico Lauria, Thomas Bice, Debbie Travers, and Shannon S Carson. 2018. The impact of visualization dashboards on quality of care and clinician satisfaction: Integrative literature review. *JMIR Hum Factors*, 5(2):e22.

V Kocaman and D Talby. 2020. Biomedical named entity recognition at scale. *Pattern Recognition. ICPR International Workshops and Challenges*.

Veysel Kocaman and David Talby. 2021. Spark nlp: Natural language understanding at scale. *Software Impacts*.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: A pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Mark N Levine, Gordon Alexander, Arani Sathiyapalan, Anjali Agrawal, and Greg Pond. 2019. Learning health system for breast cancer: pilot project experience. *JCO clinical cancer informatics*, 3:1–11.

Stéphane M Meystre, Guergana K Savova, Karin C Kipper-Schuler, and John F Hurdle. 2008. Extracting information from textual documents in the electronic health record: a review of recent research. *Yearbook of medical informatics*, 17(01):128–144.

James Mullenbach, Sarah Wiegreffe, Jon Duke, Jimeng Sun, and Jacob Eisenstein. 2018. Explainable prediction of medical codes from clinical text. *arXiv preprint arXiv:1802.05695*.

Marjan Najafabadipour, Massimiliano Zanin, Alejandro Rodríguez-González, Maria Torrente, Beatriz Nuñez García, Juan Luis Cruz Bermudez, Mariano Provencio, and Ernestina Menasalvas. 2020. Reconstructing the patient's natural history from electronic health records. *Artificial intelligence in medicine*, 105:101860.

Kalia Orphanou, Athena Stassopoulou, and Elpida Keravnou. 2014. Temporal abstraction and temporal bayesian networks in clinical domains: A survey. *Artificial intelligence in medicine*, 60(3):133–149.

Shirui Pan, Jie Wang, Chuxu Zhang, and Jiawei Han. 2024. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*.

Alvin Rajkomar, Eyal Oren, Kai Chen, Andrew M Dai, Nissan Hajaj, Michaela Hardt, Peter J Liu, Xiaobing Liu, Jake Marcus, Mimi Sun, et al. 2018. Scalable and accurate deep learning with electronic health records. *NPJ digital medicine*, 1(1):1–10.

Rasa. 2024. Duckling: Open-source library for parsing structured information from text. https://github.com/facebook/duckling.

Miriam Reisman. 2017. Ehrs: the challenge of making electronic data usable and interoperable. *Pharmacy and Therapeutics*, 42(9):572.

Guergana K Savova, James J Masanz, Philip V Ogren, Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-Schuler, and Christopher G Chute. 2010. Mayo clinical text analysis and knowledge extraction system (ctakes): architecture, component evaluation and applications. *Journal of the American Medical Informatics Association*, 17(5):507–513.

Nina Singh, Katharine Lawrence, Safiya Richardson, and Devin M Mann. 2023. Centering health equity in large language model deployment. *PLOS Digital Health*, 2(10):e0000367.

William F. Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, et al. 2014. Temporal annotation in the clinical domain. *Transactions of the association for computational linguistics*, 2:143–154.

Weiyi Sun, Anna Rumshisky, and Ozlem Uzuner. 2013. Temporal reasoning over clinical text: the state of the art. *Journal of the American Medical Informatics Association*, 20(5):814–819.

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*.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*.

Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerová, et al. 2024. Adapted large language models can outperform medical experts in clinical text summarization. *Nature Medicine*, pages 1–9.

Zhiqiang Wang, Yiran Pang, and Yanbin Lin. 2023. Large language models are zero-shot text classifiers. *arXiv preprint arXiv:2312.01044*.

Cliff Wong, Sheng Zhang, Yu Gu, Christine Moung, Jacob Abel, Naoto Usuyama, Roshanthi Weerasinghe, Brian Piening, Tristan Naumann, Carlo Bifulco, et al. 2023. Scaling clinical trial matching using large language models: A case study in oncology. In *Machine Learning for Healthcare Conference*, pages 846–862. PMLR.

Kristin Wright-Bettner, Guergana Savova, and Steven Bethard. 2020. Defining and learning refined temporal relations in the clinical narrative. In *Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis*.

Yanjun Xu, Qun Dong, Feng Li, Yingqi Xu, Congxue Hu, Jingwen Wang, Desi Shang, Xuan Zheng, Haixiu Yang, Chunlong Zhang, et al. 2019. Identifying subpathway signatures for individualized anticancer drug response by integrating multi-omics data. *Journal of translational medicine*, 17:1–16.

Jiarui Yao, Harry Hochheiser, WonJin Yoon, Eli Goldner, and Guergana Savova. 2024. Overview of the 2024 shared task on chemotherapy treatment timeline extraction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, NAACL June, 2024, Mexico City, Mexico.

# Team NLPeers at Chemotimelines 2024:
# Evaluation of two timeline extraction methods,
# can generative LLM do it all or is smaller model fine-tuning still relevant ?

**Nesrine Bannour[1], Judith Jeyafreeda Andrew[1,2], Marc Vincent[1]**

[1]Université de Paris, Imagine Institute, Data Science Platform,
INSERM UMR 1163, F-75015, Paris, France,
[2]PaRis Artificial Intelligence Research InstitutE (PRAIRIE), Paris, France
firstname.lastname[at]institutimagine.org

## Abstract

This paper presents our two deep learning-based approaches to participate in subtask 1 of the Chemotimelines 2024 Shared task. The first uses a fine-tuning strategy on a relatively small general domain Masked Language Model (MLM) model, with additional normalization steps obtained using a simple Large Language Model (LLM) prompting technique. The second is an LLM-based approach combining advanced automated prompt search with few-shot in-context learning using the DSPy framework. Our results confirm the continued relevance of the smaller MLM fine-tuned model. It also suggests that the automated few-shot LLM approach can perform close to the fine-tuning-based method without extra LLM normalization and be advantageous under scarce data access conditions. We finally hint at the possibility to choose between lower training examples or lower computing resources requirements when considering both methods.

## 1 Introduction

The advent of auto-regressive Large Language Models (LLMs) has taken the NLP field by storm and has been diffusing to more specialized domains, such as clinical NLP ever since. While the most powerful models are still only available as private owned services - oftentimes precluding their use with sensitive medical data - open source and open weight models have been catching up, mostly since the release of the LLaMA model family (Touvron et al., 2023). With such open models, in-context learning strategies became more viable. On top of those LLMs, an ecosystem of tools and frameworks has also emerged to provide more robust and efficient ways to use them. One such framework is DSPy (Khattab et al., 2023), whose ambition is to provide a principled and automated way to search LLM prompts and weights and ultimately build robust LLM pipelines.

While this latter technology still evolves, older and more established deep learning models coexist, and the comparative advantages of the two approaches are being assessed. A prominent example of those predecessors is BERT-based models, which can still be considered LLMs, although they are usually an order of magnitudes smaller than their auto-regressive counterparts. With such Masked Language Models (MLMs), fine-tuning of the model weights can be more easily performed due to their usually smaller size.

Temporal Relation Extraction (TRE) is a crucial task for several domains, particularly the clinical domain, requiring a deep understanding of natural language. With the rise of LLMs, recent research efforts attempt to apply these models to the TRE task, but results are still debatable (Han et al., 2023; Chen et al., 2023; Li et al., 2023a). In this paper, we address the clinical event-to-time expression relation extraction task by evaluating two timeline extraction methods.

The main contributions of this paper are:

- An MLM-based fine-tuning approach using a relatively light state-of-the-art MLM model.

- An automated few-shot prompting approach with an LLM using the DSPy framework.

- An evaluation and comparison of these two TRE methods, as well as two temporal expressions normalization methods: a pre-existing tool and a proposed LLM-based method.

## 2 Related Work

**Rule Based Methods.** Several research papers (Gaizauskas et al., 2006; Zhou et al., 2008; Hernández et al., 2016; Wang et al., 2016) used rule based approaches for TLINK classification. Zhou et al. (2008) and Hernández et al. (2016) used external clinical domain knowledge to improve the rules.
**Machine learning methods.** Research efforts

using Machine Learning approaches for TRE included the use of Support Vector Machine (SVMs) (Lee et al., 2016; Khalifa et al., 2016); Conditional Random Fields (CRFs) (Khalifa et al., 2016); Convolutional Neural Networks (CNNs) (Li and Huang, 2016; Chikka, 2016), and Bi-LSTMs (Tourille et al., 2017a).

**Hybrid Approaches.** Tang et al. (2013) proposed a hybrid method using a combination of SVM and CRF techniques, with rules to resolve conflicting cases. Nikfarjam et al. (2013) used a SVM with a sentence-level graph-based inference mechanism. Tourille et al. (2017b) used a SVM with word embeddings approach to extract temporal relations.

**Language Models.** Huguet Cabot and Navigli (2021) presented REBEL, which is a seq2seq model using the BART model as the base model for end-to-end relation extraction. REBEL takes as input raw input and outputs a set of triplets with relations and entities that have been linearized. Eberts and Ulges (2019) used pre-trained BERT as a base model. Entities are detected among all token spans. Entities with no relations are filtered out, and the remaining entities and their relations are classified. Lin et al. (2021) proposed the Entity-BERT model obtained with continued pre-training on PubMedBERT base uncased with MIMIC-BIG and MIMIC-SMALL using Entity-Centric masking. The authors then fine-tune EntityBERT for several tasks, including TRE.

**Prompt Learning** With the increased use of LLMs, prompt Learning has gained popularity. Within this context, several prompting techniques have been proposed using prompt templates (Jiang et al., 2020; Shin et al., 2020; Liu et al., 2023; Li and Liang, 2021; Lester et al., 2021). Few-shot prompting can be used to enable in-context learning, where we provide demonstrations of the prompt to steer the model to better performance. Frameworks such as DSPy (Khattab et al., 2023) allow for optimized few-shot prompting approaches.

## 3 Task description and data

We participated in the first subtask of the Chemotherapy Treatment Timelines Extraction Shared Task[1] (Yao et al., 2024), which aims to extract temporal relations between chemotherapy events and time expressions and then produce the final patient-level timelines by resolving duplica-

tion and conflicts in the pairwise temporal relations. The types of relations to extract are mainly CONTAINS, BEGINS-ON, and ENDS-ON. The data provided by the University of Pittsburgh/UMPC, through a Data Use Agreement during this shared task, includes a list of available de-identified Electronic Health Record (EHR) notes for patients with breast, ovarian, and melanoma cancer. Further details about the subtask and the data distribution are described in Yao et al. (2024). The organizers provide a baseline system based on the Entity-BERT model (Lin et al., 2021).

## 4 Methods

This section describes our proposed methods for the TRE task and the post-processing, normalization, and summarization steps used to construct patient-level timelines.

### 4.1 MLM fine-tuning

**Model fine-tuning** As a core model for the fine-tuning approach, we use DeBERTa-v3 base (He et al., 2021), which is a relatively light state-of-the-art 86 million parameters MLM initially trained on 160 GB of general domain text data.

The MLM was fine-tuned on a $(event, time)$ pair multi-class classification task, with processed examples coming from the gold entities and relations dataset provided in the contest training set (which -for the purpose of fine-tuning- was subdivided into a training set and validation set based on which epoch selection was done). The fine-tuned model was then tested on the contest's development set.

The finetuning was done using the huggingface's transformer library using a multiclass classification setup. Given time constraints, a single set of hyper-parameters was used for the training and given to the Trainer class of huggingface's library. Learning rate was set to 2e-5 with a weight decay of 0.01, the maximum number of epochs was set to 10 (with an epoch evaluation strategy). The label was one of: $\{begins\_at, ends\_at, contains\_1, no\_link\}$, $no\_link$ indicating an absence of a relation between the *event* and *time* entity. After the classification of each candidate pair, the ones predicted to be $no\_link$ (i.e. non-existing pairs) were discarded.

**Candidates selection** The examples themselves were either taken from the list of gold (existing) pairs of related events and time entities or from pairs made of unrelated events and time entities.

Statistics computed on the training set were used to limit the number of pairs considered: with a crude character count, it was seen that the inner distances between entities involved in temporal relations were never over 213 characters. A threshold of a maximum distance of 300 characters based on that observation was used to limit the number of candidates considered for both the training and classification process. It effectively decreased the number of those candidates to 1/3rd of the possible pairs.

**Text pre-processing**    The text was made of a window centered around the $(event, time)$ entity pair extracted from the full clinical text. Maximum margins of 200 characters before the earliest entity and after the latest entity were taken to add context. As additional pre-processing, the extracted text was modified so the time entity was preceded by a '(TIME=) ' string, while the event entity was preceded by an '(EVENT=)' string. This processing was done in order to signify to the model which terms to look at to classify the provided text based on the candidate pair corresponding to that particular text (and assuming that other pairs might exist in the same text span).

### 4.2   Automated few-shot prompting with an auto-regressive LLM

**DSPy framework**    DSPy[2] (Khattab et al., 2023) is a framework developed by the Standford NLP group, which aims to optimize LLMs prompts algorithmically. This framework offers two main concepts: Signatures and Teleprompters. The signature is a declarative specification of the input/output behavior of a DSPy module, including a simple description of the task to be solved and descriptions of the input and output fields. Teleprompters are optimizers that can learn to bootstrap and automatically select effective prompts for the program modules. Compiling a DSPy program is based on a training set, a metric to maximize for validation, and a specific teleprompter. DSPy generates new, efficient prompts to match the changes made whenever a code, data, or metric is modified. DSPy also offers several optimizers and advanced features, but due to time constraints, we focused solely on using the BootstrapFewShotWithRandomSearch optimizer while developing our approach. This optimizer self-generates complete demonstrations several times

and performs a random search over these generated demonstrations to select the best program.

**Automated few-shot prompting**    As previously stated, we use the DSPy framework to develop and prompt our LLM-based approach. We first convert our input examples into all possible candidate pairs of $(event, time)$ using the gold annotations of entities and relations. For each combination, we extract the corresponding text from the document, which only contains the mentions of these entities. The corresponding text could be a small or a large portion of the full clinical text. Using DSPy, we defined a signature with instructions specifying the three possible types of relations and a description of the expected output format. Then, to cast the TRE task into a generation task, we evaluated these two configurations:

- **Predicting the relation triplet (event, relation, time).** By asking a question with a pair of $(event, time)$ and giving the corresponding text, we prompt our model to predict exactly an ordered list containing the event, the relation type, and the temporal expression. If no relation is found, the model should return an empty list. The basic idea behind this task design is to restrict the model to generate a specified format, avoiding extensive answers and hallucinations. Moreover, this output format is intended to prevent complex postprocessing strategies required to convert expected outputs into valid structures. This is the design we followed for the official submission.

- **Predicting only the relation type.** By asking a question with a pair of $(event, time)$ and giving the corresponding text, similarly to the previous system, we prompt our model to predict solely the relation type between the two entities in the pair. If no relation type is found, an empty list should be returned. This formulation mainly aims to simplify the task to the model. This configuration is evaluated after the official submissions of the shared task.

Figure 1 illustrates our used signatures (prompts) for both configurations. DSPy adds the reasoning statement in the ChainOfThoughts setting, which the LLM will generate to explain the task and the potential steps needed to generate the final answer. This reasoning step generally starts with a general statement, "Let's think step by step in order to *answer the question* or *produce the answer*", followed

by a tailored statement that the model will generate to answer the specific question in the demonstration. For instance, *"we need to find if the chemotherapy event 'carbo' and the date '8/23' have a specific relation. In the text, it is mentioned that ... This indicates that the chemotherapy event began in 8/23"*. Moreover, the automatically selected few-shot examples will be included as in-context demonstrations. As shown in Figure 1, to help the model produce the correct answer, we modified the CONTAINS relation type to CONTAINED-BY, particularly in the first configuration, in which the model must output an ordered list as an answer.

**Experimental settings**    We conducted our experiments using the `Mixtral-8X7B-Instruct-v0.1` language model from Mistral AI (Jiang et al., 2024). We generate up to 256 tokens and set the temperature generation parameter to 0. For both configurations of our automated few-shot prompting LLM approach, we use the BootstrapFewShotWithRandomSearch optimizer to select automatically $k$ few-shot examples. These few-shot examples are either chosen from given labeled training data or self-generated based on this data. Indeed, based on the examples in the labeled training data, the DSPy program uses the LLM to produce similar generated few-shot examples. As parameters, we generated 3 candidate programs, kept the maximum labeled examples to the default value, i.e., 16 examples, and set the maximum bootstrapped demos to 4. After converting the shared task training set into possible pairs of $(event, time)$ and the corresponding text, we subdivided this set into a training and a validation set (80/20). The validation set was mainly used to optimize the selection of few-shot examples from the training set[3].

### 4.3 Normalization and patient-level summarization

The triplets of relations *(event, relation_type, time)* obtained in earlier steps had their time mention processed -if necessary- to produce a normalized $TIMEX3$ expression in the form of a date. Two methods were used in order to do so: Heideltime and a simple LLM-based query with hand-made few shot examples. Both methods could take as input the time expression of the considered $(event, time)$ pair, but also -if

present for the document containing the pair- the $document\_creation\_time$ (in the form of a date).

**HeidelTime normalization**    As a first method to normalize the temporal expressions, we use a Python wrapper for the HeidelTime tool (Strötgen and Gertz, 2013), namely *py_heildetime*[4]. HeidelTime extracts and normalizes temporal expressions according to the TIMEX3 standard. The relative temporal expressions are normalized using the $document\_creation\_time$ (DCT). Since HeidelTime did not normalize relative temporal expressions such as *currently*, we normalize it to the DCT. This method was applied to both the outputs of the LLM-based TRE approach and the MLM fine-tuning approach.

**LLM-based query normalization**    This second method was used only for the MLM fine-tuning approach of the official submission. The latest state-of-the-art 7 billion parameters, `OpenChat 3.5` model (Wang et al., 2023a), was used through a local serving of an openai compatible API. The request itself was made of three parts.

**prompt part 1** was used everytime:

> *"please normalise the following string to a date format YYYY-MM-DD or, if you can't to a YYYY-MM format"*

**prompt part 2** was appended to $prompt1$ if a document time was available (with *<doc_time_input>*, a place holder to be replaced with the document date string:

> *"(the time at which the document is redacted is <doc_time_input>)"*

**prompt part 3** was used everytime, giving the time expression to normalize. It was appended to the previous part:

> *": <time expression>"*

From the former prompt and 6 short hand-picked synthetic examples in the form of triplets *(time_expression, doctime or None, answer_date or error_string)*, a few shot strategy was implemented as a user/assistant dialog.

**Summarization**    To provide a timeline from the triplets obtained earlier, summarization was performed as follows. First, we discarded the triplets containing a *time* mention not matching the Python regular expression:

---

[3]More details about the DSPy implementation code can be found in Khattab et al. (2023) and https://github.com/stanfordnlp/dspy

[4]https://github.com/hmosousa/py_heideltime

```
Respond to the question based on the given text.
The possible answers are: 'CONTAINED-BY',
'BEGINS-ON', 'ENDS-ON'.

---

Follow the following format.

Question: ${question}

Text: ${text}

Reasoning: Let's think step by step in order to
${produce the answer}. We …

Answer: a list containing only the relation. If no
relation is found, the answer is solely an empty list.

---

Question: Given this chemotherapy event: ${EVENT}
and this temporal expression: ${TIMEX}, which is
the relation between these entities, if any?

Text: ${text}
```

(a) Prompting the model to output the *relation type* between the given $(event, time)$ pair.

```
Respond to the question based on the given text.
The possible answers are: 'CONTAINED-BY',
'BEGINS-ON', 'ENDS-ON'.

---

Follow the following format

Question: ${question}

Text: ${text}

Reasoning: Let's think step by step in order to
${produce the answer}. We …

Answer: Each answer is an ordered list, containing
the chemotherapy event, then the corresponding
answer then the temporal expression. If no relation
is found, the answer is an empty list.

---

Question: Given this chemotherapy event: ${EVENT}
and this temporal expression: ${TIMEX}, which is
the relation between these entities, if any ?

Text: ${text}
```

(b) Prompting the model to output the relation triplet $(event, relation, time)$ given the $(event, time)$ pair.

Figure 1: The two defined DSPy signatures to prompt our automated few-shot prompting LLM approach.

```
'^([0-9]{4})-([0-9]{2})-([0-9]{2})$'
```

Then, following the organizers' instructions, for groups of triplets sharing the same date and event but with different relation types, i.e., *contains-1* and a more precise type (*begins-on*, *ends-on*), only the more precise mentions were kept. At last, triplets were de-duplicated and sorted.

## 5 Evaluation metrics

For the final evaluation of patient timelines, the organizers provide an evaluation code[5]. The evaluation process covers strict and relaxed evaluation settings by calculating the average F1 score across all patients. The official score is an arithmetic mean of two types of Macro F1 measure, type A and type B, in a relaxed to-month setting. The type A evaluation includes the patients with no gold timelines, while the type B evaluation excludes the patients with no true relations. The relaxed to-month setting means only the month must match the gold annotation. More details about the evaluation process are presented in the shared task website[6]. While

selecting the different models we tried, we evaluated them based on the official score provided by the organizers.

For the optimization of our automated few-shot LLM approach and to ensure quality few-shot examples and demos, we defined a strict F1 measure. Indeed, the DSPy optimizer will only keep the few-shot examples that maximize this evaluation metric. Note that for our first configuration setting, i.e., predicting an ordered list of $(event, relation, time)$, the system prediction will not be considered a match if the model correctly predicts the relation type but fails to output the required format.

## 6 Results & discussion

To participate in this shared task, we submitted two runs. The first run is the MLM fine-tuning approach (**NLPeers 1**), and the second run is the automated few-shot prompting LLM approach (**NLPeers 2**). In this section, we begin by discussing the overall performance of our systems on the test and the development sets. Since the gold annotations of the test set will not be released, we then present a more in-depth review of each of our methods on this set and the impact of adding the LLM-based

normalization on performance.

## 6.1 Overall performance

Table 1 presents the official subtask 1 results at the patient-level of our submitted methods and the organizers' baseline system[7] on the test set, including average scores and scores per cancer type, as reported in the Leader board of the Chemotimelines shared task. The MLM fine-tuning approach outperforms the automated few-shot LLM-based approach on the test set, with an average score of 0.77 vs. 0.64. However, the best results are obtained with the baseline system, with an average score of 0.89.

In comparison with submissions from other participants for this subtask, we are the third-best team, out of eight teams, in terms of average score if we consider our submission of MLM fine-tuning approach (**NLPeers 1**). It's worth noting that only the top team outperformed the baseline system. On the Melanoma dataset, we are in second place with a score of 0.84 vs. 0.87 for both the top team and the baseline system).

Table 2 summarizes the results of our proposed approaches on the development set, including the additional experiments of LLM-based query normalization and predicting the relation type for the few-shot prompting LLM approach, which were not part of our official submissions. Similar to the results on the test set, the fine-tuned MLM approach (an average score of 0.85) outperforms the automated few-shot LLM approach in both configurations using the HeidelTime normalization (an average score of 0.61 for the relation type prediction and 0.56 for the relation triplet prediction). However, using both HeidelTime and LLM-based query normalization enhanced the results of the relation triplet prediction, hinting at the fact that performances measured on the test set could probably have been higher if the combined normalization was applied to the automated few-shot LLM approach. Interestingly, the official submission models have performances that vary in opposite directions when going from the development set to the test set: the fine-tuned MLM model performance decreases while the few-shot one increases.

## 6.2 Performance of fine-tuning MLM model

**Whole set relation type errors**   Table 3 represents a confusion matrix computed on the develop-

ment set after applying the fine-tuned MLM model, it compares the gold relation types to the predictions made. As can be seen in this table, a major source of error on the development set for this method is the mislabeling of 'false' (*no_link*) candidate triplets as *contains* triplets. It accounts for roughly 10% of the *no_link* candidates. This is not unexpected since *no_link* candidates are by far the first class present in the used development set (655 total, after filtering based on entity distance), followed by 'contains' triplets which represent roughly half of the former ones (354 total).

Next error based on absolute count are *ends-on* relations mislabeled as *begins-on* (38/83), while the converse almost never occurs (2/103), although the categories *begins-on* and *ends-on* are almost balanced (respectively 103 and 83 occurrences).

**Melanoma subset relation type errors**   As out of the of the three cancer subsets, the model seemed to perform relatively worse on the melanoma, we inspected further the errors specifically made for the melanoma subset, it appears that it responsible for the vast majority of the *no_link* candidates mislabeled as *contains* triplets made in the general set (i.e. 64 out of the 68 counted in Table 3). Interestingly, the melanoma subset also accounts for 41 out of the 49 *begins-on* relations mislabeled as *contains*. This relative concentration of errors in the melanoma subset could be explained by the lower count of melanoma examples in the training set, increasing the odds that the model learned undue associations specific to that subset.

## 6.3 Performance of automated few-shot LLM prompting

As reported in Table 2, using the HeidelTime normalization, predicting relation type yields better results than predicting relation triplets, with an average F-measure of 0.61 vs. 0.56. This could be due to the strict evaluation of the triplet configuration. Indeed, as already mentioned, no extra post-processing steps are taken for the outputs. Results per cancer are jointly discussed, along with the impact of normalization methods, in the next section.

Among all the possible candidate pairs (1287), the **relation triplet** model predicts 1046 tuples and 241 empty lists. Among the 1046 tuples, 133 are invalid, i.e., not corresponding to an ordered list $(event, relation, time)$ or not mentioning the correct $event$ or $time$ present in the input. Among

| Approach | Average Score | Breast cancer | Melanoma | Ovarian |
|---|---|---|---|---|
| Fine-tuned MLM<br>+ HeidelTime & OC normalization<br>*(NLPeers 1)* | 0.77 | 0.72 | 0.84 | 0.75 |
| Automated few-shot LLM<br>**(Relation triplet)**<br>+ HeidelTime normalization<br>*(NLPeers 2)* | 0.64 | 0.49 | 0.81 | 0.63 |
| Baseline system | 0.89 | 0.93 | 0.87 | 0.88 |

Table 1: The official results on the test set. OC refers to the LLM-based normalization using the `OpenChat` model.

| Approach | Average Score | Breast | Melanoma | Ovarian |
|---|---|---|---|---|
| Fine-tuned MLM | | | | |
| **Relation type (classification)**<br>+ HeidelTime & OC normalization<br>*(official submission, NLPeers 1)* | 0.85 | 0.84 | 0.81 | 0.88 |
| **Relation type (classification)**<br>+ HeidelTime normalization<br>*(non official submission)* | 0.74 | 0.61 | 0.85 | 0.76 |
| Automated few-shot LLM | | | | |
| **Relation triplet (generation)**<br>+ HeidelTime & OC normalization<br>*(non official submission)* | 0.72 | 0.70 | 0.74 | 0.71 |
| **Relation type (generation)**<br>+ HeidelTime normalization<br>*(non official submission)* | 0.61 | 0.57 | 0.78 | 0.48 |
| **Relation triplet (generation)**<br>+ HeidelTime normalization<br>*(official submission, NLPeers 2)* | 0.56 | 0.53 | 0.70 | 0.47 |

Table 2: The results on the development set. OC refers to the LLM-based normalization using the `OpenChat` model.

the remaining 913 valid tuples, 146 are correct (69 *begins-on*, 34 *ends-on*, 43 *contains*). As for the **relation type** model, among all the possible candidate pairs (1287), it predicts 994 relation types and 293 empty relations. Among the 994 predicted relation types, 824 respect the expected output format, and 154 are correct (90 begins-on, 47 ends-on, 19 contains).

Table 4 presents the number of semantic errors, as defined in Li et al. (2023b), as well as some semantically incorrect samples on the development set for both configurations of automated few-shot LLM approach, using the HeidelTime normalization. A semantic error is defined as a relation type that does not exist in the pre-defined set of relation types. Looking at this table, we notice that although the relation triplet prediction model produces a total of 43 errors, only 7 types of errors are generated and seem semantically "correct" but are out of the pre-defined relation type set. However, the relation type prediction model produces only a total of 16 errors, including 11 different types of relation types, which seems less precise. Indeed, the relation type prediction model tends to generate large texts containing not only the relation but also explanations and hallucinations. Though the main idea behind relation type prediction is to simplify the relation extraction task to the LLM, we believe that reformulating the task with structured instructions and input/output examples, such as our triplet prediction method, could provide better results, using the appropriate pre- and post-processing steps, as already stated in previous research works (Li et al., 2023b; Lu et al., 2022; Wang et al., 2023b).

| Gold | BEGINS-ON | CONTAINS | ENDS-ON | no_link |
|---|---|---|---|---|
| BEGINS-ON | 52 | 18 | 30 | 0 |
| CONTAINS | 49 | 328 | 30 | 68 |
| ENDS-ON | 2 | 1 | 11 | 0 |
| no_link | 0 | 7 | 12 | 587 |
| Total | 103 | 354 | 83 | 655 |

Table 3: Confusion matrix for the MLM fine-tuning approach applied on the development set.

| | Relation Type Error | Semantically incorrect samples |
|---|---|---|
| Automated few-shot LLM **Relation triplet** + Heideltime *(official submission, NLPeers 2)* | 43 | occurs on, occurs-on, contained-in, not going to occur, not related, duration, ended-on |
| Automated few-shot LLM **Relation type** + Heideltime *(non official submission)* | 16 | answer, be, beg, begins, conta, during, every-on, happening-on, happens-on, lasts-for, planned-for |

Table 4: Semantic errors and semantically incorrect samples on the development set.

### 6.4 Impact of LLM normalization on performance

It should be noted here that although we used both Heideltime and an LLM-based normalization for the official MLM fine-tuning results, due to time constraints, only the Heideltime normalization was made available for the official automated few-shot prompting results. A comparison of results with and without said LLM normalization was done after the official results on the development set. The results in Table 2 show that the additional Open Chat normalization has a strong impact on both predictors, with an increase ranging from 11 (fine-tuned MLM) to 16 points (automated few-shot prompting) on average score. This seems to suggest that such a simple prompt method can be efficient for this kind of task, where a very limited context and no specific background knowledge is needed to answer the query at hand, thus requiring no complex task description or prompt search strategy.

A more detailed look at the impact of the complementary normalization per cancer type seems to indicate that breast and ovarian subsets benefit the most regardless of the model. A detailed inspection of the differences in time expression patterns highlights that the temporal expressions of melanoma are less varied, with different pattern proportions. For example, in the few-shot LLM triplet prediction, *currently* - which is well normalized by Heideltime when the document time is given - accounts for 20% of temporal expressions of melanoma, but only 5% in other cancer. In the same way, *today* accounts for 35% of melanoma time expressions and 28% of other cancers time expressions.

While we tested two normalization approaches, the reference one provided as a scala library by the contest organizers was not tested as we failed to include it in time in our otherwise Python-based code. The effect on the measured performances (and comparison to other teams' proposals using it) is difficult to assess as - besides the respective merits of each method - the reference time normalization was used as a gold standard for the evaluation process. In effect, terms that it could not normalize were discarded, transforming potentially correct time expressions and relations to perceived incorrect ones.

### 7 Conclusion

In this work, we showed that an LLM-based automated prompting method could, with no weight fine-tuning, give good results on a temporal relation extraction task. We also showed that a smaller fine-tuned MLM likely performs better while requiring less computing resources, thus confirming that smaller model fine-tuning is still relevant for such classification tasks. Given the low number of examples retained at the end of the selection procedure by the few-shot prompting approach, it can

be inferred that a smaller set of examples could be used to reach better performances, effectively making it an interesting choice when access to annotated data is scarce. Another finding demonstrated the effective use of a simple LLM approach for a general domain task such as time normalization.

## Limitations

It is to be noted that the time devoted to developing both proposed methods was limited due to late enrollment in the shared task and access to data. The methods were also mostly developed from scratch w.r.t. the timeline prediction perspective. This strongly suggests that both approaches could be improved. As such, more work is warranted to get these proposed solutions closer to being the best-performing ones. For the automated few-shot prompting LLM solution, creating a true pipeline chaining multiple steps (e.g., verification/enrichment) could greatly increase the accuracy of the provided answers. Indeed, more evaluation steps should be included, in particular for the clinical domain, to avoid inaccuracies in the generated reasoning steps and demonstrations. Another improvement would be to use rule-based post-processing steps to deal with the inherent variability of answers produced by the LLM. Further research into using DSPy, particularly its advanced prompting and optimization features, could also be conducted. For the fined-tuned MLM approach, proper parameter selection could increase the performance and stability of the model. On a last note on the two proposed methods, we considered them as exclusive to one another to measure their respective benefits, but a combination of both could allow the final result to get even better performances. Finally, we did not compare our normalization process to the one provided as a gold standard, making it more difficult to draw definitive conclusions based on the final evaluation of our proposal and its comparison to other participants performances.

## Ethics statement

Using Large Language Models (LLMs) in the clinical domain raises several ethical concerns. First, due to the sensitive nature of clinical data, special precautions must be taken while working with it. This work uses de-identified clinical data obtained through a Data Use Agreement. Therefore, the designed prompts for our LLM-based methods do not contain identifying personal information about patients. Second, a major challenge while leveraging LLMs, particularly in clinical research, is the transparency and interoperability of results. Indeed, these models often act as 'black boxes,' making it hard to understand the generated outputs and the decisions made, which is crucial for clinicians. As a result, a human and expert evaluation is required, first for minimizing hallucinations, biases, and harmfulness outputs and then for evaluating and validating the coherence of generation. Third, LLMs are complex models with billions of parameters that necessitate lots of computational resources, thus generating a carbon footprint. This is also valid for fine-tuning the MLMs-based models. Finally, it is worth noting that the proposed methods are mainly for research purposes, and additional studies need to be conducted before integrating them into practical applications, where the goal is to help clinicians conduct a systematic analysis of large patient records.

## Acknowledgments

## References

Qingyu Chen, Jingcheng Du, Yan Hu, Vipina Kuttichi Keloth, Xueqing Peng, Kalpana Raja, Rui Zhang, Zhiyong Lu, and Hua Xu. 2023. Large language models in biomedical natural language processing: benchmarks, baselines, and recommendations. *arXiv preprint arXiv:2305.16326*.

Veera Raghavendra Chikka. 2016. CDE-IIITH at SemEval-2016 task 12: Extraction of temporal information from clinical documents using machine learning techniques. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1237–1240, San Diego, California. Association for Computational Linguistics.

Markus Eberts and Adrian Ulges. 2019. Span-based joint entity and relation extraction with transformer pre-training. *arXiv preprint arXiv:1909.07755*.

R. Gaizauskas, H. Harkema, M. Hepple, and A. Setzer. 2006. Task-oriented extraction of temporal information: The case of clinical narratives. In *Thirteenth International Symposium on Temporal Representation and Reasoning (TIME'06)*, pages 188–195.

Ridong Han, Tao Peng, Chaohao Yang, Benyou Wang, Lu Liu, and Xiang Wan. 2023. Is information extraction solved by chatgpt? an analysis of performance, evaluation criteria, robustness and errors. *arXiv preprint arXiv:2305.14450.*

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *Preprint*, arXiv:2111.09543.

Eddie Paul Hernández, Alexandra Pomares Quimbaya, and Oscar Mauricio Muñoz. 2016. Htl model: A model for extracting and visualizing medical events from narrative text in electronic health records. In *ICT4AgeingWell*, pages 107–114.

Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. REBEL: Relation extraction by end-to-end language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370–2381, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088.*

Zhengbao Jiang, Frank F. 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.

Abdulrahman Khalifa, Sumithra Velupillai, and Stephane Meystre. 2016. UtahBMI at SemEval-2016 task 12: Extracting temporal information from clinical text. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1256–1262, San Diego, California. Association for Computational Linguistics.

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714.*

Hee-Jin Lee, Hua Xu, Jingqi Wang, Yaoyun Zhang, Sungrim Moon, Jun Xu, and Yonghui Wu. 2016. UTHealth at SemEval-2016 task 12: an end-to-end system for temporal information extraction from clinical notes. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1292–1297, San Diego, California. 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.

Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei Ye, Wen Zhao, and Shikun Zhang. 2023a. Evaluating chatgpt's information extraction capabilities: An assessment of performance, explainability, calibration, and faithfulness. *arXiv preprint arXiv:2304.11633.*

Peng Li and Heng Huang. 2016. UTA DLNLP at SemEval-2016 task 12: Deep learning based natural language processing system for clinical information identification from clinical notes and pathology reports. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1268–1273, San Diego, California. Association for Computational Linguistics.

Peng Li, Tianxiang Sun, Qiong Tang, Hang Yan, Yuanbin Wu, Xuanjing Huang, and Xipeng Qiu. 2023b. CodeIE: Large code generation models are better few-shot information extractors. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15339–15353, Toronto, Canada. Association for Computational Linguistics.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190.*

Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2021. EntityBERT: Entity-centric masking strategy for model pretraining for the clinical domain. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 191–201, Online. 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*.

Yaojie Lu, Qing Liu, Dai Dai, Xinyan Xiao, Hongyu Lin, Xianpei Han, Le Sun, and Hua Wu. 2022. Unified structure generation for universal information extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5755–5772, Dublin, Ireland. Association for Computational Linguistics.

Azadeh Nikfarjam, Ehsan Emadzadeh, and Graciela Gonzalez. 2013. Towards generating a patient's timeline: Extracting temporal relationships from clinical notes. *Journal of Biomedical Informatics*, 46:S40–S47. Supplement: 2012 i2b2 NLP Challenge on Temporal Relations in Clinical Data.

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, Online. Association for Computational Linguistics.

Jannik Strötgen and Michael Gertz. 2013. Multilingual and cross-domain temporal tagging. *Language Resources and Evaluation*, 47:269–298.

Buzhou Tang, Yonghui Wu, Min Jiang, Yukun Chen, Joshua C Denny, and Hua Xu. 2013. A hybrid system for temporal information extraction from clinical text. *Journal of the American Medical Informatics Association*, 20(5):828–835.

Julien Tourille, Olivier Ferret, Aurélie Névéol, and Xavier Tannier. 2017a. Neural architecture for temporal relation extraction: A Bi-LSTM approach for detecting narrative containers. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 224–230, Vancouver, Canada. Association for Computational Linguistics.

Julien Tourille, Olivier Ferret, Xavier Tannier, and Aurélie Névéol. 2017b. Temporal information extraction from clinical text. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 739–745, Valencia, Spain. 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, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2023a. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*.

Wei Wang, Kory Kreimeyer, Emily Jane Woo, Robert Ball, Matthew Foster, Abhishek Pandey, John Scott, and Taxiarchis Botsis. 2016. A new algorithmic approach for the extraction of temporal associations from clinical narratives with an application to medical product safety surveillance reports. *Journal of Biomedical Informatics*, 62:78–89.

Xingyao Wang, Sha Li, and Heng Ji. 2023b. Code4Struct: Code generation for few-shot event structure prediction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3640–3663, Toronto, Canada. Association for Computational Linguistics.

Jiarui Yao, Harry Hochheiser, WonJin Yoon, Eli Goldner, and Guergana Savova. 2024. Shared task: Chemotherapy treatment timeline extraction. *Clinical NLP Workshop, NAACL 2024. Mexico City, Mexico.*

Li Zhou, Simon Parsons, and George Hripcsak. 2008. The Evaluation of a Temporal Reasoning System in Processing Clinical Discharge Summaries. *Journal of the American Medical Informatics Association*, 15(1):99–106.

# KCLab at Chemotimelines 2024:
# End-to-end system for chemotherapy timeline extraction – Subtask2

**Yukun Tan, Merve Dede, Ken Chen**

Department of Bioinformatics and Computational Biology,
The University of Texas MD Anderson Cancer Center, Houston, Tx, USA

{ytan1, mdede, kchen3}. mdanderson.org

## Abstract

This paper presents our participation in the Chemotimelines 2024 subtask2, focusing on the development of an end-to-end system for chemotherapy timeline extraction. We initially adopt a basic framework from subtask2, utilizing Apache cTAKES for entity recognition and a BERT-based model for classifying the temporal relationship between chemotherapy events and associated times. Subsequently, we enhance this pipeline through two key directions: first, by expanding the exploration of the system, achieved by extending the search dictionary of cTAKES with the UMLS database; second, by reducing false positives through preprocessing of clinical notes and implementing filters to reduce the potential errors from the BERT-based model. To validate the effectiveness of our framework, we conduct extensive experiments using clinical notes from breast, ovarian, and melanoma cancer cases. Our results demonstrate improvements over the previous approach.

## 1 Introduction

In recent years, the rapid development and widespread implementation of Electronic Health Records (EHRs) have created a significant demand for the clinical notes processing and information extraction within the realm of medical research. (Yanshan Wang et al., 2018) Particularly, the extraction of temporal information, encompassing temporal expressions, temporal events, and temporal relations, has created new opportunities for dynamic treatment studies(Sun et al., 2013; UzZaman et al., 2014). Among the various treatment modalities, chemotherapy stands out as one of the most critical and widely used approaches in cancer therapy. EHRs with temporal information offer a unique advantage by providing a chronological roadmap of patient-specific treatments. These timelines play a key role in understanding the effectiveness of chemotherapy, evaluating treatment responses, and identifying patterns in patient outcomes. They serve as invaluable resources for researchers investigating the interplay among treatment protocols, tumor biology, and patient characteristics. By analyzing these timelines, researchers can uncover trends, predictors of response and potential markers for treatment success or failure. As such, the construction of accurate and comprehensive chemotherapy treatment timelines is not only an academic pursuit but a practical clinical necessity in advancing cancer care and improving outcomes.

However, this task presents notable challenges due to the domain-specific nature of EHRs, namely, variations in writing style and quality, lack of text structure, and the pervasive presence of redundant information. Moreover, the creation of annotated corpora manually is a resource-intensive process, demanding substantial human effort and time. Consequently, numerous research efforts have turned to employ rule-based, machine-learning, or hybrid methods to extract general temporal information from clinical narratives (Moharasan & Ho, 2019; Najafabadipour et al., 2020; Liwei Wang et al., 2020). Notably, despite these efforts, there are currently no available tools designed specifically for extracting timelines to contextualize cancer treatment. Hence, this competition subtask aims to fill this gap by developing an end-to-end system for chemotherapy timeline extraction (*Jiarui Yao et al., 2024). This system not only addresses the urgent need for accurate and

Figure 1: System Overview - Baseline framework enhanced with clinical notes preprocessing, directional time mention filtering, and UMLS integration to extend the extraction dictionary.

comprehensive timelines but also showcases the potential of leveraging advanced computational methods to enhance cancer care practices.

## 2 System Overview

### 2.1 Baseline framework

The pipeline mainly combines three main software packages: Apache cTAKES (Savova et al., 2010), CLU Lab Timenorm (Laparra et al., 2018; Xu et al., 2019), and Huggingface transformers. cTAKES, a java-based tool, offers powerful text engineering and information extraction capabilities, particularly tailored for clinical text. It utilizes the cTAKES Python Bridge to Java (ctakes-pbj) to process text artifacts seamlessly in Python, leveraging cTAKES' modules for entity recognition and sentence tokenization. CLU Lab Timenorm is employed for identifying and normalizing date and time expressions. The pipeline has incorporated a customized version of Timenorm into the pipeline, which allows for improved handling of approximate dates, a common occurrence in clinical narratives. This step ensures consistency and standardization of temporal representations. Huggingface Transformers is a widely used deep learning library for natural language processing tasks. This pipeline employs the PubMedBERT – based model (Gu et al., 2022) (Temporal Link – TLINK model) to identify and classify the temporal relationships between chemotherapy mentions and their associated dates. The classifier determines the temporal relationship between each paired mention, whether it's "begin on", "end on", "contain-1", or "none".

### 2.2 UMLS integration

We have enhanced the capabilities of cTAKES by integrating the Unified Medical Language System (UMLS), thereby extending its dictionary to have a more comprehensive range of chemotherapy terminologies. We evaluated all medically relevant concepts in the UMLS database related to chemotherapy as well as their descendant terms to obtain a complete hierarchy. This integration allows cTAKES to recognize and extract a broader array of chemotherapy-related terms, including generic drug names, their synonyms, treatment protocols, drug brand names and so on. This approach ensures that the system captures a more exhaustive list of chemotherapy-related terms and agents, thus improving the completeness of extracted information.

### 2.3 Clinical notes preprocessing

In our clinical notes preprocessing stage, we implemented several steps to enhance the efficiency and accuracy of information extraction. We examined the provided notes carefully to evaluate their structures, and to understand the content of information in each note category. After the evaluation period, firstly, we removed files with names ending in "RAD" or "SP", as these notes often did not contain any chemotherapy information or only contained redundant chemotherapy history of the patients, which were already present in other clinical notes. For example, the files with "RAD" may contain the information related to radiation procedures or outcomes, which occasionally included descriptions of chemotherapy in the patient history statement. We determined these sections to be redundant, as more detailed and clearer descriptions were typically already found in files ending with "NOTE" or "PGN". Secondly, we eliminated the concluding portions of files containing information about the person recording the note, time, and location. While these timestamps may initially be perceived as valuable, they are usually redundant as timestamps were typically provided at the beginning of each record. Additionally, these

sections often contained abbreviations that overlapped with abbreviations in our expanded UMLS dictionary for chemotherapy agents, leading to false positives. Thirdly, we employed fuzzy recognition to filter out paragraphs related to treatment plans. Since current treatment plans are incomplete until documented as finished in subsequent notes, we can confidently exclude these sections without missing relevant information. This step effectively reduced the occurrence of false positives, as changes to treatment plans were noted elsewhere in subsequent records. These preprocessing steps not only simplified the data but also significantly enhanced the precision of our information extraction process, ensuring that extracted chemotherapy-related details are accurate and comprehensive.

## 2.4 Directional time mention filtering

After successfully extracting chemotherapy events and temporal expression pairs, we introduced a novel filter prior to the temporal relation classification step, focusing on the directionality of time mentions. This filter aims to reduce potential errors in classification by considering the ordering of temporal expressions in relation to the chemo events. Specifically, when multiple temporal expressions surround a chemo event within the same sentence and appear after the chemo event, we prioritize these temporal expressions over those occurring before the chemo event. For instance, in the sentence "He had resection in Jun 2008, last chemo was in Nov 2010," we identified the temporal expressions "Jun 2008" and "Nov 2010." In this case, we disregard the time preceding the chemotherapy event since a temporal expression already exists in the same sentence following the chemotherapy event, making it clear that "Nov 2010" pertains to the chemotherapy event. Likewise, in the sentence "She presents to the office on today's date for the chemo as per the standard FDA approved regimen. She also did radiation last week," we detected the temporal expressions "today's date" and "last week" surrounding the chemotherapy event. However, since "last week" is not in the same sentence as the chemotherapy event, we do not ignore the temporal expression "today's date." This analysis of directional cues in time mentions is crucial for our task.

While theoretically, the BERT – based model's classification (TLINK) could address this by categorizing irrelevant times as "none", our findings suggest that due to potential limitations in training data, this classification may not always be accurate, particularly in scenarios where temporal expressions occur both before and after the chemotherapy terminology. Our introduced filter significantly reduces the chances of misclassifications, thereby enhancing the accuracy and robustness of our temporal relation classification system.

## 3 Results

In Chemotimelines 2024 subtask2, our team achieved the 3rd highest rank in the average scores, with F1 scores of 0.68 for breast cancer (rank #1), 0.49 for melanoma (rank #3), and 0.45 for ovarian cancer (rank #7) (Table 3). These scores were calculated by averaging type A and type B evaluation metrics. Type A includes notes without true relations, while type B excludes such notes. Comparing our results to baseline performance, we observed an improvement of around 5%-10% for breast cancer and melanoma, while not for ovarian cancer.

Due to the unavailability of the test set, we present results from the development set to analyze the strengths and weaknesses of our pipeline. As depicted in Table 1 (Type A) and Table 2 (Type B), our system generally outperforms the baseline in terms of recall, attributed to the integration of the UMLS dictionary. However, this integration also introduces certain challenges, such as generating false positives. These false positives included synonymous terms like "vegf trap" and "aflibercept," terminologies do not present in the gold timelines such as "aldesleukin," and duplicated abbreviations with different meanings.

Our implemented preprocessing procedure and filtering step effectively reduced false positives not only from the integrated dictionary but also from cases prone to misclassification by the TLINK model. However, this also led to the exclusion of some true pairs from the gold timelines. For example, some patients only had "RAD" files, which do not pertain to chemotherapy details,

Table 1: Type A evaluation of dev set

|  |  | Prec | Recall | F1 |
|---|---|---|---|---|
| Baseline | Breast | 0.874 | 0.894 | 0.880 |
|  | Ovarian | 0.648 | 0.884 | 0.716 |
|  | Melanoma | 0.569 | 0.560 | 0.565 |
| Proposed | Breast | 0.926 | 0.897 | 0.909 |
|  | Ovarian | 0.681 | 0.851 | 0.736 |
|  | Melanoma | 0.570 | 0.627 | 0.595 |

Table 2: Type B evaluation of dev set

|  |  | Prec | Recall | F1 |
|---|---|---|---|---|
| Baseline | Breast | 0.831 | 0.885 | 0.848 |
|  | Ovarian | 0.648 | 0.884 | 0.716 |
|  | Melanoma | 0.354 | 0.340 | 0.347 |
| Proposed | Breast | 0.801 | 0.725 | 0.757 |
|  | Ovarian | 0.681 | 0.851 | 0.736 |
|  | Melanoma | 0.355 | 0.440 | 0.393 |

resulting in missed records that impact our evaluation significantly. Additionally, some chemotherapy pairs were solely mentioned in the plan section, such as "we will give chemo cycle 2 today." While we expected subsequent confirmed notes, they were not present, resulting in the omission of such pairs from our analysis.

Reviewing the test results (Table 3), we obtained the most favorable outcomes for breast cancer, which is the group with the largest sample size. Conversely, the small size of the ovarian cancer type test set poses challenges, as even slight variations in missed or additional pairs can lead to substantial variance. Furthermore, we observed that the gold timeline may not always be entirely accurate, potentially resulting in the omission of rare chemotherapy terms. Addressing these

challenges necessitates a larger and more diverse patient dataset in future evaluations.

## 4 Conclusion and future work

This paper details our efforts in the Chemotimelines 2024 subtask2, focusing on the development of an end-to-end system for chemotherapy timeline extraction. Our experiments utilizing clinical notes from breast, ovarian cancer, and melanoma cases have demonstrated the enhancements made to our pipeline. These enhancements include expanding the system's capabilities by leveraging the UMLS database and implementing preprocessing and directional filtering procedures to effectively reduce false positives.

Future works could potentially include firstly creating a more detailed and precise dictionary using the UMLS, with specific terms tailored to different cancer types, and establishing a synonymous dictionary to prevent duplication of terms. Secondly, it is crucial to exercise caution when removing files such as "RAD" and "SP," especially in cases where patients only possess these notes. Finally, exploring the use of ChatGPT and appropriate prompts as an alternative to the TLINK classifier, which is currently fine-tuned from PubMedBERT, would be a valuable exercise. ChatGPT's superior understanding of sentence context could prove beneficial for those classifications that do not require specific domain knowledge.

Table 3: Final evaluation of test set

| Average Scores | | Breast Cancer | | Melanoma | | Ovarian | |
|---|---|---|---|---|---|---|---|
| Team | Score | Team | Score | Team | Score | Team | Score |
| LAILab 2 | 0.70 | KCLab 1 | 0.68 | LAILab 2 | 0.74 | LAILab 2 | 0.74 |
| LAILab 1 | 0.56 | Wonder 2 | 0.64 | LAILab 1 | 0.57 | LAILab 1 | 0.59 |
| KCLab 1 | 0.54 | Wonder 1 | 0.63 | KCLab 1 | 0.49 | Wonder 3 | 0.55 |
| Wonder 3 | 0.53 | Wonder 3 | 0.63 | Wonder 3 | 0.39 | Wonder 2 | 0.55 |
| Wonder 2 | 0.52 | LAILab 2 | 0.62 | Wonder 1 | 0.39 | Wonder 1 | 0.53 |
| Wonder 1 | 0.52 | LAILab 3 | 0.53 | Wonder 2 | 0.39 | LAILab 3 | 0.49 |
| LAILab 3 | 0.47 | LAILab 1 | 0.52 | LAILab 3 | 0.38 | KCLAb 1 | 0.45 |
| NYULangone | 0.23 | UTSA-NLP 1 | 0.25 | NYULangone | 0.32 | UTSA-NLP 1 | 0.19 |
| UTSA-NLP 1 | 0.22 | NYULangone | 0.19 | UTSA-NLP 1 | 0.21 | NYULangone | 0.18 |
|  |  |  |  |  |  |  |  |
| Baseline | 0.58 | Baseline | 0.59 | Baseline | 0.43 | Baseline | 0.71 |

## Acknowledgments

## References

Gu, Yu, Tinn, Robert, Cheng, Hao, Lucas, Michael, Usuyama, Naoto, Liu, Xiaodong, Naumann, Tristan, Gao, Jianfeng, & Poon, Hoifung. (2022). Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing. *ACM Transactions on Computing for Healthcare*, *3*(1), 1–23. https://doi.org/10.1145/3458754

*Jiarui Yao, *Harry Hochheiser, WonJin Yoon, Eli Goldner, & Guergana Savova. (2024). Overview of the 2024 Shared Task on Chemotherapy Treatment Timeline Extraction. *Proceedings of the 6th Clinical Natural Language Processing Workshop*.

Laparra, Egoitz, Xu, Dongfang, & Bethard, Steven. (2018). From Characters to Time Intervals: New Paradigms for Evaluation and Neural Parsing of Time Normalizations. *Transactions of the Association for Computational Linguistics*, *6*, 343–356. https://doi.org/10.1162/tacl_a_00025

Moharasan, Gandhimathi, & Ho, Tu-Bao. (2019). Extraction of Temporal Information from Clinical Narratives. *Journal of Healthcare Informatics Research*, *3*(2), 220–244. https://doi.org/10.1007/s41666-019-00049-0

Najafabadipour, Marjan, Zanin, Massimiliano, Rodríguez-González, Alejandro, Torrente, Maria, Nuñez García, Beatriz, Cruz Bermudez, Juan Luis, Provencio, Mariano, & Menasalvas, Ernestina. (2020). Reconstructing the patient's natural history from electronic health records. *Artificial Intelligence in Medicine*, *105*, 101860. https://doi.org/10.1016/j.artmed.2020.101860

Savova, Guergana K., Masanz, James J., Ogren, Philip V., Zheng, Jiaping, Sohn, Sunghwan, Kipper-Schuler, Karin C., & Chute, Christopher G. (2010). Mayo clinical Text Analysis and Knowledge Extraction System (cTAKES): Architecture, component evaluation and applications. *Journal of the American Medical Informatics Association : JAMIA*, *17*(5), 507–513. https://doi.org/10.1136/jamia.2009.001560

Sun, Weiyi, Rumshisky, Anna, & Uzuner, Ozlem. (2013). Evaluating temporal relations in clinical text: 2012 i2b2 Challenge. *Journal of the American Medical Informatics Association : JAMIA*, *20*(5), 806–813. https://doi.org/10.1136/amiajnl-2013-001628

UzZaman, Naushad, Llorens, Hector, Allen, James, Derczynski, Leon, Verhagen, Marc, &

Pustejovsky, James. (2014). *TempEval-3: Evaluating Events, Time Expressions, and Temporal Relations* (arXiv:1206.5333). arXiv. https://doi.org/10.48550/arXiv.1206.5333

Wang, Liwei, Wampfler, Jason, Dispenzieri, Angela, Xu, Hua, Yang, Ping, & Liu, Hongfang. (2020). Achievability to Extract Specific Date Information for Cancer Research. *AMIA Annual Symposium Proceedings*, *2019*, 893–902. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7153063/

Wang, Yanshan, Wang, Liwei, Rastegar-Mojarad, Majid, Moon, Sungrim, Shen, Feichen, Afzal, Naveed, Liu, Sijia, Zeng, Yuqun, Mehrabi, Saeed, Sohn, Sunghwan, & Liu, Hongfang. (2018). Clinical Information Extraction Applications: A Literature Review. *Journal of Biomedical Informatics*, *77*, 34–49. https://doi.org/10.1016/j.jbi.2017.11.011

Xu, Dongfang, Laparra, Egoitz, & Bethard, Steven. (2019). Pre-trained Contextualized Character Embeddings Lead to Major Improvements in Time Normalization: A Detailed Analysis. In Rada Mihalcea, Ekaterina Shutova, Lun-Wei Ku, Kilian Evang, & Soujanya Poria (Eds.), *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)* (pp. 68–74). Association for Computational Linguistics. https://doi.org/10.18653/v1/S19-1008

# Project PRIMUS at EHRSQL 2024: Text-to-SQL Generation using Large Language Model for EHR Analysis

**Sourav Bhowmik Joy**     **Rohan Ahmed**     **Argha Pratim Saha**
**Minhaj Ahmed**     **Utsho Das**     **Partha Sarothi Bhowmik**
CSE, Shahjalal University of Science & Technology

## Abstract

This paper explores the application of the sqlcoders model, a pre-trained neural network, for automatic SQL query generation from natural language questions. We focus on the model's internal functionality and demonstrate its effectiveness on a domain-specific validation dataset provided by EHRSQL. The sqlcoders model, based on transformers with attention mechanisms, has been trained on paired examples of natural language questions and corresponding SQL queries. It takes advantage of a carefully crafted prompt that incorporates the database schema alongside the question to guide the model towards the desired output format.

## 1 Introduction

Electronic health records (EHRs), a large collection of data related to digital medical records, serve as the backbone of modern healthcare, storing a wealth of patient information. This data, encompassing diagnoses, procedures, medications, and more, offers invaluable insights for clinical decision-making and research.[3]  [4] However, effectively utilizing this vast resource is often hampered by the complexity of querying the underlying relational databases.

Traditionally, hospital staff relies on pre-defined rule conversion systems to interact with EHR databases. These systems, while functional, limit access to information beyond pre-configured rules. Modifying and extending these systems requires specialized training, creating a bottleneck for users seeking broader data access.

This paper explores the potential of natural language processing (NLP) to bridge this gap. We present a system that leverages the power of large language models (LLMs) to automatically translate natural language questions into corresponding SQL queries. This approach empowers users to directly query the EHR database using natural language, eliminating the need for complex SQL syntax and significantly streamlining data retrieval.

The core of our system lies in a pre-trained LLM, specifically the sqlcoders model. This model, trained on paired examples of natural language questions and their corresponding SQL queries, has the remarkable ability to understand the user's intent and translate it into the appropriate database query language. We delve into the inner workings of the sqlcoders model and the concept of prompt engineering, a crucial aspect of guiding the LLM towards generating accurate SQL statements.

By focusing on open-source LLMs like sqlcoders, our work contributes to the broader exploration of readily available resources for NLP tasks. We aim to demonstrate the effectiveness of supervised fine-tuning in enhancing the performance of open-source LLMs for the challenging task of text-to-SQL translation in the specific domain of healthcare.

This paper is structured as follows. First, we discuss related work on text-to-SQL translation, highlighting the advantages of LLM-based approaches and the importance of prompt engineering. Subsequently, we introduce the sqlcoders model and its methodology. We then present our approach and implementation details, followed by

422

an evaluation of the system's performance. Finally, we conclude by discussing the implications of our work and outlining future directions.

## 2 Related Work

Extracting SQL queries from natural language questions has been a well-studied area within NLP, with applications spanning various domains. This section explores relevant research directions and highlights how the sqlcoders model aligns with these methodologies.

### 2.1 Semantic Parsing and Question Answering

Question-to-SQL generation can be viewed as a sub-task of semantic parsing, where the goal is to translate natural language into a formal representation like SQL. Early approaches relied on rule-based systems or semantic parsing methods that focused on identifying the SQL structure and filling slots with relevant information from the question .These methods achieve good performance but struggle with complex queries or domain-specific terminology. [2][6] [7][8]

### 2.2 Sequence-to-Sequence Learning

Another approach leverages sequence-to-sequence (Seq2Seq) models with attention mechanisms.[5] These models encode the natural language question and decode the corresponding SQL query directly. While effective, they may struggle with order-sensitive aspects of SQL syntax and require large amounts of training data.

### 2.3 Template-Based Methods and Prompt Engineering

Some studies adopt template-based approaches where pre-defined SQL templates are filled with question elements. While this method can handle complex queries efficiently, it relies heavily on hand-crafted templates and may not generalize well to unseen scenarios. Recent work focuses on "prompt engineering," which involves carefully crafting prompts that guide large language models (LLMs) towards generating the desired output format. The sqlcoders model aligns with this approach by utilizing a comprehensive prompt that incorporates the database schema alongside the question to improve its SQL generation capabilities.

The sqlcoders model addresses the limitations of traditional semantic parsing and Seq2Seq methods by leveraging the power of LLMs. Its ability to learn from paired examples of natural language questions and their corresponding SQL queries allows it to capture complex relationships and generate accurate SQL statements. Additionally, the focus on prompt engineering ensures that the model effectively utilizes the provided database schema information. Compared to template-based methods, the sqlcoders model is more flexible and can potentially adapt to unseen scenarios. However, similar to other LLM-based approaches, it requires careful fine-tuning for optimal performance in the specific domain of healthcare.

## 3 Methodology

This section delves into the research methodology employed to investigate the effectiveness of the sqlcoders model for automated SQL query generation from natural language questions in the healthcare domain. We exploit the model's capability to learn intricate relationships between natural language and database structures, coupled with the power of prompt engineering, to achieve this goal.

### 3.1 Data Preparation

#### 3.1.1 EHR Dataset:

We utilize a well-structured Electronic Health Records (EHR) dataset, namely MIMIC-IV dataset, [1] containing various tables (e.g., patients, medications, diagnoses) and attributes (e.g., patient ID, diagnosis code, medication name) relevant to patient information. This dataset serves as the

Figure 1: Natual language question to appropriate sql

underlying data source for generating and evaluating SQL queries.

### 3.1.2 Question-SQL Pairs:

We create a collection of question-SQL pairs specific to the healthcare domain. Each pair consists of a natural language question seeking information from the EHR data and its corresponding valid SQL query that retrieves the desired answer. Here, we can introduce an image (Figure 1) to visually represent a sample question-SQL pair.

### 3.2 The sqlcoders Model

The core component of our system is the sqlcoders model, a pre-trained large language model (LLM) specifically designed for text-to-SQL translation tasks. Here, we can delve into the mathematical intuition behind the model's functionality, but due to the potentially complex nature of LLM architectures, a high-level explanation might be more suitable for this section.

### 3.3 Conceptual Framework

The sqlcoders model can be thought of as a function that maps a natural language question (q) and a database schema description (s) to a corresponding SQL query (y). We can represent this mathematically as:

$$y = f(q, s) \tag{3.1}$$

where f represents the models functionality. This function involves a complex neural network architecture, namely transformers with attention mechanisms. During training,

the model is exposed to numerous paired examples of questions, schema descriptions, and their corresponding SQL queries. This training process allows the model to develop an internal representation that captures the intricate relationships between:

- **Natural Language Semantics:** The model identifies and encodes the meaning of words and phrases within the natural language question. This includes understanding the intent of the question (e.g., retrieval, aggregation), the entities of interest (e.g., patients, medications), and the relationships between them.

- **Database Schema Knowledge:** The model learns to represent and utilize the information provided in the schema description. This includes understanding the structure of the database (tables, attributes, data types), the relationships between tables (foreign keys), and the available data elements relevant to answering the question.

- **SQL Constructs and Syntax:** Eventually The model attempts to map the extracted meaning from the question and schema to the appropriate SQL constructs. This includes generating the core components of a query like SELECT, FROM, WHERE, and JOIN, as well as populating them with relevant attributes and conditions based on the question and schema information.

424

## 3.4 Prompt Engineering

A crucial aspect of using the sqlcoders model effectively is prompt engineering. We design a comprehensive prompt that incorporates the following elements:

- **Task Description:** This clarifies the task as generating an SQL query to answer the provided question.

- **Question Placeholder:** This section is denoted by a placeholder (e.g., [QUESTION]question[/QUESTION]) where the actual natural language question is inserted during query generation.

- **Schema Description:** This section provides a representation of the database schema, including table names, attributes, and data types. This information is essential for the model to understand the available data and construct valid SQL queries. We can consider different ways to represent the schema, such as tables with columns or a more natural language-like description.

- **Instructions:** We can optionally include instructions for the model, such as handling situations where data might be unavailable or specifying calculations for revenue or cost. These instructions further guide the model towards generating accurate and relevant SQL queries.

- **Answer Placeholder:** This section (e.g., [SQL]) serves as a placeholder where the model will generate the predicted SQL query. In case of an unanswerable question, the model would generate "null" as the answer.

By effectively combining these elements within the prompt, we provide context and guide the sqlcoders model towards generating accurate and relevant SQL statements that retrieve the intended information from the EHR data.

## 3.5 Query Generation Process

1. **Iterating Through Questions:** We iterate through the collection of natural language questions in the prepared dataset.

2. **Prompt Construction:** For each question, a prompt is constructed by inserting the question into the designated placeholder within the pre-defined prompt template. The constructed prompt and schema description are fed to the sqlcoders model.

3. **Model Prediction:** The model utilizes its learned knowledge and the provided context to predict the most likely sequence of tokens representing a valid SQL query that answers the question.

## 4 Results

This section dives deeper into the model's performance based on the Reward Scoring (RS) schemes employed for evaluation.

### 4.1 Evaluation Criteria

The model's effectiveness was assessed using four RS (Reliability Score) schemes, each representing a different level of penalty for incorrect predictions:

- **RS(0):** This is the most lenient scenario where the model receives no penalty for mistakes (c=0). In the context of question answering (QAs) alone, this score essentially reflects execution accuracy in the standard text-to-SQL task.

- **RS(5):** This scenario introduces a moderate penalty (c=5). A correct prediction earns a +1 reward, while each mistake incurs a -5 penalty. In simpler terms, every five accurate predictions compensate for one incorrect prediction.

- **RS(10):** This is considered the primary evaluation metric (c=10). Each correct prediction earns a +1, whereas each

mistake results in a -10 penalty. This means ten correct predictions are needed to outweigh a single incorrect prediction.

- **RS(N):** This scenario represents the most stringent evaluation (c=N, where N is the size of the evaluation data). Here, even a single mistake can lead to a negative overall score, even if all other predictions (N-1) are correct.

## 4.2 Model Performance

The model's performance varied significantly across the different RS schemes:

1. **RS(0):** 14.14 - This positive score in the most lenient scenario indicates that the model can generate some correct SQL queries. However, the lack of penalty for mistakes doesn't provide a clear picture of its true accuracy.

2. **RS(5):** -349.61 - The substantial drop in score compared to RS(0) suggests a high number of incorrect predictions. The moderate penalty magnifies these errors, highlighting the model's sensitivity to mistakes.

3. **RS(10):** -713.37 - This significantly lower score further emphasizes the model's shortcomings. With a stricter penalty, the negative impact of errors becomes even more pronounced.

4. **RS(N):** -84885.86 - The negative score under the most stringent evaluation highlights severe limitations. Even if the model generates a large number of correct queries, a single mistake can significantly impact the overall performance.

## 4.3 Key Findings

The model's inability to achieve positive scores under most RS scenarios indicates a fundamental limitation in generating accurate SQL queries.

The significant drop in score with increasing penalty severity demonstrates the model's

| RS Scheme | Score |
|---|---|
| RS(0) (No Penalty) | 14.14 |
| RS(5) (Moderate Penalty) | -349.61 |
| RS(10) (Main Evaluation Metric) | -713.37 |
| RS(N) (Strict Penalty) | -84885.86 |

Table 1: Model Performance under Different Reliability Scoring (RS) Schemes

susceptibility to errors. Even a moderate level of penalty leads to substantial performance degradation.

The stark contrast between RS(0) and other scores emphasizes the importance of incorporating penalties into model evaluation. It provides a more realistic assessment of the model's ability to handle real-world scenarios with potential errors.

Moreover, the negative RS(N) score reveals a lack of robustness. Even a single mistake can outweigh a large number of correct predictions, indicating the model's inability to consistently generate reliable queries.

## 5 Conclusion and Future Direction

This investigation evaluated the sqlcoder model's performance in generating SQL queries using various Reliaibilty Scoring (RS) schemes. Though the model shows a basic capability to generate some correct results (evident in the positive RS(0) score), its overall accuracy and robustness require significant improvement. .By focusing on exploration of different model architectures, enhanced error handling and incorporating human expertise, future investigations hold promise for significant advancements in this domain.

### Acknowledgments

# References

[1] Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J. Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, Li-wei H. Lehman, Leo A. Celi, and Roger G. Mark. Mimic-iv, a freely accessible electronic health record dataset. *Scientific Data*, 10(1):1, 2023.

[2] Dongjun Lee, Jaesik Yoon, Jongyun Song, Sanggil Lee, and Sungroh Yoon. One-shot learning for text-to-sql generation, 2019.

[3] Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. Ehrsql: A practical text-to-sql benchmark for electronic health records. 35:15589–15601, 2022.

[4] Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico, June 2024. Association for Computational Linguistics.

[5] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.

[6] Xiaojun Xu, Chang Liu, and Dawn Song. Sqlnet: Generating structured queries from natural language without reinforcement learning, 2017.

[7] Tao Yu, Zifan Li, Zilin Zhang, Rui Zhang, and Dragomir Radev. Typesql: Knowledge-based type-aware neural text-to-sql generation, 2018.

[8] Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from natural language using reinforcement learning, 2017.

# NYULangone at Chemotimelines 2024: Utilizing Open-Weights Large Language Models for Chemotherapy Event Extraction

**Jeff Zhang[1], Yindalon Aphinyanaphongs[1], Anthony Cardillo[1],**
[1]NYU Langone Health, MCIT Department of Health Informatics,
**Correspondence:** Yin.A@nyulangone.org

## Abstract

The extraction of chemotherapy treatment timelines from clinical narratives poses significant challenges due to the complexity of medical language and patient-specific treatment regimens. This paper describes the NYULangone team's approach to Subtask 2 of the Chemotimelines 2024 shared task, focusing on leveraging a locally hosted Large Language Model (LLM), Mixtral 8x7B (MistralAI, France), to interpret and extract relevant events from clinical notes without relying on domain-specific training data. Despite facing challenges due to the task's complexity and the current capacity of open-source AI, our methodology highlights the future potential of local foundational LLMs in specialized domains like biomedical data processing.

## 1 Introduction

The extraction of structured information from unstructured clinical narratives is a crucial task in healthcare informatics, enabling better patient care and clinical decision-making. The Chemotimelines 2024 shared task focuses on extracting chemotherapy treatment timelines from clinical narratives, a challenging task for understanding oncology patients' treatment paths. Our team, NYULangone, participated in Subtask 2, aiming to leverage the general reasoning capabilities of large language models (LLMs) for this purpose.

## 2 Related Work

Clinical narrative processing traditionally relies on rule-based systems or machine learning models trained on domain-specific annotated data. Recent advances in NLP have seen the rise of transformer-based models and LLMs, offering powerful general-purpose language understanding capabilities. However, their application in domain-specific tasks like chemotherapy timeline extraction remains in the infancy of exploration.

## 3 System Description

Our system builds upon a locally deployed instance of Mixtral, an open-weights LLM. The system comprises two rounds of text inference: the first round is an extraction of chemotherapy events from individual notes, and the second round is the aggregation of events from multiple notes to a single timeline.

---

**Algorithm 1** Patient Chemotherapy Summary Algorithm

---

1: **for each** patient **do**
2:      **for each** note of the patient **do**
3:          Prompt Mixtral to read the note and extract chemotherapies
4:      **end for**
5: **end for**
6: Prompt Mixtral to combine the extracted chemotherapies from every note to create a patient-level summary of all chemotherapies

---

### 3.1 Architecture

We employed Mixtral 8x7B v0.1, an open-weights LLM originally published by Mistral AI in December 2023. The system leverages its pre-trained weights without further domain-specific fine-tuning. The system processes clinical narratives as raw text files, uses the LLM to extract relevant events and dates, and structures them into the required JSON format for output.

### 3.2 Implementation

The system was hosted on NYU Langone's high-performance cluster "Ultraviolet." Using SLURM, a compute instance was requisitioned using three NVIDIA A100s with 128GB of system RAM. The model weights for Mixtral 8x7B were downloaded from Hugging Face, and inference was performed with the Transformers library for Python.

### 3.3 Prompts

For the first inference used to extract chemotherapy events from notes, we used the following Markdown-style prompt:

[INST] **GOAL and PURPOSE:** You are an experienced medical annotator with special expertise in natural language processing of oncology documents. You will be given a list of JSON objects to turn into a list of lists.

**INSTRUCTIONS:** Read the patient's note in its entirety, given in the section "# PATIENT NOTE" below. Use THYME guidelines to create "events"; every mention of a chemotherapeutic drug or component should have: the name of the drug, an associated date, the temporal_relation between the use of that drug and the associated date. Each event must be in the form ['chemo drug name', 'temporal_relation', 'YYYY-MM-DD']. If a drug is associated with multiple dates, or a date is associated with multiple drugs, break them into separate events. 'temporal_relation' must be one of ["contains-1", "begins-on", "ends-on", "before"].

**EXAMPLES:** ['herceptin', 'begins-on', '2013-06-17'], ['taxol', 'contains-1', '2013-09']

**OUTPUT:** Use only well-formatted JSON. Only output the timeline of chemotherapy events; place it under "# TIMELINE". Do not make any additional notes or comments, only JSON under "# TIMELINE". [/INST] **PATIENT NOTE** <insert patient note here> **TIMELINE**

This first inference accomplishes the extraction of each chemotherapy event in each note. However, the events are not organized by patient yet. For the second inference used to aggregate chemotherapy events from multiple notes into patient timelines, we used the following prompt:

[INST] **GOAL and PURPOSE:** You are an experienced medical annotator with special expertise in natural language processing of oncology documents. You will be given a JSON list of lists. Your job is to output a list of lists for each patient.

**EXAMPLE OUTPUT:**

```
patient_01:
['taxol', 'begins-on', '2013-06-17']
['taxol', 'ends-on', '2013-09']
...
patient_02:
```

[/INST]

## 4 Results

On the dev set, our system achieved an average F1 score of 0.35. On the validation set, our system achieved an average F1 score of 0.23 across different cancer types, as shown in Table 2 of the competition results.

## 5 Discussion

While our performance was well below the baseline and leading teams, it provided valuable insights into the challenges and potential of using locally hosted LLMs in clinical NLP tasks without domain-specific training.

The opaque inner workings of LLMs preclude an exact understanding of why certain chemotherapy events are more easily extracted than others. The errors our system demonstrates could largely be grouped into several types:

- Confabulation of drugs not mentioned (e.g. extracting "herceptin" from a patient radiology report without any mention of chemotherapy)

- Inclusion of non-chemotherapeutic drugs, especially steroids (e.g. extracting "prednisone" for a patient on immunosuppression)

- Failure to include clearly mentioned drugs (e.g. failing to extract "aflibercept" when it was well documented in a patient note)

Despite the objectively poor performance, our results highlight a future potential for LLMs to be used in biomedical NLP tasks. Local LLMs that can competently perform general reasoning are still a new technology, with expert opinion suggesting that local models like Mixtral currently perform at a GPT-3 (OpenAI, United States) level of performance.

# 6  Conclusion and Future Work

Our exploration into using local LLMs for chemotherapy treatment timelines extraction offers a starting point for further research in this area. Future work will focus on enhancing model understanding of clinical contexts through retrieval augmented generation (RAG) and ensemble prompting techniques such as "tree of thought."

# 7  Acknowledgments

# References

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scibert: A pretrained language model for scientific text. 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 3615–3620, Hong Kong, China. Association for Computational Linguistics.

Tom Brown and et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Jacob Devlin and et al. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Alvin Rajkomar, Eyal Oren, Kai Chen, Andrew M Dai, Nissan Hajaj, Michaela Hardt, Peter J Liu, Xiaobing Liu, Jake Marcus, Mimi Sun, Patrik Sundberg, Hector Yee, Kun Zhang, Yonghui Zhang, Gerardo Flores, Gavin E Duggan, Jamie Irvine, Quoc Le, Kurt Litsch, Alexander Mossin, Justin Tansuwan, De Wang, James Wexler, Jimbo Wilson, Dana Ludwig, Samuel L Volchenboum, Katherine Chou, Michael Pearson, Srinivasan Madabushi, Nigam H Shah, Atul J Butte, Michael D Howell, Claire Cui, Greg S Corrado, and Jeff Dean. 2018. Scalable and accurate deep learning with electronic health records. *NPJ Digital Medicine*, 1.

Guergana K Savova, James J Masanz, Philip V Ogren, Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-Schuler, and Christopher G Chute. 2010. Mayo clinical text analysis and knowledge extraction system (ctakes): architecture, component evaluation and applications. *Journal of the American Medical Informatics Association*, 17(5):507–513.

Ashish Vaswani and et al. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.

Junyuan Yao, Harry Hochheiser, Woosub Yoon, Erin Goldner, and Guergana Savova. 2024. Overview of the 2024 shared task on chemotherapy treatment timeline extraction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop, NAACL June 2024*, Mexico City, Mexico.

# AIRI NLP Team at EHRSQL 2024 Shared Task:
# T5 and Logistic Regression to the Rescue

**Oleg Somov[1,2], Aleksei Dontsov[1], Elena Tutubalina[1,3,4]**

[1]AIRI, [2]MIPT, [3]Sber AI, [4]Kazan Federal University

**Correspondence:** somov@airi.net

## Abstract

This paper presents a system developed for the Clinical NLP 2024 Shared Task, focusing on reliable Text-to-SQL modeling on Electronic Health Records (EHRs). The goal is to create a model that accurately generates SQL queries for answerable questions while avoiding incorrect responses and handling unanswerable queries. Our approach comprises three main components: a query correspondence model, a Text-to-SQL model, and an SQL verifier. For the query correspondence model, we trained a logistic regression model using hand-crafted features to distinguish between answerable and unanswerable queries. As for the text-to-SQL model, we utilized T5-3B as a pre-trained language model, further fine-tuned on pairs of natural language questions and corresponding SQL queries. Finally, we applied the SQL verifier to inspect the resulting SQL queries. During the evaluation stage of the shared task, our system achieved an accuracy of 68.9% (metric version without penalty), positioning it at the fifth-place ranking. While our approach did not surpass solutions based on large language models (LMMs) like ChatGPT, it demonstrates the promising potential of domain-specific specialized models that are more resource-efficient. The code is publicly available at https://github.com/runnerup96/EHRSQL-text2sql-solution.

## 1 Introduction

Electronic health records (EHRs) play a critical role in storing comprehensive medical histories within hospital settings, capturing everything from patient admissions to treatment and discharge. However, efficiently retrieving relevant information from these records remains a significant challenge, particularly from complex medical relational databases.

This paper focuses on enhancing a text-to-SQL system specifically designed for the medical domain. The aim is to improve the retrieval pro-



Figure 1: The schema intersection algorithm. We match normalized n-grams of an input natural language question "How much is the cost for the drug nystatin cream?" against normalized database (DB) content. As we can see on schema intersection count distribution, the NULL questions have much less schema intersection elements in comparison to SQL questions. For more details, please refer to Section 3.

cess of patient information and enable better clinical decision-making. The objective is to develop a model capable of accurately generating SQL queries for answerable questions while effectively handling unanswerable queries and avoiding incorrect responses. In other words, when faced with unanswerable questions, the model should refrain from generating any SQL prediction and indicate the absence of an answer by returning NULL.

To conduct experiments, we utilize the Text-to-SQL benchmark (Lee et al., 2022) provided by the organizers of the Clinical NLP 2024 task (Lee et al., 2024). This benchmark consists of pairs of input utterances and expected SQL queries, including cases where generating an SQL query is impossible for a given question. This dataset is linked to two open-source EHR databases—MIMIC-III (Johnson et al., 2016) and eICU (Pollard et al., 2018). This benchmark includes questions that address the actual needs of a hospital and incorporate various time expressions crucial to daily healthcare work.

In this paper, we describe our solution for the

Clinical NLP 2024 shared task on reliable Text-to-SQL. As shown in Figure 2, our system consists of three components - query correspondence model, Text-to-SQL model, and SQL result inspector. To sum up, the system takes the user's query as input and goes through the following steps: feature extraction, query scoring using the query correspondence model and alignment check, SQL generation using a Text-to-SQL model with question and schema input representation, SQL results inspection, execution of the generated query, and checking if the execution result meets the requirements. If the requirements are met, the system returns the result to the user.

The paper is organized as follows. Section 2 presents related work on Text-to-SQL corpora and state-of-the-art (SoTA) models. We describe our model with three components in Section 3. Experiments with baselines and our model are presented in Section 4. Finally, we discuss errors and conclude the work in Sections 5 and 6, respectively.

## 2 Related work

Text-to-SQL currently is one of the most developing and promising research areas in the field of semantic parsing. Well-known public leaderboard Spider (Yu et al., 2018) popularized the task, and Text-to-SQL domain developed many directions and specializations. Spider dataset is a complex and cross-domain Text-to-SQL dataset which consists of 10181 questions with 5693 SQL queries on 200 databases. The main goal of the dataset is to generalize to new databases. However, the Spider dataset does not contain unanswerable questions. Spider also gave rise to more complex datasets like BIRD (Li et al., 2024), which paid attention not only to SQL query complexity (introduction of window functions, etc.) but also to the optimality of the generated query. BIRD databases are close to real-world examples, with tables consisting of millions of data rows; hence, the optimal SQL is required.

Another dataset named CoSQL (Yu et al., 2019) raised questions about ambiguity and the system's ability to handle such questions. It is a dialogue-based Text-to-SQL benchmark, which consists of the following dialogue acts - answering user questions with SQL, double checking the user intent if the questions are ambiguous, or the system reminder to the user that the question is not related to the database.

Spider leaderboard gave rise to specialized Text-to-SQL architectures. Naturally, SoTA solutions adapted the following Text-to-SQL solutions - schema linking stage, encoding of question and schema, and subsequent decoding. Starting from most notable solutions like BRIDGE, which induced database content into training process (Lin et al., 2020) and RAT-SQL, which modified transformer architecture for question with schema interaction and specialized grammar-based decoding process (Wang et al., 2021) coming to the fine-tuning approaches which reached its peak in RES-DSQL (Li et al., 2023) approach and PICKARD (Scholak et al., 2021). LLMs are also present in the leaderboard in the form of in-context learning few-shot approaches(Gao et al., 2023a) and SQL debugging stages (Pourreza and Rafiei, 2024). Most solutions utilize ChatGPT-4 as a core model and experiment with different prompt strategies for stages of schema linking, query generation, and SQL debugging.

Increased attention towards Text-to-SQL domain detected the problem of generalization in semantic parsing. The Spider dataset focused on cross-domain generalization, but the work of (Suhr et al., 2020) made the challenges more visible, introducing the challenges of single database split compared to cross-database setting. Recently, Somov and Tutubalina (2023) evaluated the generalization capabilities of supervised models on the original, multilingual, and target length splits of the improved version of the Spider dataset called PAUQ (Bakshandaeva et al., 2022). Results indicate that the models can generalize well to unseen simple SQL's, while multilingual split shows that some models benefit from learning on the translated task.

Overall, the ongoing progress in dataset development and the advancement of specialized architectures have significantly contributed to a deeper understanding of the Text-to-SQL task and its applications across various domains, including medicine.

## 3 Main method

Our final solution consists of 3 components - query correspondence module, fine-tuned Text-to-SQL model, and SQL result inspector, which checks the result of the generated query. The system pipeline is presented in Figure 2. The system output can be NULL if the system considers the query unanswerable or results if the system can answer the query. This section will describe our validation schema

Figure 2: The system overview. The user query inputs into the Text-to-SQL system. The feature extractor extracts features for the query correspondence model. The query correspondence model scores the query by extracted features. If the question is aligned with our system, we pass the input question into Text-to-SQL generation model. It consists of question and schema input representation component and Text-to-SQL model. The generated query is passed to the SQL results inspector, which checks weather the query can be executed and checks the result of the execution. If the query execution result meets the requirement, we return the result to the user.

and all the system components in detail.

## 3.1 Validation schema

For our method evaluation, we have developed our validation schema. Our solution consists of two machine learning models - Query Correspondence model and Text-to-SQL model. The leader board submission of NULL revealed that the evaluation and test sets consist of approximately 20% NULL's while our training set has approximately 9% of NULL's. For the Query Correspondence model, we have prepared a similar test distribution - the training set has 10% of NULLs while the validation set has 20% of nulls. Since we do not observe the distribution shift for SQL question, for Text-to-SQL model, we prepared an i.i.d. splitting for evaluation.

## 3.2 Query correspondence model

The query correspondence model (QCM) is a component that analyzes input questions and discards them if they look like questions that can not be answered based on database content. It consists of two components - a feature extractor and a machine learning model. We get the input question and run preprocessing. The preprocessing steps include - punctuation cleaning, stop-word exclusion, lemmatization, and lowercase casting. Then, we extract 3 features from the processed question - schema intersection feature, first-word feature, and query length feature.

- Schema intersection feature is the number of elements from the database(attributes, tables, values) in the question. We extract and preprocess database content with punctuation

cleaning, lemmatization, and lowercase casting. We merge attributes, tables, and values into one set. The processed question is tokenized by spaces and transformed into another set. The intersection between these two sets is the result feature value. On the public test set, the feature for detection of NULL scored 18.77%, detecting 94% of NULL questions. Since the feature proved to be important for NULL question discarding, we later used it in our experiment as a decision component with a manually selected threshold. The schema intersection algorithm and corresponding intersection count distribution is presented in Figure 1.

- We examined the intersection of the processed NULL questions beginning (first 2 words) with processed SQL questions beginning and found out that there is only an 8% intersection between sentence beginnings. Therefore, we matched all the NULL first 2 words against the input question. If the input processed sentence is matched against processed null sentence beginnings, the feature value is True and False otherwise.

- We have analyzed NULL question length and SQL question length and saw that the average length of NULL questions is 11 ($\sigma = 3$), and the average length of SQL questions is 15 ($\sigma = 6$). Due to such differences, we also utilized question process length as a feature.

During our experiments, we have also used other features, like pre-trained language model maximum entropy score, SVM classifier (Vapnik and

Chervonenkis, 1974) score based on TF-IDF encoder, pre-trained Transformer (Devlin et al., 2018) encoder-based retrieval features (distance to closest question with SQL, distance to closest question with NULL, number of NULL candidates @5/@10/@100 - but these features made our results only worse on public test set, so we have discarded them.

We pass these selected features through normalization and then pass them to a logistic regression model. We trained this model on the binary task on standardized extracted features and predicted SQL vs. NULL for every question. We evaluate our solution based on two metrics: sensitivity (Se) and specificity (Sp).

$$Se = \frac{TP}{TP + FN}, \quad Sp = \frac{TN}{TN + FP}$$

On our split, the model gets the average of sensitivity and specificity equal to $0.91$ on our validation set with a $0.5$ threshold. If the question is predicted as NULL, we do not pass the question further.

### 3.3 Text-to-SQL model

The next step is the Text-to-SQL model. We chose the T5-3B (Raffel et al., 2020) model, and our evaluation showed its high performance. If the query correspondence model evaluates the question as answerable, we pass the question to the text-to-SQL model. We wanted the model to learn only Text-to-SQL task. Therefore, we have trained the model on only text-to-SQL pairs. On our validation set, the execution match from the benchmark evaluation script with this model was Acc0 = 99%, Acc5= 92%, Acc10=86%, AccN=−501% on Text-to-SQL pairs only. We use classic input representation 1 as a concatenation of database name, question, and linearized schema representation of tables $T$ and columns $C$ (Shaw et al., 2021).

$$X = Database\ name : Question\ |$$
$$[T_1] : [C_{11}], ...[C_{1|T_1|}] \mid [T_2] : [C_{21}], ... \quad (1)$$

We normalized the target SQL query with classic Spider Text-to-SQL fine-tuned model preprocessing as in RESDSQL (Li et al., 2023). The target representation during training was the following:

$$Y = Database\ name \mid Query \quad (2)$$

We have trained our T5-3B model for 16 epochs(approximately 4000 iterations) with a training batch size of 2 and gradient accumulation batch



Figure 3: T5-3B training process on Text-to-SQL custom i.i.d. validation split from section 3.1. On the right exact match training plot, we see that the model decently learns to correctly align novel questions to SQL queries as the validation loss decreases.

size of $8$. The learning rate was 5e-5. Our input maximum length was $800$, and the target length was $514$. As demonstrated in Figure 3 we see that the model successfully converged and reached decent exact match accuracy.

### 3.4 SQL result inspector

After generation, we pass the result SQL to the SQL inspector. We rely on the hypothesis that the user must be very specific in his question to correctly match the elements of the schema in his question and get an answer to his question. Therefore, if the query fails and returns a None or 0 value for aggregate queries - we treat it as a false and exit with NULL. We examined the training SQL query outputs and discovered that approximately $90\%$ of training queries return some meaningful result, meaning not None or 0 value for aggregate queries. We evaluate the SQL inspector on the EHRSQL train set - we run T5 prediction through it and evaluate how many generated SQL queries the inspector will discard and how many will approve. As in the Query Correspondence model, we measure sensitivity and specificity as in Equations 3.2. We get the average of sensitivity and specificity of $96\%$.

## 4 Experiments

This section describes our most successful attempts on the test leaderboard. Solutions feature schema intersection algorithm explained in Sec. 3.2 and SQL inspector explained in Sec. 3.4. In Table 1, we present official evaluation scores[1]: Accuracy0, Accuracy5, Accuracy10, AccuracyN. These metrics differ in penalty strategy for wrong predictions. In particular, Accuracy0 does not penalize any mis-

---

[1]For details, see https://www.codabench.org/competitions/1889/

434

| | Method | Accuracy0 | Accuracy5 | Accuracy10 | AccuracyN |
|---|---|---|---|---|---|
| 1 | Schema intersection@2 + ChatGPT ICL 5-shot | 55 | -41.8 | -138.6 | -22545 |
| 2 | Schema intersection@2 + T5-3B + ChatGPT debugger + SQL inspector | 53.3 | -27.2 | -107.8 | -18746.7 |
| 3 | Schema intersection@2 + T5-3B + SQL inspector | 64.4 | 53.3 | 42.2 | **-2535.6** |
| 4 | QCM + T5-3B + SQL inspector | **68.9** | **56.5** | **44** | -2831.1 |

Table 1: Experimental results of our systems on the official test set.

takes, while Accuracy5 counts a -5 penalty for each mistake result.

## 4.1 ChatGPT: in-context learning with few-shot examples

In 2020, the paradigm of in-context learning, introduced by Radford et al. (2019), emerged as a powerful technique that enables Language Model Models (LLMs) to solve problems without requiring fine-tuning. We effectively leverage the potential of few-shot learning by exposing the model to a few examples from the training set along with their corresponding solutions. To facilitate this process, we create an index of training questions with corresponding SQL query by extracting embedding of the question using SentenceBERT[2] (Reimers and Gurevych, 2019). To identify the most similar matches, we calculate the Euclidean distance between the index question vectors and the embedding of the natural language question. These selected questions, along with their corresponding SQL queries, are included in the prompt.

Furthermore, to provide the LLM with an understanding of the database's structure, we append a textual representation of the entire database schema and question at the end of the prompt. This approach mirrors the methodology employed in the DAIL-SQL technique (Gao et al., 2023b).

The final prompt is further passed into OpenAI API[3], model version `gpt-3.5-turbo`.

After gathering the results, we filtered out queries that did not pass our schema intersection manual threshold of 2.

## 4.2 ChatGPT debugger

Recent advances on the Spider leaderboard showed that the ChatGPT can not only work as an in-context learning algorithm but can also refine a given query. We have utilized the DIN-SQL (Pourreza and Rafiei, 2024) approach for self-correction. To address this, DIN-SQL proposed

a self-correction module where the model is instructed to correct those minor mistakes. This is achieved in a zero-shot setting, where only the buggy code is provided to the model, and it is asked to fix the bugs. DIN-SQL proposed two different prompts for the self-correction module: generic and gentle. The generic prompt, DIN-SQL requests the model to identify and correct the errors in the "BUGGY SQL". The gentle prompt, on the other hand, does not assume the SQL query is buggy; instead, it asks the model to check for any potential issues and provides some hints. Since our Text-to-SQL T5-3B model performed well on our validation split, we have utilized a gentle approach for the model to fix potential bugs. We have used the original implementation [4]. Also, DIN-SQL experiments showed that a gentle prompt is more effective for the GPT-4 model, which proved to be better at this task. After gathering the results, we filtered out queries that did not pass our schema intersection manual threshold of 2. The query debugger algorithm resulted in quality deterioration in comparison to our final solution. The GPT-4 debugger of generated T5 queries usually just deleted some comparisons or conditions from the final query, making more false positive predictions.

## 4.3 RESDSQL fine-tuning

We have also tried to fine-tune the RESDSQL solution to the EHRSQL task. RESDSQL is fine-tuned SoTA on the Spider leaderboard. It consists of two training phases - cross-encoder classifier for question-relevant columns and tables detection and query generation stage via a pre-trained language model. During training in the query generation stage, the decoder input is prefixed with SQL skeleton, forcing the model to generate a correct SQL template and then fill it with schema elements and values. Although the cross-encoder component had a validation AUC score for detection of tables and columns 97.7%, the result execution accuracy0 on

Figure 4: The distribution of embedding in two-dimensional space via t-SNE of our custom split of original training EHRSQL data of train and validation and along with test questions as in 3.1 for evaluation of Query Correspondence Model. **train**$_{SQL}$ stands for train question embeddings which have SQL, **train**$_{NULL}$ which do not.**val**$_{SQL}$ stands for val question embeddings which have SQL, **val**$_{NULL}$ which do not. **test** stands for EHRSQL test question embeddings.

our Text-to-SQL validation set was 77.5%. We suspect the problem is the inconsistency of the SQL skeleton with the target SQL query.

For example, almost every second SQL EHRSQL query contains strftime function, which includes two attributes; however, the RESD-SQL SQL skeleton, which is prefixed in the query generation model, contained only one attribute.

### 4.4 Experimental results

As shown in Table 1, we have conducted 4 submission experiments. All of our experiments use the schema intersection feature. The first three use it as a decision feature, while the final solution model uses it as a feature in QCM. We see that the first experiment was the worst. There was no SQL inspector phase, and ChatGPT itself generated incorrect queries. Then, we enhanced our solution with the SQL inspector component and used ChatGPT as a debugger for our T5 predictions. Unfortunately, the debug mode worsened SQL predictions, but due to the SQL inspector, we had fewer false positive predictions, as we can see in the accuracies of a penalty. We concentrated on purely T5-3B and other components in our following experiments. At first, we used only the schema intersection feature to discard unanswerable queries, but after careful

exploratory data analysis, we found more features to be a good signal for our decision - so we developed a query correspondence model, which gave us the highest score.

Although we have good accuracy across all of our components on our validation split, we have much lower results on the test leaderboard. In Figure 4, we have plotted reduced SentenceBERT questions embeddings via the t-SNE algorithm on a coordinate plane. We see that the testing questions are shifted relatively to training data in terms of SQL and NULL questions. Although we also mimic data drift in our validation schema as in section 3.1, our validation questions are still closer to training questions than test questions. The solution to that problem might be running a solution in production mode with activated data markup for online and offline metrics alignment.

## 5 Error analysis

We manually checked the errors of our final system components on our validations sets from 3.1. The Text-to-SQL T5-3B errors mostly consist of regular errors - ASC to DESC mismatch, wrong column, missed comparison expression. Sometimes the model shows the overfitting signs - looping the prediction output, adding wrong syntactic constraints (like adding not needed GROUP BY) or extra symbol to value ('10-31' $\rightarrow$ '10-31\'). 

The query correspondence errors come mostly from the starting word feature - although it helps to identify questions with starting phrases that were in the training set, it does not help to combat novel starting phrases that occur in the test set.

As we pointed out, we evaluated our query inspector on sensitivity and specificity metrics. Specificity is 99%, and sensitivity is 94%. We can see that we are stricter than necessary to generated queries, and sometimes correct SQL can return the result of None or 0, but we will not return it to the user.

## 6 Discussion and conclusion

In this work, we have built a reliable Text-to-SQL solution. We have developed our validation schema for model evaluation and submitted our final system results to the EHRSQL leaderboard. Our solution consists of 3 components - query correspondence model, Text-to-SQL generation model, and SQL inspector. During validation, we measured our performance based on sensitivity and specificity

metrics to account for NULL queries and execution match for the query generation model. The sensitivity and specificity metrics for the query correspondence model and SQL inspector are $81\%/99\%$ and $94\%/99\%$ accordingly. The execution accuracy of Text-to-SQL model is $99\%$. Our components are independent of each other; therefore, we can calculate the product probability that the NULL question will be discarded is $98\%$, while the probability that the SQL question will be answered correctly is $75\%$.

The advantages of our system are the following:

- The solution discards unanswerable queries with high precision while keeping a decent execution accuracy.

- Our components can be independently optimized.

- The solution is interpretable because every component has its single responsibility.

- Our model can be used on-premise without confidential data leaks to external language models.

The disadvantages of the system are:

- The cascading effect of the system leads to lower execution accuracy.

- Weak out-of-distribution robustness.

- We employ a heavy Text-to-SQL T5-3B pre-trained language model, which needs significant resources for deployment.

As a future work direction, we see the necessity of developing reliable, robust, and lightweight specialized solutions. These solutions can be run and maintained on-premise without exposing personal data to external LLMs.

## Acknowledgments

## Ethics Statement

One limitation of using databases for retrieval is that these sources may not be complete and can include errors. T5-3B, like any language model, may be subject to representation biases and potentially misleading results, which is a critical concern in the healthcare domain.

All pre-trained language models and datasets used in this work are publicly available for research purposes.

We honor and support the ACL Code of Ethics.

## References

Daria Bakshandaeva, Oleg Somov, Ekaterina Dmitrieva, Vera Davydova, and Elena Tutubalina. 2022. Pauq: Text-to-sql in russian. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2355–2376.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2023a. Text-to-sql empowered by large language models: A benchmark evaluation. *arXiv preprint arXiv:2308.15363*.

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2023b. Text-to-sql empowered by large language models: A benchmark evaluation.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Haoyang Li, Jing Zhang, Cuiping Li, and Hong Chen. 2023. Resdsql: Decoupling schema linking and skeleton parsing for text-to-sql. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 13067–13075.

Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. 2024. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36.

Xi Victoria Lin, Richard Socher, and Caiming Xiong. 2020. Bridging textual and tabular data for cross-domain text-to-sql semantic parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4870–4888.

Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. 2018. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data*, 5(1):1–13.

Mohammadreza Pourreza and Davood Rafiei. 2024. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Advances in Neural Information Processing Systems*, 36.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

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. *The Journal of Machine Learning Research*, 21(1):5485–5551.

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.

Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. PICARD: Parsing incrementally for constrained auto-regressive decoding from language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9895–9901, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Peter Shaw, Ming-Wei Chang, Panupong Pasupat, and Kristina Toutanova. 2021. Compositional generalization and natural language variation: Can a semantic parsing approach handle both? 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 922–938, Online. Association for Computational Linguistics.

Oleg Somov and Elena Tutubalina. 2023. Shifted pauq: Distribution shift in text-to-sql. In *Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation in NLP*, pages 214–220.

Alane Suhr, Ming-Wei Chang, Peter Shaw, and Kenton Lee. 2020. Exploring unexplored generalization challenges for cross-database semantic parsing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8372–8388, Online. Association for Computational Linguistics.

Vladimir Vapnik and Alexey Chervonenkis. 1974. Theory of pattern recognition.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2021. Rat-sql: Relation-aware schema encoding and linking for text-to-sql parsers.

Tao Yu, Rui Zhang, He Yang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, et al. 2019. Cosql: A conversational text-to-sql challenge towards cross-domain natural language interfaces to databases. *arXiv preprint arXiv:1909.05378*.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. *arXiv preprint arXiv:1809.08887*.

# IKIM at MEDIQA-M3G 2024: Multilingual Visual Question-Answering for Dermatology through VLM Fine-tuning and LLM Translations

**Marie Bauer**[1]   **Constantin Marc Seibold**[1]   **Jens Kleesiek**[1,2,3,4]   **Amin Dada**[1]

[1]Institute for AI in Medicine (IKIM), University Hospital Essen (AöR), Essen, Germany
[2]Cancer Research Center Cologne Essen (CCCE), West German Cancer Center Essen
University Hospital Essen (AöR), Essen, Germany
[3]German Cancer Consortium (DKTK, Partner site Essen), Heidelberg, Germany
[4]Department of Physics, TU Dortmund, Dortmund, Germany

## Abstract

This paper presents our solution to the MEDIQA-M3G Challenge at NAACL-ClinicalNLP 2024. We participated in all three languages, ranking first in Chinese and Spanish and third in English. Our approach utilizes LLaVA-med, an open-source, medical vision-language model (VLM) for visual question-answering in Chinese, and Mixtral-8x7B-instruct, a Large Language Model (LLM) for a subsequent translation into English and Spanish. In addition to our final method, we experiment with alternative approaches: Training three different models for each language instead of translating the results from one model, using different combinations and numbers of input images, and additional training on publicly available data that was not part of the original challenge training set.

## 1 Introduction

Over the past 25 years, various studies have discussed the shortage of dermatologists in the US (Kimball, 2003; Kimball and Resneck Jr, 2008; Ehrlich et al., 2017). At the same time, machine learning methods offer potential relief for the limited time available to dermatologists (Fogel and Kvedar, 2018) and, on some tasks, even exceed expert capabilities (Esteva et al., 2017). Recently introduced vision-language models (VLMs) showed promising capabilities in radiology and pathology visual question-answering (VQA) tasks (Moor et al., 2023; Wu et al., 2023; Thawkar et al., 2023; Liu et al., 2023; Chen et al., 2024). Therefore, it can be assumed that they also provide relief in the field of dermatology. However, there are no existing dermatology VQA datasets (Lin et al., 2023). Yet, VLMs need fine-tuning datasets to achieve the high accuracy required for medical tasks (Liu et al., 2023).

A possible data source for such tasks are telemedical records. Telemedicine describes triag-

ing, diagnosing, and monitoring patients remotely through digital images and text messages (Waller and Stotler, 2018). Shortly after the outbreak of the COVID-19 pandemic, the availability of telemedicine services increased in parts of China (Hong et al., 2020; Song et al., 2020), providing new opportunities to create VQA datasets. Following these developments, the MediQA-M3G challenge (wai Yim et al., 2024a) is based on data from one of these telemedical platforms. The participants are offered photos of skin diseases and textual interactions between patients and medical professionals. While the original data is in Chinese, automated translations into English and Spanish were also provided. This raises several questions that we examined in the course of the challenge. First, there is the question of which model should be used on the Chinese data since all medical VLMs were trained in English. Another question is how helpful the training on the translated English and Spanish data is or whether problems such as translation errors and cultural differences are a hindrance.

To answer these questions, this paper compares various fine-tuning methods in preparation for our challenge submission. We first evaluate the usefulness of additional imaging data from two publicly available dermatological classification datasets in solving the challenge. We then compare multi-image training to training with a single image per data entry. Finally, we also test if training three different models for each language outperforms training a single model and translating its predictions into the other two target languages.

## 2 Related Work

### 2.1 VLMs

With the rapid development of LLMs (Hoffmann et al., 2022; Touvron et al., 2023a,b; Peng et al., 2023), various approaches have been pursued to extend these models to vision-language models

(VLMs) (Alayrac et al., 2022; Li et al., 2023b; Liu et al., 2023). This usually involves combining pre-trained LLMs and image models using dedicated architectures and training them on multimodal data. A notably straightforward yet effective architecture that has emerged from these efforts is LLaVA (Liu et al., 2023). Within this approach, a basic feed-forward network, comprising two layers is employed to map the image embeddings to the language embedding space of the LLM. Similarly to the development of specialized biomedical LLMs (Chen et al., 2023; Labrak et al., 2024; Xie et al., 2024), a modified version of LLaVA designed for biomedical applications, known as LLaVA-med (Li et al., 2023a), has been introduced. All of our solutions to this challenge are based on LLaVA-med.

## 2.2 Translation

Shortly after the release of ChatGPT and the subsequent focus on LLMs, their translation ability was explored (Hendy et al., 2023; Jiao et al., 2023; Bawden and Yvon, 2023). In contrast to previous neural machine translation (NMT) approaches that revolved around specialist language models trained on parallel translation corpora (Tiedemann and Thottingal, 2020; Costa-jussà et al., 2022), LLMs learn translation through vast pre-training and instruction tuning. Improvements over traditional NTM models include smoother translations (Hendy et al., 2023). However, these improvements are accompanied by higher translation error rates (Yao et al., 2024). An interesting observation by Hendy et al. (2023) is that GPT produces more accurate translations of noisy Chinese texts than traditional NMT models. Since the data in this challenge consists of Chinese consumer health questions, a translation with LLMs seems reasonable in this context. However, it also makes sense to evaluate a traditional NMT model due to the higher error rates of LLMs. Following its promising performance on medical downstream tasks (Dada et al., 2024), we used Mixtral-8x7B-Instruct (Jiang et al., 2024) for LLM-based translation of Chinese predictions and OPUS (Tiedemann and Thottingal, 2020) as the NMT model.

## 2.3 Consumer Health Question-answering

Previous works focused on consumer health question-answering but were only text-based (Ben Abacha et al., 2019; Ben Abacha and Demner-Fushman, 2019). Existing VQA datasets do not include consumer health inquiries and are based

on radiology (Lau et al., 2018; Liu et al., 2021; Hu et al., 2023) and pathology images (He et al., 2020). Since no datasets are based on multimodal dermatology consumer health questions, there are currently no existing approaches for this task. Furthermore, using Chinese texts translated into English and Spanish is a novel approach that requires methods to address this setting adequately.

## 3 Challenge Dataset

The given dataset (wai Yim et al., 2024b) consists of clinical history and patient query examples. Along these textual inputs, one or multiple photos of the described skin disease were attached to the query. The gold labels consisted of one or multiple answers by Chinese dermatologists. All texts were machine-translated into English and Spanish without further information on which model was used for translation. One exception is the test set, which was translated manually. For the validation and test sets, annotators also provided a score indicating how complete an answer is concerning the query. Possible completeness scores were $0.0$, $0.5$, and $1.0$, ranging from incomplete to entirely complete. As a metric, the deltaBLEU score (Galley et al., 2015) was computed between predictions and dermatologists' answers using the completeness score for weighting.



Figure 1: Histogram of number of words per dermatologist answer in English.

The training set consists of $842$ patient queries with an average of $2.94$ images per query. Additional $56$ examples were provided as a validation set and $100$ examples as a test set. Figure 1 shows the histogram of the number of words per dermatologist answer for the English training set. Most answers consist of only a few words, usually the diagnosis. However, some outliers are considerably longer, with over $315$ words. These answers contain lengthy descriptions of the treatment and

follow-up steps for the patient. While we manually analyzed the data, we could not find a consistent relationship between the type of request and the verbosity of the response.



Figure 2: Histogram of number of images per patient query.

Figure 2 shows the histogram of the number of images per sample in the training set. Like the number of words per answer, queries usually have few images attached.

## 4 Methodology

The following section describes the different approaches for the challenge. We describe multi-image training, additional non-challenge data used, and our methods of translating LVM predictions into new target languages.

### 4.1 Training on multiple input images

The challenge data often provided multiple images for each input text (see Figure 2). This led to the question of whether all of them should be used together in a single text prompt, decreasing the number of training examples but potentially increasing the information available to the model in each case, or if each image should be used as input separately, thus increasing the number of training examples but potentially decreasing the quality of the input.

### 4.2 Additional fine-tuning data

In addition to the data presented by the challenge, we attempted to train the model on additional publicly available dermatological image datasets. For this, we employed Fitzpatrick17k (Groh et al., 2022), which contains approximately 17,000 labeled dermatological images and Dermnet[1], adding an additional 19,500 images. The aim was to increase the model's overall domain knowledge and

---

[1] https://dermnet.com/

to improve its performance in identifying common dermatological illnesses before training it on challenge data. We prompted the model to identify the illness in the picture using the image label as the prediction target.

### 4.3 Translation or language-specific fine-tuning

A central question for us was whether we should fine-tune three different models, one for each challenge language, or train a single model and translate the resulting predictions into the other two target languages. The first attempt had the potential to yield good results, especially in English, since LLaMA, which provided the base weights for LLaVA-Med, was only peripherally trained on Chinese and Spanish. On the other hand, the quality of training data was the highest in Chinese, since this was the language the data was sourced in, and translations were automatic and, in some places, inaccurate. This could lead to the model learning inaccurate terms, reflecting poorly in the test set because it was translated manually. In this case, translating the generated answers would be the preferred option. When translating with Mixtral, we prompted the model to generate an accurate translation of a Chinese forum post with medical content in Spanish and English, respectively. Figure 3 shows these prompts. To achieve higher-quality translations and to ensure the model would adhere to our instructions, we constructed 3 few-shot examples containing fictional example sentences that were similar in style but not originally contained in the training data. Finally, we post-processed with simple regex expressions to exclude additional remarks Mixtral often made, which were not part of the translation.

## 5 Results

Our best results were achieved by training LLaVA-med exclusively on Chinese challenge data, for only a single epoch, as more epochs to decrease performance. The learning rate was $2e-5$, with an overall batch size of $4$ and $16$ gradient accumulation steps. We did not make use of validation data in fine-tuning for our final submission. The resulting predictions were then translated into Spanish and English using Mixtral-8x7b-instruct. This method achieved a score of 7.05 BLEU for Chinese, 2.66 for English, and 1.36 for Spanish. (see Table 1). These represent the highest scores

Figure 3: The system prompts used to generate translations from Chinese into English and Spanish



Figure 4: The left-hand side shows the dataset collection process. It consists of chat interactions between Chinese dermatologists and patients. Each patient inquiry contains a text and multiple photos of their skin disease. We train a VLM on the original Chinese examples. For the application of this model in other languages, we translate the model answers from Chinese to English and Spanish using an LLM.

achieved in the challenge for Chinese and Spanish. As mentioned in the previous section, this represents a fairly simple approach compared to other experiments we performed, which is visualised in its entirety in Figure 4.

# 6 Discussion

In addition to the main result described above, we performed several additional experiments with differing approaches, which in most cases led to worse performance than in the version we submitted. Table 1 gives an overview of these results. The following section gives some reasons for why additional training might have harmed model performance in this case and why a simple approach ended up achieving the highest scores.

## 6.1 Analysis of fine-tuning methods

It becomes apparent that additional datasets that were not originally part of the challenge worsen results by 0.61 points in the case of English and 1.08 points in the case of Spanish. Following up with fine-tuning on challenge data improved the

score again slightly, but it does not come close to reaching the scores of training exclusively on challenge data. It is possible that this was due to the incompatibility of datasets, meaning that diagnoses contained in challenge data were not represented in Fitzpatrick or Dermnet. Additionally, challenge data often contained more complex tasks than correctly identifying what could be seen in the image, e.g., answering questions about potential treatments. Finally, the particular writing style of many entries in the challenge data and differing translations may also have played a role.

Mixtral-8x7b-instruct seems to outperform Opus as a translation option despite Opus being a group of models designed specifically for translation between set language pairs. One constraint expected to lead to Opus's poorer performance was that this model family only contains a model for Chinese to English and English to Spanish translations, but none for Chinese to Spanish, thus necessitating a translation first to English and then to Spanish. However, our results show this is not the case since the Spanish Opus translation outperforms the En-

Table 1: This table contains the various results we achieved with different fine-tuning methods. Datasets used: 1. M3G: original challenge data 2. DN: Dermnet 3. FP: Fitzpatrick17k

| ID | Datasets | Training Language | Translation Method | Score (ZH / EN / ES) |
|----|----------|-------------------|--------------------|----------------------|
| 1 | M3G | Chinese | Mixtral | **7.05** / **2.66** / 1.36 |
| 2 | M3G | English | - | - / 2.05 / - |
| 3 | M3G | Spanish | - | - / **1.58** / - |
| 4 | M3G | Chinese | Opus | 7.05 / 0.60 / 0.99 |
| 5 | FP | English | - | - / 0.47 / - |
| 6 | FP + M3G | English | - | - / 0.94 / - |
| 7 | DN | English | - | - / 0.57 / - |
| 8 | DN + M3G | English | - | - / 1.44 / - |
| 9 | DN + FP | English | - | - / 0.77 / - |
| 10 | DN + FP + M3G | English | - | - / 1.41 / - |

glish translation.

The answer is unclear regarding whether one should train designated models for each language or translate results using a translator model. Translation outperforms designated training in the case of English but not in the case of Spanish.

In addition to the variables discussed up to this point, we trained LLaVA-Med by using multiple input images instead of a single one, which also worsened the results. Compared to our best result, training exclusively on Chinese challenge data, multi-image training only scored 0.63 BLEU. In general, increasing the number of input images during training seems to decrease LLaVA's performance. This might be because LLaVA models are usually only pre-trained using single images, so the model can not properly handle multi-image input. Another potential reason performance decreases with more input is the way multimodality is implemented in LLaVA models: image features are projected into the embedding space of the language model, thus effectively increasing input length. Longer contexts have been shown to decrease language model performance. (Levy et al., 2024)

## 6.2 Error analysis

Looking at the model predictions, it becomes clear that there are still several issues with its performance. Firstly, as the model mirrors challenge training data, it tends to gravitate towards very short and concise answers, simply stating a presumed diagnosis, see for example Figure 5.

These answers could, in some cases, agree with expert opinions but did not contain the same amount of additional information and did not directly answer the original question. There were

**Example**

**Prompt:** Male, 16 years old. Got Pustule for 10 years at the hands and feet. Previously diagnosed as chronic eczema. Some of our colleagues suspected it is scabies. And hence was treated once as such, with no success. Later on, with a follow up on the medical history, it was revealed the lesion at the foot was the result of a scald. The lesions in the hands and feet have been there for 10 years. There was some relief after the use of anti-biotics. Now it mainly happens on the palms and side of the hands and feet. Once again, does it look like eczema? How to account for these small pustules? Picture uploaded before: The latest status:

**LLaVA-med**: Impetigo.

**Reference 1**: Impetigo.

**Reference 2**: Tinea pedis combined with fungal rash. Suggest to run a test on fungal infection.

**Reference 3**: What appear on the hand is eczema, impetigo. Exfoliative keratolysis at the feet

Figure 5: Example prediction by LLaVA-Med, translated into English using Mixtral.

also instances in which the model predicted a different but similar illness to the reference, often with a different descriptive adjective, see Figure 6.

Regarding translations, there were several differences between Mixtral and Opus. Opus tended to mistranslate medical terms. A possible reason is the added context that an LLM receives through the system prompt, clarifying that the input text has medical content. Additionally, it often failed to identify two connected Chinese characters as a single word, thus translating single characters liter-

Table 2: This table contains some examples that illustrate the differences in translation quality between Mixtral and Opus.

| LLaVA-med prediction | Mixtral translation | Opus translation | Example ground truth |
|---|---|---|---|
| 慢性湿疹 | Chronic eczema. | Chronic wet rashes. | Could be chronic eczema. |
| 季节性蜕皮症 | Seasonal skin exfoliation. | Seasonal skinness. | Exfoliative keratolysis. |
| 虫咬皮炎吧 | Must be an insect bite dermatitis. | Worm bites skin inflammation. | Seborrheic dermatitis. |
| 考虑毛囊炎 | Considering folliculitis. | Think of cyborgitis. | [...] folliculitis or summer dermatitis [...] |

**Example**

**Prompt**: See picture of the disease -maculopapular rash. The patient is male, 16 years old. Impression: appendicitis. Got influenza when in the hospital. After injection of antondine by the night shift physician, the skin in the affected area becomes itchy. They come to me for consultation the next day.

**LLaVA-med**: Skin rash disease?

**Reference 1**: Should be Measles type drug rash.

**Reference 2**: Drug rash.

Figure 6: Example prediction by LLaVA-Med, translated into English using Mixtral.

ally, missing the actual meaning. Table 2 contains a typical example for this: The Chinese word for 'eczema' consists of the characters for 'wet' and 'rash'. Opus interpreted these as distinct characters instead of a single word and thus reached an inaccurate translation. Opus also tended to choose general terms for some words instead of the correct scientific term. (E.g., simply 'inflammation' instead of 'dermatitis', see also Table 2) On the other hand, Mixtral achieved a relatively high quality of translations, given that it is neither officially trained on Chinese nor specifically biomedical data.

## 7 Limitations

The model we submitted has some limitations, excluding it from clinical use in its current state. Most importantly, even though our results scored the highest in two languages, the overall scores were very low. Significantly higher diagnosis accuracy must be achieved to make it useful in a clinical

setting.

Secondly, due to the nature of the training data, the model often mimics the writing style of the forum posts it was trained on, leading to fewer professional expressions than expected in a clinical setting.

Similarly, since training data was obtained from Chinese sources containing frequent suggestions for using Traditional Chinese Medicine, the model made similar recommendations in some cases. This might not meet the standards of care in other countries. Thus, regional differences in care methods have to be considered when training similar models intended for clinical use in the future.

## 8 Conclusion

We present our submission to the Multilingual & Multimodal Medical Answer Generation task of the MediQA 2024 challenge. Our results compare well with other submitted approaches, but their quality is still insufficient for clinical use. This was partly because our method could not overcome hurdles presented by the challenge, such as short target predictions, translation issues, and regional differences in care methods. VLMs with better analytic capabilities in the medical domain must be created to achieve scores high enough for real-world applications. Nonetheless, the increased availability of telemedical records and the increasing availability of data from a variety of countries also presents an opportunity for medical LVM research.

# References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. Flamingo: a visual language model for few-shot learning. In *Advances in Neural Information Processing Systems*.

Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of BLOOM. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 157–170, Tampere, Finland. European Association for Machine Translation.

Asma Ben Abacha and Dina Demner-Fushman. 2019. On the summarization of consumer health questions. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2228–2234, Florence, Italy. Association for Computational Linguistics.

Asma Ben Abacha, Chaitanya Shivade, and Dina Demner-Fushman. 2019. Overview of the MEDIQA 2019 shared task on textual inference, question entailment and question answering. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 370–379, Florence, Italy. Association for Computational Linguistics.

Zeming Chen, Alejandro Hernández-Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023. Meditron-70b: Scaling medical pretraining for large language models. *Preprint*, arXiv:2311.16079.

Zhihong Chen, Maya Varma, Jean-Benoit Delbrouck, Magdalini Paschali, Louis Blankemeier, Dave Van Veen, Jeya Maria Jose Valanarasu, Alaa Youssef, Joseph Paul Cohen, Eduardo Pontes Reis, et al. 2024. Chexagent: Towards a foundation model for chest x-ray interpretation. *arXiv preprint arXiv:2401.12208*.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Amin Dada, Marie Bauer, Amanda Butler Contreras, Osman Alperen Koraş, Constantin Marc Seibold, Kaleb E Smith, and Jens Kleesiek. 2024. Clue: A clinical language understanding evaluation for llms. *Preprint*, arXiv:2404.04067.

Alison Ehrlich, James Kostecki, and Helen Olkaba. 2017. Trends in dermatology practices and the implications for the workforce. *Journal of the American Academy of Dermatology*, 77(4):746–752.

Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. 2017. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118.

Alexander L Fogel and Joseph C Kvedar. 2018. Artificial intelligence powers digital medicine. *NPJ digital medicine*, 1(1):5.

Michel Galley, Chris Brockett, Alessandro Sordoni, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 445–450, Beijing, China. Association for Computational Linguistics.

Matthew Groh, Caleb Harris, Roxana Daneshjou, Omar Badri, and Arash Koochek. 2022. Towards transparency in dermatology image datasets with skin tone annotations by experts, crowds, and an algorithm. *arXiv preprint arXiv:2207.02942*.

Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. 2020. Pathvqa: 30000+ questions for medical visual question answering. *arXiv preprint arXiv:2003.10286*.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and Laurent Sifre. 2022. An empirical analysis of compute-optimal large language model training. In *Advances in Neural Information Processing Systems*, volume 35, pages 30016–30030. Curran Associates, Inc.

Zhen Hong, Nian Li, Dajiang Li, Junhua Li, Bing Li, Weixi Xiong, Lu Lu, Weimin Li, and Dong Zhou. 2020. Telemedicine during the covid-19 pandemic: Experiences from western china. *J Med Internet Res*, 22(5):e19577.

Xinyue Hu, Lin Gu, Qiyuan An, Mengliang Zhang, Liangchen Liu, Kazuma Kobayashi, Tatsuya Harada, Ronald M. Summers, and Yingying Zhu. 2023. Expert knowledge-aware image difference graph representation learning for difference-aware medical visual question answering. KDD '23, page 4156–4165, New York, NY, USA. Association for Computing Machinery.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. Preprint, arXiv:2401.04088.

Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, and Zhaopeng Tu. 2023. Is chatgpt a good translator? a preliminary study. arXiv preprint arXiv:2301.08745, 1(10).

Alexa B Kimball. 2003. Dermatology: a unique case of specialty workforce economics. Journal of the American Academy of Dermatology, 48(2):265–270.

Alexa Boer Kimball and Jack S Resneck Jr. 2008. The us dermatology workforce: a specialty remains in shortage. Journal of the American Academy of Dermatology, 59(5):741–745.

Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. Preprint, arXiv:2402.10373.

Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. 2018. A dataset of clinically generated visual questions and answers about radiology images. Scientific data, 5(1):1–10.

Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. Same task, more tokens: the impact of input length on the reasoning performance of large language models. Preprint, arXiv:2402.14848.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023a. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. arXiv preprint arXiv:2306.00890.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. BLIP-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 19730–19742. PMLR.

Zhihong Lin, Donghao Zhang, Qingyi Tao, Danli Shi, Gholamreza Haffari, Qi Wu, Mingguang He, and Zongyuan Ge. 2023. Medical visual question answering: A survey. Artificial Intelligence in Medicine, 143:102611.

Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. 2021. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI), pages 1650–1654.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning.

Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zakka, Eduardo Pontes Reis, and Pranav Rajpurkar. 2023. Med-flamingo: a multimodal medical few-shot learner. In Machine Learning for Health (ML4H), pages 353–367. PMLR.

Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. arXiv preprint arXiv:2304.03277.

Xuan Song, Xinyan Liu, and Chunting Wang. 2020. The role of telemedicine during the covid-19 epidemic in china—experience from shandong province.

Omkar Thawkar, Abdelrahman Shaker, Sahal Shaji Mullappilly, Hisham Cholakkal, Rao Muhammad Anwer, Salman Khan, Jorma Laaksonen, and Fahad Shahbaz Khan. 2023. Xraygpt: Chest radiographs summarization using medical vision-language models. arXiv preprint arXiv:2306.07971.

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.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.

Wen wai Yim, Asma Ben Abacha, Velvin Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen, and Martin Krallinger. 2024a. Overview of the mediqa-m3g 2024 shared task on multilingual and multimodal medical answer generation. In Proceedings of the 6th Clinical Natural Language Processing Workshop, Mexico City, Mexico. Association for Computational Linguistics.

Wen wai Yim, Velvin Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Morgan Waller and Chad Stotler. 2018. Telemedicine: a primer. *Current allergy and asthma reports*, 18:1–9.

Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Towards generalist foundation model for radiology. *arXiv preprint arXiv:2308.02463*.

Qianqian Xie, Qingyu Chen, Aokun Chen, Cheng Peng, Yan Hu, Fongci Lin, Xueqing Peng, Jimin Huang, Jeffrey Zhang, Vipina Keloth, Huan He, Lucila Ohno-Machido, Yonghui Wu, Hua Xu, and Jiang Bian. 2024. Me llama: Foundation large language models for medical applications. *Preprint*, arXiv:2402.12749.

Binwei Yao, Ming Jiang, Diyi Yang, and Junjie Hu. 2024. Benchmarking llm-based machine translation on cultural awareness. *Preprint*, arXiv:2305.14328.

# A   Code availability

The code used to perform all experiments listed in this paper is available in this repository. [2]

---

[2] https://github.com/Shiniri/MediQA-M3G-Submission

# NEUI at MEDIQA-M3G 2024: Medical VQA through consensus

**R. Omar Chávez García**[†][*] and **Oscar Lithgow-Serrano**[†][*]

[†] Dalle Molle Institute for Artificial Intelligence Research (IDSIA), USI-SUPSI, Switzerland

`{omar.chavez, oscarwilliam.lithgow}@idsia.ch`

## Abstract

This document describes our solution to the MEDIQA-M3G: Multilingual & Multimodal Medical Answer Generation. To build our solution, we leveraged two pre-trained models, a Visual Language Model (VLM) and a Large Language Model (LLM). We fine-tuned both models using the MEDIQA-M3G and MEDIQA-CORR training datasets, respectively. In the first stage, the VLM provides singular responses for each pair of image & text inputs in a case. In the second stage, the LLM consolidates the VLM responses using it as context among the original text input. By changing the original English case content field in the context component of the second stage to the one in Spanish, we adapt the pipeline to generate submissions in English and Spanish. We performed an ablation study to explore the impact of the different models' capabilities, such as multimodality and reasoning, on the MEDIQA-M3G task. Our approach favored privacy and feasibility by adopting open-source and self-hosted small models and ranked 4th in English and 2nd in Spanish.

## 1 Introduction

Recent visual iterations of Large Language Models (LMM) explore a central concept that deals with multimodal inputs, known as visual instruction tuning. These studies result in sizable Visual Language Models (VLM) such as VisualBERT (Li et al., 2019), LLaVA (Haotian et al., 2023), MiniGPT-4 (Zhu et al., 2023) that demonstrate impressive results on natural instruction-following and visual reasoning capabilities.

The need for multimodal models is particularly pronounced in the medical domain. Medical Visual Question Answering (VQA) can assist in clinical decision-making, provide reliable and user-friendly answers to free-form questions, serve as a diagnostic tool, or act as a knowledgeable assistant,

potentially alleviating the burden on the healthcare system and enhancing the efficiency of medical professionals. A mature and comprehensive medical VQA system could directly review patients' images and answer any questions, especially relevant when medical professionals may not be immediately available.

The MediQA-M3G task, which focuses on clinical dermatology multimodal query response generation, exemplifies this need. This task aims to automatically generate clinical responses given textual clinical history, user-generated images, and queries. The common challenges of VQA are amplified in the medical domain, where highly specialized knowledge must be leveraged in coordination with specific visual features from the images. This task is further complicated by the fact that the query, content, and images are provided by patients, which implies a highly heterogeneous text style, varying levels of description details, and, in the case of images, highly variable light, focus, zoom, and quality conditions.

We utilized a compact (1.86B parameters) Visual Language Model (VLM) named Moondream (Moondream AI, 2024) to evaluate the performance of small Language-Image Models (LIMs) on the M3G multimodal task. Moondream is built upon a Sigmoid loss for Language-Image Pre-training (SigLIP) and the Phi-1.5 language model. We fine-tuned the VLM using the provided training data, extending each case title and description to all the provided images.

The output of VLM might contain redundancies and short answers that deviate from the provided context in the query. We implemented a post-processing step of the VLM output to address this issue by constructing a new query for a fine-tuned BioMistral LLM. This step relies on the idea that we already have the context to improve the VLM answer. The context consists of the original query title and content from the test dataset cases and the

---

[*]All authors contributed equally.

VLM response, which we refer to as image analysis. Along with the context, we used the general query *"What is the disease present in the photo? What is the treatment?"* We use the same pipeline for both English and Spanish submission entries.

## 2 Task definition

The *MEDIQA-M3G: Multilingual & Multimodal Medical Answer Generation* task focuses on the problem of clinical dermatology multimodal query response generation, a first of its kind, aiming to automatically generate clinical responses given textual clinical history, user-generated images, and queries (wai Yim et al., 2024a). This shared task is motivated and very in line with the rapid development of telecommunication technologies and the increased demand for remote clinical diagnosis and treatment. Unlike previous works focusing only on text or specific types of images, this task incorporates text and one or more images. Participants were given textual inputs, including clinical history and a query, along with associated images, and they were expected to generate a relevant textual response. The training data for this task was translated and adapted from Chinese datasets, and participants could opt to work in Chinese (Simplified), English, or Spanish for the test set (wai Yim et al., 2024b).

## 3 Related work

### 3.1 Large Language Models (LLM)

Integrating generative large language models (LLMs) has been pivotal in medical question-answering systems. Recent advancements have seen the adaptation of generalist LLMs like GPT-4 and Gemini to more specialized domains. However, the proprietary nature of such models limits their accessibility to research. This challenge has been addressed by the open-source movement, with models such as Llama 2 (Touvron et al.), Vicuna (Chiang et al.), and Mistral (Jiang et al.) providing a foundation for further innovations in medical LLMs. Multiple open-source LLMs based on decoder-only architecture have recently been developed for the medical domain, e.g., BioGPT (Luo et al.) and PMC-LLaMA (Wu et al.). Two notable recent contributions in this space are MediTron (Chen et al.) and BioMistral (Labrak et al.). MediTron, leveraging Llama-2, has been pre-trained on a vast corpus of medical literature to offer medical reasoning capabilities. In parallel, BioMistral

adapts the Mistral model to the biomedical domain, showing the potential of merging techniques (Yu et al.) on pre-trained models to enhance performance and out-of-domain generalization. In particular, BioMistral, through techniques like DARE, has demonstrated improved robustness in multi-lingual settings, a key factor in real-case global medical applications.

The massive increase in the size of large language models and, by extension, visual language models to hundreds of billions of parameters has unlocked various emerging capabilities that have redefined the landscape of natural language processing and a plethora of downstream tasks. A common challenge remains whether such emergent abilities can be achieved at a smaller scale using strategic choices for training, e.g., data selection. Proposals such as the Phi family models aim to answer this question by training small language models (SLMs) that achieve performance on par with models of much larger (yet still far from the frontier models) (Javaheripi and Bubeck, 2024). Their success relies upon training data quality and the scalability of their smaller models.

### 3.2 Multimodality

The recent progress of multimodal models in the medical domain is highlighted by the progress in large vision language models such as Flamingo (Alayrac et al.), GPT-4V, and Gemini (Gemini Team et al.), which have demonstrated remarkable capabilities in executing instructions, engaging in dialogues, and managing image-based tasks. Such advancements show the potential of fusing vision encoders with large language models (LLMs) to create systems that can interpret and respond to complex queries involving textual and visual inputs. However, increased hardware demands, longer test times, slower inference speeds, and privacy concerns when used as cloud services are challenges to their use in real-case scenarios, especially for the case of medical applications.

**End-to-end Vision-Language Pre-training.** End-to-end vision-language pre-training (VLP) has been used to develop multimodal foundation models that excel in various vision-and-language tasks. Despite the effectiveness of these models, thanks to the evolution of architectures, learning objectives, and strategies such as contrastive learning and image-text matching, their use is hindered by requiring substantial computational resources

for end-to-end training in large image-text pair datasets. Another limitation is the lack of leverage in existing unimodal pre-trained models. (Faria et al.; Lin et al.)

**Modular Vision-Language Pre-training.** In contrast, this approach involves modular VLP methods that utilize off-the-shelf pre-trained models, keeping them frozen during VLP. This includes techniques that freeze the image encoder, leveraging pre-trained models like CLIP (Radford et al.), and methods that freeze the language model to harness the knowledge from LLMs for vision-to-language tasks. A challenge in this approach is aligning visual features with the text space. BLIP-2 (Li et al.) is a successful recent approach that efficiently uses frozen image encoders and LLMs for various vision-language tasks with reduced computational costs.

**Multimodal Instruction-following Agents.** Instruction tuning has been crucial in reducing complexity and costs by training the model to handle various tasks represented by different instructions, thus eliminating the need for separate models for each application. Common architectures for instruction-following Large Language Models (LLMs) include a pre-trained visual backbone, a pre-trained LLM, and a vision-language cross-modal connector. Notable recent implementations include BLIP-2 (Li et al.) and LLaVA (Liu et al., 2023b,a) models. These represent significant steps in leveraging pre-trained models and visual instruction-tuning to enhance the capabilities of multimodal systems. The introduction of LLaVA-Phi (Zhu et al.) further exemplifies the trend toward creating efficient and compact models capable of delivering high performance in real-time applications. These developments point to AI systems' growing capabilities in understanding and acting upon instructions involving both visual and textual information.

**Medical Visual Question Answering.** Medical VQA can potentially transform clinical decision-making and patient engagement (Lin et al.). The unique challenges of the medical domain, such as privacy requirements, the need for expert annotation, and the integration of medical knowledge bases, are part of the complexity of developing effective Medical VQA systems. Dataset quality and diversity are among the most impactful limitations that must be addressed to advance the field. Although the LLMs and LMMs are adapted to the medical domain and already trained for instruction-following, it is often observed that their zero-shot and few-shot performance can be further enhanced by performing a complementary, focused SFT stage on specific tasks. Notably, task-specific models trained on carefully curated datasets have frequently outperformed generalist models of similar size, especially in highly specialized domains such as medicine.

## 4 Methodology

We utilized a compact (1.83B parameters) Visual Language Model (VLM) named Moondream (Moondream AI, 2024) to evaluate the performance of small Language-Image Models (LIMs) on the M3G multimodal task. Moondream is built upon a Sigmoid loss for Language-Image Pre-training (SigLIP) (Beyer et al., 2022) and the Phi-1.5 language model, a Transformer with 1.3B parameters (Li et al., 2023; Microsoft Research, 2023). In such a setup, a contrastively pre-trained model provides significantly more useful tokens than one classification pre-trained model (Zhai et al., 2023). Figure 1 shows the schematic of the proposed solution involving the VLM and the BioMistral-7B-DARE (Labrak et al.) LLM as a specialized stage for final response consolidation.

### 4.1 Training

**Fine-tuning VLM.** We fine-tuned the VLM using the whole provided training data, extending each case title and description to all the provided images (see Table 1). We employed the flash attention algorithm to mitigate memory issues during training and inference. Our hardware setup was limited to a single NVIDIA RTX 3090 GPU for fine-tuning and inference.

The motivation behind this training dataset is to increase the number of training samples, given the reduced number of clinical cases in the provided training data. The caveat of this approach is that although we consider each augmented sample as valid, there might be responses that overlap, complement, or contradict a valid clinical response.

**Fine-tuning LLM.** Our team, having participated in the MEDIQA-CORR (Ben Abacha et al., 2024a) task, leveraged the LLM fine-tuned for that task. Specifically, we instruction-tuned the BioMistral-7B-DARE on the MEDIQA-CORR dataset (Ben Abacha et al., 2024b). For this, we

Figure 1: Overview of the proposed solution[1]. The contrastively pre-trained SigLIP vision model encodes the image into visual tokens individually. These visual tokens are passed along with a query to the Phi 1.5 LLM, producing responses for individual images. A consolidation response stage is performed via the fine-tuned Biomistral LLM using the VLM responses and context from the original query.

mapped the labeled dataset into three types of instructions: classifying if the statement had an error or not, detecting the culprit sentence, and correcting a given erroneous sentence to ensure consistency with the rest of the clinical text. The Supervised Fine-Tuning (SFT) was performed using the parameter-efficient method LORA on an NVIDIA A100-80G for four epochs. Without further training, we then used this MEDIQA-CORR fine-tuned model in the M3G task.

## 4.2 Inference strategies

**Strategy-1: Direct inference with VLM.** We constructed the output by performing inference on each image of each case in the test dataset. This step means that for one case, we request the fine-tuned VLM with our query and each of the case's images. Finally, all VLM responses for a case, as many as images in the test case, were concatenated as the final response (see Table 2:2). The results of this strategy outperform the baseline obtained using the non-fine-tuned VLM (see Table 2:1).

**Strategy-2: Two-stage inference (VLM + LLM)** The output of the previous approach might contain redundancies and short answers that deviate from the provided context in the query. To address this issue and to harness knowledge from a bigger specialized model, we implemented a two-stage strategy that augmented the previously described *Direct inference* strategy with a post-processing step. This step relies on the idea of leveraging the arguably better reasoning capabilities of a bigger specialized

LLM to better harness the provided case information, i.e., query title and content, along with the VLM answers to generate a final response. Specifically, we requested the LLM with a prompt consisting of the query: "What is the disease present in the photo? What is the treatment?"; the context: dataset query title and content; and the image analysis: list of VLM responses (see Table 2:3). Table 3 shows examples of the composite input of this step and the resulting consolidated response.

Regarding multilingualism, Strategy-1 was built considering only one language data stream, English. The VLM was fine-tuned using only the English queries, content, and responses. However, as the LLM we employed in Strategy-2 has multilingual capabilities (see sec. 3.1), we also applied the post-processing step of Strategy-2 to the Spanish version of the cases. We provided a prompt with the query and context in Spanish but with the English image analysis. We added additional prompt instructions to the model to request Spanish responses exclusively. As a result of this change, we could provide output for the Spanish version of the task (see Table 2:3).

## 5 Results and analysis

**Results during competition.** From the official results during the competition (Table 2 ids: 1-3), we observe that by fine-tuning the VLM (id: 2), we obtained a significant improvement, with a `deltableu` of 0.595, which is more than a two-fold enhancement over the baseline non-fine-tuned version (id: 1) that held a `deltableu` of 0.231. Furthermore, the implementation of Strategy-2 (id: 3) marked a substantial leap, exhibiting a quadruple

---

[1]MEDIQA-M3G dataset contains images of medical conditions that may be sensitive and/or graphic in nature.

| Original sample (single language) | Training sample |
|---|---|
| **case:** `ENC00018` | **sample:** `ENC00018_image1_response1` |
| (image1, image2) | (image1) |
| **from:** human; **value:** (<u>title</u>) `View image` (<u>content</u>) `Female, 19 years old, has had a hard lump in her ear for three months, as hard as a wooden board, with no sense of fluctuation. After incision, a white dense substance was found. What kind of cyst could this be?` | **from:** human; **value:** (<u>title</u>) `View image` (<u>content</u>) `Female, 19 years old, has had a hard lump in her ear for three months, as hard as a wooden board, with no sense of fluctuation. After incision, a white dense substance was found. What kind of cyst could this be?` (<u>augmented query</u>) `What is the disease in the photo? What is the best treatment?` |
| **from:** response 1; **value:** `Erythema annulare centrifugum? Use licorice decoction with corticosteroid ointment.` | **from:** agent ; **value:** `Erythema annulare centrifugum? Use licorice decoction with corticosteroid ointment.` |
| **from:** response 2; **value:** `I think it still looks like urticaria, continue with the anti-allergy treatment.` | |
| **from:** response 3; **value:** `I think the likelihood of urticaria is the highest, but the skin lesions at the root of the thigh are hard to explain, so erythema annulare cannot be ruled out either...` | |

Table 1: Training example used for fine-tuning the VLM. We augment the training query (represented by the title and content case) with the standard query from the challenge description. We generate a training sample per each image and response combination. Hence, each case in the training dataset will generate $I \times R$ training samples, where $I$ is the number of images in the case, and $R$ is the number of responses for the selected language.

increase in performance for the English language tasks, as indicated by a `deltableu` of 2.133 compared to the 0.595 achieved by Strategy-1 (id: 2). When applied to Spanish, Strategy-2 (id: 3) show a significant drop but still got a competitive performance with a `deltableu` of 0.974. Our best runs (id: 3) were placed at the 4th and 2nd positions for English and Spanish, respectively.

### 5.1 Ablation study

We conducted an ablation study to assess the impact of various components in our best strategy (Strategy-2). We can isolate and understand their contributions to the strategy's effectiveness by systematically removing or altering specific model elements. Our analysis focuses on three primary objectives: investigating the Unimodal Bias phenomenon, assessing the extra reasoning capacity contribution of the Large Language Model (LLM), and evaluating the impact of training the LLM on a specialized dataset for error detection and correction in clinical notes.

**Investigating the Unimodal Bias Phenomenon.** To explore the Unimodal Bias and the impact of incorporating visual modality, we performed experiments 4 and 5 (see Table 2). We follow the same pipeline as in Fig. 1 without using the input images for the unimodal experiments. Thus, the VLM only sees the test case's title and content text inputs as prompts, i.e., we remove the references to "image" from the prompt. In experiment id:4 our strategy involved employing the Visual Language Model (VLM) without providing visual inputs, relying solely on textual content. This setup mimics Strategy-1 but aims to quantify the absence of visual modality. Experiment id:5 followed a similar approach, utilizing both the VLM and LLM without visual inputs, akin to Strategy-2. As seen in Table 2, the results – `deltableu` scores of 1.418 and 0.968 for English and Spanish, respectively in id: 4 and 0.328 for English in id: 5– indicate the positive impact of using both modalities in this task. The decrements in `BERTscore` and `deltableu` metrics suggest that relevant information in the encoded & tokenized image input is helping, along with textual case input, to determine the test case queries.

**Assessing the Extra Reasoning Capacity of the LLM.** The comparison of Strategy-2's performance under varying conditions—specifically when the LLM is provided with both the case context and VLM responses versus when it only receives the VLM responses for summarization—sheds light on the LLM's reasoning ability. Experiment id:6 explores this, allowing us to distill the LLM's added value in synthesizing and

452

| | id | Strategy | EN | ES |
|---|----|----------|-----|-----|
| Test | 1 | Moondream | 0.231 | - |
| | 2 | Moondream-FT | 0.595 | - |
| | 3 | Moondream-FT + BioMistral-FT | **2.133** | 0.974 |
| Test_after | 4 | Moondream-FT :: w/o visual | 0.328 | - |
| | 5 | Moondream-FT + BioMistral-FT :: w/o visual | 1.418 | 0.968 |
| | 6 | Moondream-FT + BioMistral-FT :: w/o context | 1.183 | - |
| | 7 | Moondream-FT + BioMistral :: w/o FT-LLM | 1.963 | **1.745** |

Table 2: Official scores (`deltableu`) of the different submitted strategies for English (EN) and Spanish (ES). Stages, Test: during competition, Test_after: after the end of competition. The best scores by language appear in bold.

reasoning over the provided information. With a `deltablue` score of 1.183 in English, this experiment shows how much the LLM's reasoning capabilities, beyond mere summarization, contribute to generating more correct and contextually aware responses.

**Evaluating the LLM's Training on Error Detection and Correction.** BioMistral LLM utilizes Mistral as its foundation model. It is further pretrained on PubMed Central (a dataset containing citations and abstracts of biomedical literature), making it a top performer in medical question-answering benchmarks in English. Experiment id:7 investigates the relevance of the ability of error detection and correction within clinical notes by exploring the original BioMistral against one fine-tuned on the CORR dataset. This experiment examines the hypothesis that an LLM trained for error detection&correction could better integrate VLM responses with the textual case content, especially in correcting inconsistencies in VLM responses ("diagnostic"). The results from this experiment, 1.963 for English and 1.745 for Spanish, demonstrate the potential benefits of specialized fine-tuning for tasks out of the LLM's immediate domain expertise, highlighting the enhanced capability for error correction and the generation of coherent and accurate clinical responses.

## 5.2 Discussion

Analyzing the results of the competition phase and of the ablation study (see Figure 2 & Table 2), we observe that when using the non-fine-tuned version of BioMistral, we obtain the smallest drop in performance, a mere 7%. In contrast, a more significant drop in scores, a 33% degradation, was observed when the visual input was removed. Interestingly, the loss was even higher, at 44%, when neglecting the reasoning capabilities of the LLM. This suggests that the analysis and synthesis, i.e., **consensus generation** capacity of the LLM, is a key component of the strategy.

All the ablation experiments in Table 2, except for the experiment id:7, for Spanish, resulted in lower scores. Interestingly, the Spanish version in experiment id:7 scored the highest and surpassed any published run in the leaderboard to the best of our understanding. We hypothesize that by fine-tuning BioMistral on the CORR dataset (which is only in English), we not only steered the LLM towards a narrow set of tasks, specifically clinical error detection and correction but also disrupted the model's capacity to handle other languages due to the monolingual nature of the training set and the prevalent use of English in the pretraining corpora. This leads to an intriguing inquiry: how robust is the multimodal capacity attributable to merging methods like DARE (Yu et al., 2024) when subjected to monolingual posterior fine-tunings? This question warrants further investigation. Furthermore, removing the visual input has almost no impact on the performance in the Spanish version (id:3 vs. id:5). We hypothesize that this is also the result of our BioMistral-FT version's degraded multilingual capacity, which makes it unable to integrate the case context in Spanish with the VLM responses in English.

The competition results and ablation analysis clearly indicate that in Strategy-2, all components work collaboratively for the better. Even the fine-tuned version of BioMistral, which had the lowest impact, positively contributed to the final score. This shows the key role of consensus generation.

Figure 2: Deltableu scores of our submissions. The strategy used is represented on the X-axis. Scores for English are in blue and for Spanish in orange. The shaded area represents the submissions in the after-test stage.

By integrating multiple independent responses from the multimodal model and re-analyzing the case context, the strategy generates a revised final response, which is more contextually accurate.

LLMs, and by extension VLMs, differ significantly from prior deep learning methods regarding their scale, capabilities, and broad potential impact. For instance, these models are trained on massive datasets and use billions of parameters, resulting in considerable complexity. Models of this scale require significant hardware resources for training, fine-tuning, and, some, even inference. Privacy in medical applications of LLMs is paramount, and the possibility of training and testing this kind of model on-site is critical. Relying on third-party hardware providers to store or process medical data becomes a privacy risk (Khullar et al., 2024; Meskó and Topol, 2023). Our proposed pipeline considers this requirement when dealing with the delicate na-

ture of the images in the training dataset. To do so, we explored the use of compact Visual Language Models and their performance in the M3G task. Our results provided a promising perspective, even with the limited data we utilized for training and the conservative score we obtained in the challenge.

## 6 Limitations

Our proposed approach holds significant promise for the VQA problem in clinical dermatology , but several limitations associated with deploying VLMs and LLMs in real-world medical settings necessitate careful consideration.

We are optimistic about the potential of our 2-step method, which is designed to consolidate multiple responses from image-text query pair analysis into a single, consensus-based solution. This approach allows us to utilize simpler VLMs, initially intended for single-image scenarios, in settings

454

with an unknown number of multiple images. However, we acknowledge that this flexibility comes with the cost of posing as many queries to the VLM as there are images in a single case. This could increase computational costs, potentially making the solution computationally impractical for real-world deployment.

Another limitation is that VLMs, which were aligned with domain-specific images and texts during pretraining, are observed better to leverage domain-specific training examples during the instruction tuning phase. However, in our setting, the VLM lacks alignment for visual-medical texts and relies solely on fine-tuning to generate the most appropriate answers. This lack of alignment makes the VLM more demanding for instruction tuning data.

Another crucial problem with VLMs based on the pre-trained vision encoder is resolution. They are trained and also expect to analyze the full image input. However, for some specific cases, and even if the input is big enough, the focus of the query relies on certain zones of the image –in extreme cases, these zones are tiny compared to the image resolution. The M3G dataset showcases this very problem. Most dermatology-related images in either training, validation, or test datasets contain wide shots of the patient's limbs, and only a tiny region of the image provides valuable visual cues. We envisioned exploring techniques such as Visual Search (Wu and Xie, 2023) and Visual Cropping (Zhang et al., 2024) that can help us tackle this issue without compromising the size of our affordable VLM.

Regarding the LLM component, even if we are within the considered "small" scale, the computational demands of these models are substantial. Operating such models requires significant computational resources, which may not be feasible in all clinical environments. This issue can hinder our proposed solution's scalability and practical deployment in resource-limited settings. Moreover, LLMs are prone to generating "hallucinations" or outputs that may include incorrect or misleading information. This phenomenon is particularly concerning in the medical field, where accuracy is crucial to avoid misdiagnoses or inappropriate treatments. Intrinsic hallucinations, where outputs logically contradict known facts, and extrinsic hallucinations, where outputs cannot be verified, both pose serious risks in clinical applications.

Additionally, data bias and patient privacy are critical. LLMs trained on biased data can perpetuate or amplify these biases, leading to skewed or unfair medical advice. For example, the competition dataset observed a frequent recommendation based on traditional Chinese medicine. Thus, a model trained on this dataset may exhibit a predisposition, favoring a certain type of recommendation, irrespective of local or user-specific preferences. Given the sensitive nature of medical data, ensuring patient privacy while using such models is also paramount. Furthermore, updating these models with new medical knowledge remains a complex and resource-intensive process. This is problematic in the fast-evolving field of medicine, where staying current with the latest research and clinical findings is essential. For example, if a new Adverse Drug Effect is discovered, it is vital to update the models' knowledge promptly.

Finally, it is necessary to be aware that although this shared task is a crucial step towards better understanding and addressing the relevant task of automatically generating clinical responses given the textual clinical history and user-generated images, similar to existing benchmarks and metrics, it does not imply a comprehensive assessment of the performance of the system in real-medical contexts. Metrics such as trustworthiness, helpfulness, explainability, and faithfulness are crucial for clinical applications, and addressing these issues involves not only technical advancements in the architecture and training of LLMs but also close collaboration with medical professionals to ensure the clinical validity and ethical deployment of these technologies.

In conclusion, while our solution shows promise in addressing the MEDIQA-M3G task, limitations must be addressed to make it suitable for clinical use. Further exploration of optimization strategies, evaluation with other metrics, and collaboration with medical professionals are necessary to improve our approach's clinical relevance and effectiveness in real-world healthcare settings.

## 7 Future directions

We plan to incorporate a broader array of medical and health-related datasets into our training regimen to enhance our models' domain-specific accuracy and relevance. Specifically, we aim to utilize the Skin Condition Image Network (SCIN) dataset (Ward et al., 2024) focused on dermatology, including structured and unstructured textual data.

Moreover, we are interested in exploring the potential benefits of integrating data from various clinical specialties into our training process to see how this affects the model's performance and applicability across different medical fields.

We are particularly keen on incorporating retrieval-augmented generation (RAG) strategies related to the challenges of model knowledge updating and mitigating hallucinations. These strategies leverage existing related medical knowledge during the inference phase to enhance the factuality of the generated responses. By doing so, we expect to improve the reliability and accuracy of the model outputs, which is crucial for clinical applications.

Finally, we recognize the importance of interdisciplinary collaboration in developing medical VLMs and LLMs. Therefore, we are already in plans to initiate partnerships with medical professionals who can provide valuable insights, contribute relevant training data, and help define the desired outcomes for these technologies. Their involvement is critical not only in the development phase but also in testing these models in real-world clinical scenarios to ensure they meet both practical clinical needs and high standards of medical care.

# 8 Conclusions

We explored a solution to the clinical dermatology multimodal query response generation task and proposed a pipeline that can be expanded to similar multimodal tasks. We leverage performant pre-trained language models, fine-tuning the small VLM to adapt to the clinical task. We also show how the pipeline adapts to the multilingual complementary problem by relying on the multilingual capabilities of the pre-trained LLM. Participating in this challenge represented a feasibility study and opened several work perspectives for multimodal medical applications.

# References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: A Visual Language Model for Few-Shot Learning. 35:23716–23736.

Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.

Galen Andrew and Jianfeng Gao. 2007. Scalable training of l1-regularized log-linear models. In *International Conference on Machine Learning*.

Asma Ben Abacha, Wen wai Yim, Yujuan Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Yujuan Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. 2022. Big vision.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. MEDITRON-70B: Scaling Medical Pretraining for Large Language Models. *Preprint*, arxiv:2311.16079.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing GPT-4 with 90%* ChatGPT quality.

Ana Cláudia Akemi Matsuki Faria, Felype de Castro Bastos, José Victor Nogueira Alves Silva, Vitor Lopes Fabris, Valeska de Sousa Uchoa, Décio Gonçalves de Aguiar Neto, and Claudio Filipi Goncalves Santos. Visual Question Answering: A Survey on Techniques and Common Trends in Recent Literature. *Preprint*, arxiv:2305.11033.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, and Others. Gemini: A Family of Highly Capable Multimodal Models. *Preprint*, arxiv:2312.11805.

Liu Haotian, Li Chunyuan, Wu Qingyang, and Jae Lee Yong. 2023. Visual instruction tuning. In *NeurIPS23*, volume 6.

Mojan Javaheripi and Sébastien Bubeck. 2024. Phi-2: The surprising power of small language models.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, de las Casas, Diego, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7B. *Preprint*, arxiv:2310.06825.

D. Khullar, X. Wang, and F. Wang. 2024. Large language models in health care: Charting a path toward accurate, explainable, and secure ai. *J GEN INTERN MED*.

Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. BioMistral: A Collection of Open-Source Pretrained Large Language Models for Medical Domains. *Preprint*, arxiv:2402.10373.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. *Preprint*, arxiv:2301.12597.

Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *Preprint*, arXiv:1908.03557.

Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Textbooks are all you need ii: phi-1.5 technical report.

Zhihong Lin, Donghao Zhang, Qingyi Tao, Danli Shi, Gholamreza Haffari, Qi Wu, Mingguang He, and Zongyuan Ge. Medical Visual Question Answering: A Survey. 143:102611.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. *Preprint*, arxiv:2310.03744.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. *Preprint*, arxiv:2304.08485.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie Yan Liu. BioGPT: Generative pre-trained transformer for biomedical text generation and mining. 23(6):1–12.

B. Meskó and E. J. Topol. 2023. The imperative for regulatory oversight of large language models (or generative ai) in healthcare. *NPJ Digit Med*.

Microsoft Research. 2023. The language model phi-1.5.

Moondream AI. 2024. moondream: a computer-vision model can answer real-world questions about images. Visited, March 2024.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *Preprint*, arxiv:2103.00020.

Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015. Yara parser: A fast and accurate dependency parser. *Computing Research Repository*, arXiv:1503.06733. Version 2.

Hugo Touvron, Thibaut Lavril, Xavier Martinet, Marie-anne Lachaux, Timothee Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Guillaume Lample, and Meta Ai. LLaMA : Open and Efficient Foundation Language Models.

Wen wai Yim, Asma Ben Abacha, Yujuan Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen, and Martin Krallinger. 2024a. Overview of the mediqa-m3g 2024 shared task on multilingual and multimodal medical answer generation. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wen wai Yim, Yujuan Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Abbi Ward, Jimmy Li, Julie Wang, Sriram Lakshminarasimhan, Ashley Carrick, Bilson Campana, Jay Hartford, Pradeep Kumar S, Tiya Tiyasirichokchai, Sunny Virmani, Renee Wong, Yossi Matias, Greg S. Corrado, Dale R. Webster, Dawn Siegel, Steven Lin, Justin Ko, Alan Karthikesalingam, Christopher Semturs, and Pooja Rao. 2024. Crowdsourcing dermatology images with google search ads: Creating a real-world skin condition dataset. *Preprint*, arXiv:2402.18545.

Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. PMC-LLaMA: Towards Building Open-source Language Models for Medicine. *Preprint*, arxiv:2304.14454.

Penghao Wu and Saining Xie. 2023. V*: Guided visual search as a core mechanism in multimodal llms. *Preprint*, arXiv:2312.14135.

Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language Models are Super Mario: Absorbing Abilities from Homologous Models as a Free Lunch. *Preprint*, arxiv:2311.03099.

Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. *Preprint*, arXiv:2311.03099.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. *Preprint*, arXiv:2303.15343.

Jiarui Zhang, Mahyar Khayatkhoei, Prateek Chhikara, and Filip Ilievski. 2024. Towards perceiving small visual details in zero-shot visual question answering with multimodal llms. *Preprint*, arXiv:2310.16033.

457

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *Preprint*, arXiv:2304.10592.

Yichen Zhu, Minjie Zhu, Ning Liu, Zhicai Ou, Xiaofeng Mou, and Jian Tang. LLaVA-Phi: Efficient Multi-Modal Assistant with Small Language Model. *Preprint*, arxiv:2401.02330.

## A  Appendix

### A.1  Strategy-2

### A.2  Submission results

| - | Content |
|---|---|
| Case | **ENC00932** |
| Prompt | ### Instruction: Given the CONTEXT and IMAGE ANALYSIS What is the disease and treatment?<br><br>### Input:<br>CONTEXT: Please help take a look Is this onychomycosis? The patient is female, 60 years old. Got onychomycosis. Use Levothyroxine Sodium tablet daily. The nail is thin in the front, but rather thick in the rear part.<br><br>IMAGE ANALYSIS: Onychomycosis, nail detachment? Onychomycosis, nail detachment? Onychomycosis, nail detachment? |
| Response | Onychomycosis, nail detachment? |
| Official responses | - It doesn't look like onychomycosis.<br>- Should be onychomycosis. |
| Case | **ENC00952** |
| Prompt | ### Instruction: Given the CONTEXT and IMAGE ANALYSIS What is the disease and treatment?<br><br>### Input:<br>CONTEXT: Please help see what disease is this. What is the treatment plan for that? Appreciate your help. Hand itchy, dry and cracking for more than 2 months. Small blisters showed up two months ago with unknown cause. Unbearably itchy. Later on, the small blisters fused together to form a patch, cracking and peeling with pain that is unbearable. Prescribed in the hospital oral medicines ( Compound glycyrrhizin capsules, Cimetidine, Cyproheptadine) with no succes. Even with help from many people, still not knowing the cause, treatment unsuccessful. Urgently requesting experts here can help diagnose and treat. Many thanks!<br><br>IMAGE ANALYSIS: It should be eczema, in my personal opinion. It should be eczema, but because your location is in the hospital, it's hard to treat, especially if it's related to infections. It should be eczema, but because your location is in the hospital, it's hard to treat, especially if it's related to infections. |
| Response | Eczema, Intravenous prednisolone |
| Official responses | - Chapped eczema<br>- Chronic eczema. Take special care on contact with specific material or work related. |
| Case | **ENC00925** |
| Prompt | ### Instruction: Given the CONTEXT and IMAGE ANALYSIS What is the disease and treatment?<br><br>### Input:<br>CONTEXT: Papules and lumps at the lower left side of the abdomen for 7 months. The infant patient is 10 months old, and got the disease for 7 months. Got white flat papules and plaque on lower left abdomen. (see pictures) There is no pathological examination yet. Would experts here help diagnose.<br><br>IMAGE ANALYSIS: Consider whether there is a lesion under the skin, the specific direction is difficult to consider. |
| Response | Langerhans cell histiocytosis |
| Official responses | - Could be Verrucous nevus<br>- Verrucous nevus |

Table 3: Examples of prompts given to the LLM for the Strategy-2. Note: In this table, prompts are formatted in the Alpaca style with added line breaks for improved readability.

| run_id | lang | stage | shortdesc | deltaBLEU | BERTScore |
|--------|------|-------|-----------|-----------|-----------|
| 52859 | en | test | Moondream | 0.231 | 0.810 |
| 52872 | en | test | Moondream-FT | 0.595 | 0.851 |
| 52897 | en | test | Moondream-FT + BioMistral-FT | 2.133 | 0.850 |
| 54076 | en | test_after | Moondream-FT :: w/o visual | 0.328 | 0.842 |
| 54086 | en | test_after | Moondream-FT + BioMistral-FT :: w/o visual | 1.418 | 0.846 |
| 54091 | en | test_after | Moondream-FT + BioMistral | 1.963 | 0.829 |
| 54092 | en | test_after | Moondream-FT + BioMistral-FT :: w/o context | 1.183 | 0.860 |
| 52899 | es | test | Moondream-FT + BioMistral-FT | 0.974 | 0.814 |
| 52908 | es | test | Moondream-FT + BioMistral-FT | 0.974 | 0.814 |
| 54085 | es | test_after | Moondream-FT + BioMistral-FT :: w/o visual | 0.968 | 0.810 |
| 54173 | es | test_after | Moondream-FT + BioMistral | 1.745 | 0.809 |

Table 4: All team submissions by language and in chronological order.

# VerbaNexAI at MEDIQA-CORR 2024: Efficacy of GRU with BioWordVec and ClinicalBERT in Error Correction in Clinical Notes

**David Villate[1,*], Laura Tinjaca[2,*], Laura Estrada[3,*], Edwin Puertas[4,+], Juan Pajaro[5,*]**

*Pontificia Universidad Javeriana
+Universidad Tecnologica de Bolívar

[1]juand.villate@javeriana.edu.co, [2]tinjacac.l@javeriana.edu.co
[3]l-estrada@javeriana.edu.co, [4]epuerta@utb.edu.co, [5]juanpajaro@javeriana.edu.co

## Abstract

The automatic identification of medical errors in clinical notes is crucial for improving the quality of healthcare services.LLMs emerge as a powerful artificial intelligence tool for automating this task. However, LLMs present vulnerabilities, high costs, and sometimes a lack of transparency. This article addresses the detection of medical errors through the fine-tuning approach, conducting a comprehensive comparison between various models and exploring in depth the components of the machine learning pipeline. The results obtained with the fine-tuned ClinicalBert and Gated recurrent units (Gru) models show an accuracy of 0.56 and 0.55, respectively. This approach not only mitigates the problems associated with the use of LLMs but also demonstrates how exhaustive iteration in critical phases of the pipeline, especially in feature selection, can facilitate the automation of clinical record analysis.

## 1 Introduction

Large language models (LLMs) demonstrate promise in tackling unseen tasks with notable competencies. However, these models exhibit a fundamental vulnerability. LLMs are costly to train and utilize: their cost has increased 10 to 100-fold since 2015 and must be run on giant compute clusters. The training data used for corporate models is a closely guarded secret that lacks transparency [3]. Additionally, the success of LLMs has led to certain online content being generated entirely by these models, which are susceptible to producing non-factual information. In specialized domains, online information can be unreliable, detrimental, and contain logical inconsistencies that impede the models' reasoning ability. Nevertheless, most prior research on common sense detection has concentrated on the general domain. [1].

In this context, our study focuses on the challenge of identifying common sense errors in clinical notes. Unlike correcting these errors, which

requires a deep understanding and specific knowledge of the medical field, identification is a crucial first step that demands the models' ability to recognize inaccuracies and anomalies in the text. This work explores how advanced natural language processing (NLP) technologies, such as GRU with BioWord-Vec, and especially ClinicalBERT[5], can be useful for analyzing unstructured medical texts. Our methodological approach involves a comprehensive comparative analysis among these models, highlighting their efficiency in identifying errors in clinical notes, underscoring the relevance of adapting model training to the peculiarities of medical data. We seek to demonstrate that, through specialization and fine-tuning of these LLMs models, it is possible to significantly improve their ability to detect erroneous or missing information, crucial for diagnosis and treatment in the clinical setting. This study not only aims to demonstrate the capabilities and limitations of advanced models in specialized medical contexts but also to emphasize the importance of integrating specialized knowledge within LLMs to optimize the reliability and usefulness of clinical notes in medical practice.

This document is described as part of our participation in the Shared Task Medical Error Detection and Correction of the Association for Computational Linguistics [1].

## 2 Related Work

In recent advances in the field of NLP, the ability to identify common sense errors in clinical notes poses a significant challenge and represents an opportunity to improve the quality of healthcare. The relevance of this study lies in exploring the applicability of advanced NLP models for the accurate detection of inaccuracies in medical records. These models constitute a promising advance over the inherent limitations off LLMs especially those arising from the quality and diversity of their training

461

datasets [4]. LLMs often require domain-specific adaptations to effectively handle specialized tasks due to these limitations [4]. Moreover, models like ClinicalBERT have been shown to significantly improve their performance in interpreting clinical language by adapting to specific contexts [11].

NLP has the capacity to derive meaningful insights from unstructured data, specifically in the domain of categorizing incident reports and adverse events. Understanding the nature and reasons behind these incidents is crucial for analyzing adverse events. If NLP can enable the extraction of these insights from larger datasets, it has the potential to enhance learning from adverse events in the healthcare field. [13].

Given the complexity of clinical notes and the necessity for a high degree of precision in their analysis, this study is grounded in the review of previous research that has addressed similar issues in the domain of medical text classification. A relevant study focused on clinical text classification using rule-based features and knowledge-guided convolutional neural networks, leveraging trigger phrases and Unique Medical Concepts (CUIs) from the unified Medical Language System (UMLS) to enhance classification accuracy in class-imbalanced situations [12]. This approach underscores the effectiveness of integrating deep learning with explicit medical knowledge, emphasizing the importance of adapting model training to the specificities of clinical data.

Additionally, a comparative investigation evaluated various deep learning models, including Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), GRU, Long Short-Term Memory (LSTM), Bidirectional Long Short-Term Memory (Bi-LSTM), and a Transformer encoder, in their ability to handle unstructured medical note texts affected by different levels of class imbalance [6]. This analysis provides a critical perspective on the variability in model performance in the face of the unique challenges posed by medical data, highlighting the need for more specialized and adaptive approaches.

These studies and similar efforts outline the current state of using advanced NLP technologies in medical text classification. The present study draws inspiration from these research endeavors to advance understanding of the application of specific NLP models in error identification in clinical notes. In doing so, we aim to contribute to the field by providing valuable insights for future research and practices in this essential domain.

## 3 System Description

In the system description of our study, we address the implementation of an advanced predictive model specifically designed for detecting errors in clinical notes. This model relies on two fundamental pillars of NLP: GRU and the ClinicalBERT architecture [5]. The formulation of our central hypothesis questions the effectiveness of lexical and contextual features obtained through these NLP technologies to identify inaccuracies within clinical texts.

We propose two main methodological strategies. The first strategy implements GRU to extract lexical features, leveraging its ability to process complex temporal dependencies in the data [8]. This aspect is reinforced by the use of BioWordVec, which provides detailed vector representations of medical terms, thereby facilitating the capture of the semantic complexities of clinical language. The adaptability of GRU models to variable-length sequences proves particularly useful for analyzing medical texts, where critical information may be irregularly distributed throughout the document [9].

The second strategy focuses on harnessing ClinicalBERT, a model known for its ability to weigh the relevance of words through attention mechanisms, thereby enabling a deep understanding of the context in which medical terms are embedded. This approach significantly benefits from transfer learning, adapting previously acquired knowledge from extensive medical text corpora to fine-tune the model for our specific task. The synergy between GRU and ClinicalBERT enables a comprehensive analysis of the texts, evaluating not only coherence but also the semantic accuracy of the clinical content [6].

ClinicalBERT exhibits superior performance in identifying significant connections between medical concepts, a validation corroborated by medical experts [6].. This model has surpassed several benchmarks in predicting 30-day hospital readmissions, using discharge summaries and notes from early intensive care units, covering multiple clinically relevant metrics [6]. The attention weights generated by ClinicalBERT facilitate the interpretation of predictions, providing a deeper understanding of the context in which medical terms are embedded. We have released the model parameters and training scripts to encourage further research

in this field. Thanks to its flexible structure, ClinicalBERT can be easily adapted to other predictive tasks with minimal engineering effort, making it ideal for studies requiring detailed analysis of clinical language [6].

Based on the outlined strategies, we configure a detailed Training System as depicted in Figures 1 and 2 of the study. This system unfolds through a sequence of well-defined stages: data ingestion and preliminary cleaning, generation of training instances, extraction of both lexical and contextual features, followed by the classification phase, and finally, model evaluation. This process ensures comprehensive treatment of clinical notes, optimizing error detection through the joint evaluation of long temporal dependencies and detailed contextual analysis.

This approach highlights not only the relevance of incorporating advanced NLP tools in the assessment of clinical texts but also the potential of these technologies to progress towards a higher degree of accuracy and reliability in medical documentation.

## 4  Data Description

The dataset provided by MEDIQA-CORR @ NAACL-ClinicalNLP 2024 [2] offers a comprehensive collection of medical texts, each corresponding to a clinical case report. This dataset stands out for its structured and detailed content, tailored for facilitating the analysis and identification of medical errors. Below are the key features of this dataset:

This dataset represents a valuable tool for research in the field of NLP applied to medicine, especially in tasks related to the identification and correction of errors in clinical texts. The richness and specificity of the data facilitate the development and evaluation of advanced NLP models, as addressed in this study, providing a solid foundation for detailed analysis and improvement of the quality of clinical notes.

## 5  Embeddings

In the process of generating embeddings for our analysis, we applied meticulous preprocessing to the provided data. This preprocessing consisted of a series of essential steps to ensure the quality and uniformity of the text, including correcting encoding errors and normalizing medical terms and units of measurement. This preliminary treatment of the texts is crucial to mitigate variations and ensure the integrity of the analyzed data.

Subsequently, we focused on transforming these normalized texts into vector representations using the BioWordVec model. This model, specifically trained on extensive medical corpora, was selected for its ability to accurately capture the semantics and clinical context of the terms used in the notes. By converting the texts into 200-dimensional vectors, representations of unrecognized words were adjusted using the <OOV> token, following a standardized approach for sequence length. This text-to-embeddings transformation procedure is essential for subsequent analysis using NLP techniques.

We used BioWordVec based a previous study, which findings across five models utilizing various word embeddings indicate that BioWordVec embeddings marginally enhanced the Bi-LSTM model's performance for certain datasets. Overall, models incorporating BioWordVec embeddings exhibited slightly superior performance compared to those utilizing GloVe embeddings[9] .

Through tokenization and sequence adjustment, we prepared the data for processing by advanced models such as GRU and ClinicalBERT. These models require structured and coherent inputs to effectively interpret the information contained in the clinical notes and thus identify possible errors. The meticulousness of this approach highlights the importance of preprocessing in NLP-supported clinical research. By transforming clinical notes into contextualized embeddings, we facilitate deep and accurate analysis by NLP models, enhancing error detection. This process not only enhances the analytical capability of the models but also underscores the value of rigorous data preparation in the field of artificial intelligence applied to medicine.

## 6  Data Transformation

After normalizing the data, we proceeded with its segmentation into training and test sets, adjusting this split according to the specific model to be used and experimenting with different partitions to always seek optimal accuracy. For the analysis with GRU, we selected an 80-20 split for training and testing, respectively, while for the evaluation using BioWordVec and ClinicalBERT, the distribution was adjusted to 70-30. This differentiation allowed us to adapt the learning and validation process to the peculiarities of each model, optimizing their ability to analyze and understand complex clinical texts.

This meticulous preparation and segmentation

Figure 1: Model GRU



Figure 2: Model Bio Clinical

of the data reflect the rigor with which we approach the implementation of advanced NLP techniques. By establishing solid foundations for the training and evaluation of models such as GRU, BioWord-Vec, and ClinicalBERT, our goal is to maximize their effectiveness in the precise identification of errors in medical documentation. This commitment to a detailed and adaptive methodology underscores our objective to advance the application of artificial intelligence to improve the accuracy and reliability of clinical documentation.

## 7 Feature Extraction

The process of extracting lexicographic features involved analyzing fundamental aspects of the text, such as the use of specific terms and the overall semantic structure of the clinical notes. This included evaluating polarity and the frequent use of parts of speech, which are indicative of the tone and intention of the medical text. Through this analysis, we sought to better understand how lexicographic features can influence the presence of errors within the notes.

For the GRU-based model, we adjusted the class weights to address the imbalance in our data, using the number of unique classes derived from the training set. This adjustment was crucial for training a balanced model capable of effectively classifying texts based on the presence or absence of medical errors. The GRU model was configured with layers specifically designed to capture and analyze complex temporal dependencies within clinical texts, including regularization layers to prevent overfitting and optimize overall performance.

Simultaneously, for the implementation based on ClinicalBERT, we proceeded with data tokenization and preparation using the AutoTokenizer from 'emilyalsentzer/BioClinicalBERT'. This preparation was essential to adapt our clinical notes to the format required by ClinicalBERT, allowing the model to process and classify the texts efficiently. The training of the model focused on binary classification of texts, training on contextualized representations generated to identify the presence of errors with high precision.

The training of the GRU and ClinicalBERT models was conducted under carefully selected configurations to optimize their learning and evaluation on the dataset. These configurations included defining the number of epochs, batch size, and learning rate, which are fundamental elements for the success of the training.

## 8 Settings

In the setup of the study, specific adjustments were made to the hardware and software parameters to optimize the analysis of the GRU and Clinical-BERT models. These adjustments included the optimization of processors and the allocation of execution threads, essential for the efficient processing of the clinical dataset.

Additionally, differentiated configurations were implemented in the software environment to adapt to the peculiarities of each model. This involved optimizing data loading, preprocessing, and embedding generation, ensuring that both GRU and ClinicalBERT operated under optimal conditions for text analysis. Adapting the computational environment allowed for maximizing the capabilities of each model, facilitating a thorough and precise analysis of clinical texts.

The computational infrastructure was also configured to log and store the highest performing features and classifiers during the experimental phase. This systematic approach allowed for continuous monitoring of model performance, providing a solid foundation for iteration and enhancement of analysis strategies.

This detailed setup reflects the methodical and rigorous approach adopted for the preparation and execution of the NLP models. By optimizing computational resources and adapting the software, the necessary conditions were established for an effective evaluation of the models' ability to identify errors in medical documentation.

## 9 Experiments and Analysis of Results

Comprehensive evaluations of multiple natural language processing models were conducted using the dataset provided by MEDIQA-CORR @ NAACL-ClinicalNLP 2024, with the goal of identifying those with the best performance in detecting errors in clinical notes . These experiments not only allowed for the adjustment of model configurations but also served to identify optimal techniques that significantly contribute to the analysis of medical texts. Among the evaluated models, GRU and BioClinicalBERT proved to be the most effective across various metrics and scenarios, which is why they were selected for further detailed analysis.

During the initial evaluations, models such as RF, RoBERTa, BERT, and CNN were also tested.

Hyperparameters for these models were adjusted to obtain better results, revealing their potential when dealing with larger datasets [7]. The implementation of RF and CNN models highlighted the importance of feature identification and automatic feature extraction, respectively [10]. Moreover, the use of BERT models leveraged the transformer architecture to pre-train language representations, enhancing the understanding of context and semantics in clinical terms [7]. This extensive evaluation facilitated the refinement of strategies and parameters for each model, aiming to maximize their accuracy in classifying texts based on the presence of errors.

Throughout multiple iterations in the preevaluation phase, strategies and parameters for each of these selected models were refined with the goal of maximizing their ability to classify texts accurately based on the presence of errors. Standard competition metrics, with a special emphasis on accuracy (ACC), were employed to measure the performance and effectiveness of the developed systems.

The experiments revealed notable differences in the efficacy of the GRU and BioClinicalBERT models for analyzing the medical corpus. While GRU excelled in its ability to process text sequences and capture temporal dependencies, BioClinicalBERT proved to be particularly effective in understanding the context and specific semantics of clinical terms. This distinction underscores the complementarity of the models in handling complex medical texts.

The results, summarized in Table 1, provide a clear view of the performance of the models under study. Compared to other traditional classification algorithms, GRU and BioClinicalBERT provided a deeper and more nuanced analysis of clinical notes, demonstrating their superiority in identifying inaccuracies and textual anomalies.

This detailed analysis reinforces the importance of adopting advanced NLP approaches in the realm of clinical documentation. The findings not only demonstrate the viability of these models to improve error detection in medical texts but also open new avenues for future research in the field of NLP applied to health, marking a step forward in the goal of elevating the quality and reliability of medical information through technology.

## 10 Result Test

Table 2 summarizes the performance of various classifiers in terms of accuracy during the training and testing phases, showing both the absolute accuracy (Training Accuracy, Testing Accuracy) as well as the accuracy differences between these phases for each evaluated model. This initial evaluation allowed us to identify models with promising performance.

Among the evaluated models, ClinicalBERT and GRU stood out for their robust performance across various metrics and were selected for further detailed analysis. After rigorous validation, which included reviewing performance and learning curves, Table 3 details the accuracy of these models on the validation set, confirming their efficacy.

The selection of ClinicalBERT and GRU was based on a rigorous analysis of their capacity to process and analyze complex clinical texts, showing notable superiority in identifying errors in medical documentation. The validation of these models confirms the effectiveness of our selection strategy and highlights the importance of exploring in depth how these models can contribute to improving the analysis of clinical notes in the future.

## 11 Discussion and Conclusion

The meticulous selection of embeddings and NLP models, specifically GRU and ClinicalBERT, is crucial for the accurate analysis of clinical texts, as evidenced in our findings. These decisions are vital for optimizing error detection in clinical notes. However, there is a need to expand experimentation with a broader spectrum of models and embeddings to validate their effectiveness in specific clinical contexts. The analysis of the results, presented in Table 3, compares models from the most basic to the more complex ones (excluding large language models), revealing a progression and the importance of a detailed methodology and the adaptation of models to clinical textual peculiarities. This approach underscores the urgency of increasing experimentation to enhance precision and applicability in improving clinical documentation.

Despite considering the use of advanced LLMs like Gemini or ChatGPT-4, this study highlights the efficacy of alternative models such as ClinicalBERT and GRU. This preference is due not only to their competent performance but also to their specific adaptability to the demands of clinical texts. This approach is crucial in environments

| Modelos | Train | | Test | |
|---|---|---|---|---|
| | **F1** | **Acc.** | **F1** | **Acc.** |
| Roberta | 0.71 | 0.55 | 0.71 | 0.56 |
| Roberta_Sobremuestreo | 0.64 | 0.53 | 0.61 | 0.47 |
| Roberta_96_warmup_steps_9_epochs | 0.64 | 0.56 | 0.55 | 0.45 |
| Roberta_48_warmup_steps_15_epochs | 0.86 | 0.86 | 0.5 | 0.46 |
| Roberta_15_epochs | 0.88 | 0.87 | 0.52 | 0.48 |
| Roberta_20_epochs | 0.66 | 0.49 | 0.67 | 0.5 |
| Roberta_25_epochs | 0.99 | 0.99 | 0.51 | 0.48 |
| Roberta_30_epochs | 0.99 | 0.99 | 0.51 | 0.48 |
| Roberta_35_epochs | 0.99 | 0.99 | 0.44 | 0.5 |
| Roberta_40_epochs | 0.99 | 0.99 | 0.48 | 0.51 |
| Roberta_sobremuestro_steps_45_epochs | 1 | 1 | 0.41 | 0.43 |
| Bio_medical_sobremuestro_5_epochs | 0.67 | 0.67 | 0.53 | 0.51 |
| Bio_medical_96_warmup_steps_5_epochs_8_batch | 0.69 | 0.7 | 0.48 | 0.48 |
| Bio_medical_96_warmup_steps_10_epochs_8_batch | 0.95 | 0.95 | 0.45 | 0.45 |
| Bio_medical_sobremuestro_10_epochs_16_batch | 0.94 | 0.94 | 0.48 | 0.47 |
| Bio_medical_sobremuestro_7_epochs_16_batch | 0.8 | 0.78 | 0.49 | 0.45 |
| Gpt2-medium_1_batch | 0.66 | 0.5 | 0.7 | 0.54 |
| Longformer-base-4096 | 0.66 | 0.5 | 0.73 | 0.58 |
| Random Forest_split_vectorizar | 0.36 | 0.56 | 0.35 | 0.54 |
| Random Forest_vectorizar_split | 0.36 | 0.56 | 0.35 | 0.53 |
| Random Forest_vectorizar_split_10_leaf | 0.35 | 0.56 | 0.35 | 0.53 |
| Random Forest_80_train_20_test | 0.35 | 0.56 | 0.33 | 0.51 |
| Stacking RL, SVC y RF | 0.3609 | 0.3597 | 0.7514 | 0.7511 |
| Stacking RL, SVC, RF, GB y DT | 0.023 | 0.021 | 0.822 | 0.8219 |
| GRU_No_Embbeding | 0.665 | 0.5 | 0.69 | 0.53 |
| GRU_glove-wiki-gigaword-200 | 0.47 | 0.52 | 0.33 | 0.42 |
| GRU_glove-wiki-gigaword-200_DropOut | 0.98 | 0.97 | 0.54 | 0.51 |
| GRU_BioWordVec_PubMed_MIMICIII_d200 | 0.67 | 0.56 | 0.6 | 0.49 |
| GRU_BioWordVec_MIMICIII_d200_desbalanceo | 0.71 | 0.56 | 0.69 | 0.53 |
| GRU_BioWordVec_MIMICIII_d200_L2 | 0.72 | 0.56 | 0.7 | 0.54 |
| GRU_BioWordVec_MIMICIII_d200_L1_L2 | 0.53 | 0.57 | 0.52 | 0.53 |
| LR | 0.2836 | 0.2841 | 0.2836 | 0.2841 |
| CNN | 0.4043 | 0.5619 | 0.3746 | 0.5365 |
| RNN | 0.2668 | 0.4380 | 0.2935 | 0.4634 |
| LSTM | 0.4043 | 0.5619 | 0.3746 | 0.5365 |

Table 1: Detailed Model Results

| Classifiers | Train Acc. | Diff. Train Acc. | Test Acc. | Diff. Test Acc. |
|---|---|---|---|---|
| ClinicalBERT | 0.77 | 0.00 | 0.51 | 0.00 |
| GRU | 0.57 | -0.20 | 0.53 | 0.02 |
| Random Forest | 0.56 | -0.11 | 0.54 | 0.03 |
| CNN | 0.56 | -0.11 | 0.53 | 0.02 |
| LSTM | 0.56 | -0.11 | 0.53 | 0.02 |
| RoBERTa | 0.55 | -0.12 | 0.56 | 0.05 |
| GPT-2 | 0.50 | -0.17 | 0.54 | 0.03 |
| Longformer | 0.50 | -0.17 | 0.58 | 0.07 |
| Stacking RL, SVC y RF | 0.35 | -0.32 | 0.75 | 0.24 |
| RL | 0.28 | -0.39 | 0.28 | -0.23 |
| Stacking RL, SVC, RF, GB y DT | 0.21 | -0.46 | 0.82 | 0.31 |

Table 2: Model and Results Record

| Classifiers | Acc Validation | Diff. Acc. |
|---|---|---|
| ClinicalBERT | 0.56 | 0.00 |
| GRU | 0.55 | -0.01 |

Table 3: Selected Models Validation Set Results

where data security, privacy, and time and resource constraints are primary considerations. In such contexts, the need for efficient yet less demanding models makes specialized alternatives surpass more generalist LLMs, aligning better with practical limitations and data protection imperatives in clinical research

Throughout the experiments conducted, it was observed that specific models such as GRU and ClinicalBERT demonstrate significant potential in processing medical text, emphasizing that, with proper data preparation and model tuning, it is possible to effectively manage the complexities inherent in clinical texts. Although the highest accuracy percentages obtained do not significantly exceed the random decision threshold, these results do not detract from the effectiveness of the models employed, but rather underline the importance of a meticulous selection and configuration of modeling features and parameters.

This study demonstrates that advancements in NLP can significantly contribute to the clinical field, although it also highlights the ongoing challenge of adapting these technologies to the specificities of medical language and data. NLP models, even in the face of accuracy challenges, prove to be valuable tools when carefully adjusted based on a deep understanding of the context and specific objectives of the task.

For future research, feature selection is highlighted as the primary strategy. It is suggested to focus on the development and application of advanced methodologies for feature extraction and selection, with the aim of refining the analytical capabilities of models for the precise processing and understanding of medical texts. This methodological approach not only anticipates an increase in the accuracy of models for anomaly detection and error identification in clinical documentation but also promises to deepen our understanding of the adaptation and optimization of NLP techniques for specific needs within the healthcare domain.

In conclusion, this study significantly contributes to the field of NLP applied to the medical domain, promoting the continuous innovation and optimization of models that, through meticulous choice and configuration of features, have vast potential to elevate the quality of clinical documentation. A notable finding is the moderate impact that pre-trained embeddings have on model performance, indicating that the integration and thorough exploration of these pre-trained tools can be crucial for amplifying the effectiveness of NLP in clinical contexts. This constant adaptation and improvement of technologies promise to advance towards optimizing the practical utility of NLP models, thereby contributing to improving the standards of care and documentation in the healthcare sector.

# References

[1] Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

[2] Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin.

[3] Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, and Christopher D. Manning. BioMedLM: A 2.7b parameter language model trained on biomedical text. *Preprint*, arxiv:2403.18421 [cs].

[4] Matt Casey. 2023. Large language models: their history, capabilities and limitations.

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint*. ArXiv:1810.04805 [cs].

[6] Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. 2020. ClinicalBERT: Modeling Clinical Notes and Predicting Hospital Readmission. *arXiv preprint*. ArXiv:1904.05342 [cs].

[7] Jyoti Kumari and Abhinav Kumar. 2023. JA-NLP@LT-EDI-2023: Empowering Mental Health Assessment: A RoBERTa-Based Approach for Depression Detection. In *Proceedings of the Third Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 89–96, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

[8] Lishuang Li, Jia Wan, Jieqiong Zheng, and Jian Wang. 2018. Biomedical event extraction based on GRU integrating attention mechanism. *BMC Bioinformatics*, 19(9):285.

[9] Hongxia Lu, Louis Ehwerhemuepha, and Cyril Rakovski. A comparative study on deep learning models for text classification of unstructured medical notes with various levels of class imbalance. 22(1):181.

[10] Yuanren Tong, Keming Lu, Yingyun Yang, Ji Li, Yucong Lin, Dong Wu, Aiming Yang, Yue Li, Sheng Yu, and Jiaming Qian. Can natural language processing help differentiate inflammatory intestinal diseases in china? models applying random forest and convolutional neural network approaches. 20(1):248.

[11] Yuqing Wang, Yun Zhao, and Linda Petzold. 2023. Are Large Language Models Ready for Healthcare? A Comparative Study on Clinical Language Understanding. *arXiv preprint*. ArXiv:2304.05368 [cs].

[12] Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Clinical text classification with rule-based features and knowledge-guided convolutional neural networks. *BMC Medical Informatics and Decision Making*, 19(3):71.

[13] Ian James Bruce Young, Saturnino Luz, and Nazir Lone. 2019. A systematic review of natural language processing for classification tasks in the field of incident reporting and adverse event analysis. *International Journal of Medical Informatics*, 132:103971.

# HSE NLP Team at MEDIQA-CORR 2024 Task: In-Prompt Ensemble with Entities and Knowledge Graph for Medical Error Correction

**Airat Valiev**
HSE University / Moscow, Russia
aa.valiev@hse.ru

**Elena Tutubalina**
Kazan State University / Kazan, Russia
HSE University / Moscow, Russia
tutubalinaev@gmail.com

## Abstract

This paper presents our LLM-based system designed for the MEDIQA-CORR @ NAACL-ClinicalNLP 2024 Shared Task 3, focusing on medical error detection and correction in medical records. Our approach consists of three key components: entity extraction, prompt engineering, and ensemble. First, we automatically extract biomedical entities such as therapies, diagnoses, and biological species. Next, we explore few-shot learning techniques and incorporate graph information from the MeSH database for the identified entities. Finally, we investigate two methods for ensembling: (i) combining the predictions of three previous LLMs using an AND strategy within a prompt and (ii) integrating the previous predictions into the prompt as separate 'expert' solutions, accompanied by trust scores representing their performance. The latter system ranked second with a BERTScore score of 0.8059 and third with an aggregated score of 0.7806 out of the 15 teams' solutions in the shared task.

## 1 Introduction

Medical records play a crucial role in healthcare systems as they capture essential patient information, including diagnoses, treatments, and outcomes. Medical texts are characterized by complex terminology, context-specific knowledge, and significant implications. Detecting and rectifying errors within clinical notes necessitates domain expertise and reasoning. This task presents a complex challenge that demands precise analysis and understanding of the medical domain.

In recent years, Large Language Models (LLMs) have revolutionized the field of natural language processing (NLP) by demonstrating unprecedented performance across a wide range of tasks. These models, often based on Transformer (Vaswani et al., 2017; Devlin et al., 2018), have become the cornerstone of modern NLP research (Pan et al., 2023). LLMs excel in key areas such as semantic un-

derstanding and contextualization (Radford et al., 2018), multimodal capabilities (Livne et al., 2023), few-shot and zero-shot learning (Dang et al., 2022), as well as various medical applications including disease diagnosis (Schubert et al., 2023), drug discovery (Livne et al., 2023), and medical records processing (Guevara et al., 2024).

Automated fact-checking has garnered significant attention due to the escalating challenge posed by misinformation. Traditionally, fact-checking has relied on manual verification conducted by human experts, primarily focusing on general-domain texts like Wikipedia articles and news reports (Zhang and Gao, 2023; Quelle and Bovet, 2024). Recently, LLMs have offered the capability to analyze false statements and provide an assessment of their factual accuracy by leveraging their pre-trained knowledge and contextual understanding (Wang and Shu, 2023; Guan, 2021; Lewis et al., 2020; Chen et al., 2021). Several methodologies have been proposed to enhance the overall performance in LLMs, and the most notable ones are Chain of Thought (CoT) (Zhang, 2023).

In this work, we utilize several key approaches for medical records correction using LLMs (see Figure 1). These approaches include entity extraction and normalization (Miftahutdinov et al., 2020, 2021; Sung et al., 2022), few-shot learning techniques (Brown et al., 2020), graph-based knowledge incorporation (Fei et al., 2021), and ensembling strategies (Wang et al., 2022). We investigate the application of these approaches to enhance the accuracy of medical error correction.

The paper is organized as follows. Section 1 presents shared task and data overview. We describe our approach with three key components and state-of-the-art (SoTA) models in Section 2. Experiments with baselines and our model are presented in Section 3.3.4. Finally, we discuss the results and conclude the work in Sections 4 and 5, respectively.

Figure 1: The system overview. The process can be described as follows: the system begins by receiving medical text as input. Initially, a prompt template is utilized, supplemented with a small number of few-shot examples (either 2 or 5). The Named Entity Recognition (NER) model is then employed to identify and extract named entities within the large language model's context. Subsequently, potential replacements for these extracted entities are sought within the Medical Subject Headings (MeSH) thesaurus. The prompt, enriched with these replacements, is passed to the selected OpenAI model. Finally, the model's output is returned and stored in the prediction file. This constitutes the overall operation of the system.

## 2 Task and Data Overview

The MEDIQA-CORR 2024 shared task (Ben Abacha et al., 2024) focuses on analyzing snippets of clinical text to address specific subtasks related to medical error detection and correction. These subtasks include:

1. Binary Classification: The first subtask involves determining whether the given clinical text contains a medical error. This step requires evaluating the accuracy, consistency, and factual correctness of the information presented in the text.

2. Span Identification: If a medical error is identified, the goal is to locate the specific text span associated with the error. This step is crucial for precisely pinpointing the erroneous segment within the clinical text.

3. Natural Language Generation (Correction): Once the medical error is identified and its location is determined, the task is to generate a free-text correction for the identified error. The generated correction should be contextually appropriate, accurate, and concise, effectively addressing the error in the clinical text.

We focus on the latter subtask, which encompasses all three subtasks mentioned.

The dataset provided by the organizers, known as the MS Training Set, consists of 2,189 clinical texts. Additionally, there is the 'MS' Validation Set comprising 574 clinical texts and the 'UW' Validation

Set comprising 160 clinical texts (Abacha et al., 2024). The test portion of the dataset is formed by combining clinical texts from both collections. Each clinical text in the dataset is labeled as either correct or containing one error. More formally, the task involved in this dataset is as follows:

1. Predicting whether a given text contains an error or not. The error flag is represented by 1 if the text contains an error and 0 if it is error-free.

2. For texts flagged as containing errors, extract the sentence that contains the error.

3. Generating a corrected version of the identified error sentence.

## 3 Method

The error correction method of Figure 1, proposed in the current work, is straightforward and consists of three major steps: we first prepare the data to make predictions: extract named entities from texts, and search for the term replacements. Then we form the prompt for the model from the template, add a few examples and additional data (NER results, MeSH terms), and then use LLMs to make predictions.

### 3.1 Data preparation

Let us first discuss the first step. Before predicting, some preparations were made with the input texts, including Named Entity Recognition (NER), and

Figure 2: An example of the result obtained from Named Entity Recognition (NER).



Figure 3: An example for the term D014883 (water-electrolyte imbalance) related entities, extracted from the MeSH database.

possible term replacement data extraction from the MeSH thesaurus.

### 3.1.1 Biomedical entities

Biomedical concepts, such as diseases, symptoms, drugs, genes, and proteins, are critical for many biomedical applications, including drug discovery (Khrabrov et al., 2022), clinical decision making (Sutton et al., 2020; Peiffer-Smadja et al., 2020), and biomedical research (Lee et al., 2016; Tutubalina et al., 2017; Soni and Roberts, 2021; Sakhovskiy et al., 2021; Sakhovskiy and Tutubalina, 2022; Miftahutdinov et al., 2020, 2021).

For NER, we use the BERN2 (Advanced Biomedical Entity Recognition and Normalization) model (Sung et al., 2022) is a neural biomedical named entity recognition and normalization tool. BERN2 significantly improves upon its predecessor (Kim et al., 2019) by employing a multi-task NER model and neural network-based entity linking (EL) models, resulting in faster and more accurate inference.

Using this tool, we extracted named entities (with MeSH identifiers) such as diagnosis, therapy, biological species, and medical entities. You can see an example of such extraction in Figure 2.

### 3.1.2 MeSH: Medical Subject Headings

MeSH is a hierarchically organized and concept-based vocabulary produced by the National Library of Medicine (NLM) (Mao and Lu, 2017). Its primary purpose is to facilitate indexing, cataloging, and searching of biomedical and health-related information. MeSH plays a crucial role in various NLM databases, including MEDLINE/PubMed and the NLM Catalog. MeSH consists of standardized keywords that describe the subject matter of journal articles, clinical notes, and other biomedical texts. These terms are carefully curated and organized to ensure consistency and accuracy. Researchers, librarians, and information specialists use MeSH to index and retrieve relevant literature. By assigning MeSH terms to documents, they enhance search precision and recall. MeSH thesaurus could be applied to perform Biomedical Literature Indexing (like in MEDLINE/PubMed (von Korff, 2022)), Concept Mapping, and Synonyms (MeSH provides a standardized way to map synonyms and related terms, for different synonyms of a medical condition to be linked to a single MeSH term), and investigating cross-lingual clinical entity linking using MeSH concepts. Highlights the importance of MeSH in linking biomedical entities across languages. MeSH serves as a foundational resource for organizing and accessing biomedical knowledge. Its controlled vocabulary ensures consistency and precision, benefiting researchers, clinicians, and information professionals.

In the presented work, we use the MeSH database to perform the knowledge graph search - for the extracted entities with available MeSH IDs, we've found their possible replacements (the example is in Figure 3) (other entities on the same relation level with the parent term node) to present them to the LLM as clues about possible errors in a text.

## 3.2 Dataset description

The statistical data about the dataset can be seen from the Table 1. In total, 2,923 texts (2,189 texts

Figure 4: The solution ensembling overview. In this approach, we use previous predictions of different models for each input text and resulting prediction scores, and a new template. We evaluated three major ensembling strategies, including AND (all three models found an error), majority of votes, and weighted approach (weight prediction by each prediction score), in the validation stage, but decided to make the final prediction using AND strategy.

|           | Train | Val MS | Val UW | Test |
|-----------|-------|--------|--------|------|
| Texts     | 2 189 | 574    | 160    | 925  |
| NER ent.  | 3,3   | 3,3    | 3,3    | 6,5  |
| MeSH terms| 2,1   | 2,1    | 2,1    | 2,2  |

Table 1: Dataset statistics by the number of texts and found entities.

in the train part + 574 texts in the MS validation part + 160 texts in the UW validation part), the BERN2 model found 9,682 named entities with MeSH IDs, an average of 3.3 entities per single text. An average of 2.1 MeSH term replacements were found using MeSH graph search. The train part included 1,219 texts with errors and 970 correct entries, the MS validation part consisted of 80 correct and 80 with errors, and the UW validation included 319 entries with errors and 255 correct.

The test data part consisted of 925 text entries. During the test part processing, the BERN2 model extracted 6,032 MeSH IDs (avg. 6.5 terms per text), with an average of 2.2 replacements extracted from the MeSH thesaurus.

## 3.3 Making predictions

After the preparation step, we move forward to make the predictions and find the texts with medical errors. We have studied and used three general LLM-based approaches for prediction making:

1. Ordinary prompting (2-shot and 5-shot)

2. Prediction ensembling (ensemble of 3 solutions)

3. In-prompt ensembling (expert opinions with trust scores)

In this section, we first discuss the ordinary solution with different OpenAI models (GPT3.5-turbo, GPT4, GPT4-turbo preview) (Yenduri et al., 2022) and simple prompts. These models continually improve the instruction following ability and have broader general knowledge and advanced reasoning capabilities. The solution idea is simple, as we discussed earlier: the model receives the prompt prefix containing the behavior rules for the model (see Appendix 1), 2 of 5 examples (texts and expected output from the training dataset part), and the text to analyze along with the NER information (found named entities) and replacement entities from the MeSH graph. All significant parts of the template are highlighted in color. A few shot examples fixed set (2 or 5) were selected from the Train data split to present the data with and without corrections needed equally.

### 3.3.1 Ordinary prompting (2-shot)

The first solution (as illustrated in Figure 1), with the 2-shot template, consists of a prefix (2-shot prompt prefix from Appendix A1), text to predict, and additional data: a list of found named entities with additional info from the BERN2 model.

### 3.3.2 Ordinary prompting (5-shot)

The second solution with the 5-shot prompt template, is constructed from a 5-shot prompt prefix from Appendix A.1.3, the text to predict and additional data: NER results and MeSH graph data with possible entity replacements. The process scheme is illustrated in Figure 1.

### 3.3.3 Prediction ensembling

Decision ensembling is different from previously discussed approaches. In this variant, as it is

Figure 5: The system overview. The medical text inputs into the system. First of all, we use the prompt template and add 2 or 5 few-shot examples. The NER model finds named entities for the large language model. Then we find possible replacements for extracted entities in the MeSH thesaurus. Here we also use previous predictions of different models for each input text and resulting prediction scores, adding these 'expert' opinions with expert trust scores, to the prompt. The generated prompt is passed to the openAI model of our choice, and we return the result to the prediction file.

shown in Figure 4, we simply construct the prediction from the three top-score previous predictions, based on AND strategy: for each text entry we decide the error exists, if only the error is found in all the three previous predictions - in this case we include the error sentence number and correction from the previous prediction with the highest score. If at least one model has predicted this sentence as correct, we count it as containing no errors. This strategy slightly improved the resulting score: 0.62 -> 0.64.

### 3.3.4 In-prompt ensembling

In this approach, we have combined the idea of basic prompting, few-shot learning, and an ensemble of experts. We again add information about NER entities and MeSH graph replacements, but because of the ensembling approach evaluated, we also include predictions from the top three previous submissions (model predictions with the highest score), calling it 'expert's solutions'. We also add three expert trust scores - these are the test scores for these submissions, to help the model estimate the expert opinion correctness indirectly.

We added the test predictions of the three previous models. Still, in the case of a real data evaluation, this ensemble could be formed from the three different models and their predictions, and trust scores could be obtained from the validation scores.

The result of this ensemble addition could be the following: "Expert 1 with trust score (weight1): (outputs1), expert 2 with trust score (weight2): (outputs2), expert 3 with trust score (weight): (outputs3)." Prompt prefix (Appendix A.2.3) and process scheme 5 are included. The ensemble example

with the real data is the following:

- Expert 1 with trust score 0,72: "Error exists: |||Yes||| Correction: ||| Patient's symptoms are suspected to be due to acute gastroenteritis.||| Error sentence number: |||10|||",

- Expert 2 with trust score 0,69: "Error exists: |||Yes||| Correction: ||| Patient's symptoms are suspected to be due to typhoid fever.||| Error sentence number: |||10|||",

- Expert 3 with trust score 0,68: "Error exists: |||No||| Correction: |||None||| Error sentence number: |||None|||"

## 4 Baselines

During the model development and preparation, we explored various baselines. In addition to the above-mentioned methods, we initially investigated a simpler BERT-based approach (Devlin et al., 2018) and utilized other LLMs such as self-hosted LLaMA-based Med42 70b (Christophe et al., 2023) and Meditron 7b (Chen et al., 2023).

The BERT model, specifically the PubMedBERT-base checkpoint (Gu et al., 2021), was trained for 10 epochs on a subset of the training data. However, it performed poorly on the validation data, achieving a score of approximately 0.57 even on the task of text classification for error presence, which is a binary classification problem. This subpar performance can be attributed to a limited number of training examples and the wide variation in replaceable terms and diverse themes found in medical texts. Due to these unsatisfactory results in the validation phase, we decided not

Table 2: Evaluation results. Here 'ens' stands for an ensemble of 3 previous solutions and these predict scores, 'NER' - for named entities from the text, and 'MeSH' - for the related terms from the MeSH thesaurus. The general approach is shown in Figure 1, the prediction ensemble - in Figure 4, and an ensemble of experts in Figure 5.

| Base model name | Prompt | Additional data | AggrScore | R1F | BERTScore | BLEURT | AggrC |
|---|---|---|---|---|---|---|---|
| gpt-3.5t | General 2-shot | NER | 0.31 | 0.35 | 0.38 | 0.34 | 0.24 |
| gpt-4-t-1401-preview | 5-shot | NER | 0.55 | 0.55 | 0.55 | 0.55 | 0.41 |
| gpt-4-t-preview-0125 | 5-shot | NER + meSH | 0.62 | 0.62 | 0.60 | 0.62 | 0.53 |
| - | - | Ens. of 3 predicts | 0.64 | 0.64 | 0.62 | 0.63 | 0.54 |
| gpt-4-t-preview-0125 | Ensemble prompt | NER+ens | 0.68 | 0.68 | 0.67 | 0.68 | 0.52 |
| gpt-4-t-0125-preview | Ensemble prompt | NER+ens+MeSH | 0.69 | 0.71 | 0.67 | 0.69 | 0.51 |
| gpt-4 | Ensemble prompt | NER+ens+MeSH | 0.72 | 0.74 | 0.69 | 0.72 | 0.55 |
| gpt-4-t-0125-preview | Ensemble prompt | NER+ens+MeSH | **0.78** | **0.81** | **0.76** | **0.78** | **0.51** |

to proceed with evaluating the model's precision on the test data and instead moved on to explore alternative solution methods.

The LLaMA-based models exhibited better performance and were successful in identifying and correcting misplaced terms, achieving an aggregated score of approximately 0.43 on the validation data. However, these models disregarded certain in-prompt rules and ensemble solutions. Consequently, despite not showing any positive performance improvements with the addition of NER data and graph entities, they were excluded from the test submission.

## 5   Experiments and Results

The evaluation results of our error correction systems are shown in Table 2. The aggregate score is the main evaluation score to rank the participating systems. We've used the following scripts[1] for evaluation. More specifically about the metrics used for evaluation:

- NLG (Natural Language Generation) metrics: ROUGE(Lin, 2004), BERTScore (Zhang et al., 2019), BLEURT(Sellam et al., 2020), their Aggregate-Score (Mean of ROUGE-1-F, BERTScore, BLEURT-20), and their Composite Scores (AggrC) for the evaluation of Sentence Correction.

- The Composite score is the mean of individual scores computed as follows for each text:

  - 1 point if both the system correction and the reference correction are "NA";
  - 0 point if only one of the system or the reference is "NA".

- NLG metrics value in [0, 1] range (e.g., ROUGE, BERTScore, BLEURT, or Aggregate-Score) if both the system correction and reference correction are non-"NA" sentences.

- The Aggregate score is the main evaluation score to rank the participating systems(Abacha et al., 2023).

As we can see from table 2, we can observe that the more powerful language model, 'sophisticated' prompting, and additional data presented to the language model lead to better results: results improved from 0.62 to 0.72 and finally to 0.78, which is a Top-3 solution of an entire competition. One also can see that additional examples (2 vs 5 texts) in the few-shot section also increase performance: 0.31 vs 0.55. Also, the in-prompt ensembling technique improves final results greatly because the model can see the solutions from previous runs along with the scores for these runs, and correct the current prediction, which leads to more stable and reliable predictions and error corrections. We also could see the obvious trend of better performance with more complicated models: GPT 4 outperforms GPT 3.5 Turbo, and GPT 4 Turbo preview beats the ordinary GPT 4: 0.31 vs 0.55 vs 0.62, respectively.

The methodology delineated herein possesses the potential for expansion and further refinement through the incorporation of techniques such as the

Knowledge Graph, PromptKG (Xie et al., 2022), the meta-prompting approach, and the Chain of Thought (CoT) approach. Additionally, the integration of specialized models, specifically designed for error detection and error span identification, into the model pipeline could be achieved directly by utilizing the chaining techniques (e.g. langchain). This would serve to enhance the robustness and accuracy of the overall system.

# 6 Conclusion

In this work, we have addressed the issue of identifying and resolving error text in biomedical texts. We have proposed a system for the MEDIQA-CORR shared task by utilizing prompting, ensembling techniques, and LLMs. Our approach demonstrates that the problem can be solved using ordinary GPT models without pre-training, relying solely on in-context learning, along with the NER model and additional MeSH knowledge graph data. By employing an in-prompt ensemble of LLMs as experts and incorporating data from the MeSH knowledge graph and NER results, we achieved a high task aggregated score of 0.78059, securing the 3rd position on the official competition leaderboard. Our results highlight the effectiveness of our proposed prompting approach while also indicating areas for future improvement. Utilizing more advanced tools like full-scale RAG and fine-tuned biomedical LLMs could potentially enhance the quality of error correction. In addition, we plan to make all our code and data publicly accessible shortly after the publication of our paper.

## Acknowledgments

## Limitations

Large Language Models (LLMs) have emerged as powerful tools for natural language understanding and generation. However, their effectiveness hinges on the quality and diversity of their training data. LLMs are typically trained on vast corpora of text from the internet, making manual curation infeasible. Consequently, they inherit any biases, inaccuracies, or limitations in their training data. Additionally, the success of LLMs has led to the generation of online content by these models, which may introduce hallucinated information.

## Ethics Statement

Using databases for retrieval carries a drawback in that these sources may lack comprehensiveness and contain inaccuracies. LLMs are not immune to representation biases and the risk of producing potentially misleading outcomes, particularly in the healthcare sector.

All pre-trained language models and datasets utilized in this study are openly accessible for research purposes.

We honor and support the ACL Code of Ethics.

## References

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Asma Ben Abacha, Wen-wai Yim, George Michalopoulos, and Thomas Lin. 2023. An investigation of evaluation metrics for automated medical note generation. *arXiv preprint arXiv:2305.17364*.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Yiran Chen, Pengfei Liu, and Xipeng Qiu. 2021. Are factuality checkers reliable? adversarial meta-evaluation of factuality in summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2082–2095, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. 2023. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*.

Clément Christophe, Avani Gupta, Nasir Hayat, Praveen Kanithi, Ahmed Al-Mahrooqi, Prateek Munjal, Marco Pimentel, Tathagata Raha, Ronnie Rajan, and

Shadab Khan. 2023. Med42 - a clinical large language model.

Hai Dang, Lukas Mecke, Florian Lehmann, Sven Goller, and Daniel Buschek. 2022. How to prompt? opportunities and challenges of zero-and few-shot learning for human-ai interaction in creative applications of generative models. *arXiv preprint arXiv:2209.01390*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Hao Fei, Yafeng Ren, Yue Zhang, Donghong Ji, and Xiaohui Liang. 2021. Enriching contextualized language model from knowledge graph for biomedical information extraction. *Briefings in bioinformatics*, 22(3):bbaa110.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23.

Shuo Guan. 2021. Knowledge and keywords augmented abstractive sentence summarization. In *Proceedings of the Third Workshop on New Frontiers in Summarization*, pages 25–32, Online and in Dominican Republic. Association for Computational Linguistics.

Marco Guevara, Shan Chen, Spencer Thomas, Tafadzwa L Chaunzwa, Idalid Franco, Benjamin H Kann, Shalini Moningi, Jack M Qian, Madeleine Goldstein, Susan Harper, et al. 2024. Large language models to identify social determinants of health in electronic health records. *NPJ digital medicine*, 7(1):6.

Kuzma Khrabrov, Ilya Shenbin, Alexander Ryabov, Artem Tsypin, Alexander Telepov, Anton Alekseev, Alexander Grishin, Pavel Strashnov, Petr Zhilyaev, Sergey Nikolenko, and Artur Kadurin. 2022. nablaDFT: Large-Scale conformational energy and hamiltonian prediction benchmark and dataset. *Phys. Chem. Chem. Phys.*, 24(42):25853–25863.

Donghyeon Kim, Jinhyuk Lee, Chan Ho So, Hwisang Jeon, Minbyul Jeong, Yonghwa Choi, Wonjin Yoon, Mujeen Sung, and Jaewoo Kang. 2019. A neural named entity recognition and multi-type normalization tool for biomedical text mining. *IEEE Access*, 7:73729–73740.

Sunwon Lee, Donghyeon Kim, Kyubum Lee, Jaehoon Choi, Seongsoon Kim, Minji Jeon, Sangrak Lim, Donghee Choi, Sunkyu Kim, Aik-Choon Tan, et al. 2016. Best: next-generation biomedical entity search tool for knowledge discovery from biomedical literature. *PloS one*, 11(10):e0164680.

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.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Micha Livne, Zulfat Miftahutdinov, Elena Tutubalina, Maksim Kuznetsov, Daniil Polykovskiy, Annika Brundyn, Aastha Jhunjhunwala, Anthony Costa, Alex Aliper, and Alex Zhavoronkov. 2023. nach0: Multimodal natural and chemical languages foundation model. *arXiv preprint arXiv:2311.12410*.

Yuqing Mao and Zhiyong Lu. 2017. Mesh now: automatic mesh indexing at pubmed scale via learning to rank. *Journal of biomedical semantics*, 8:1–9.

Zulfat Miftahutdinov, Ilseyar Alimova, and Elena Tutubalina. 2020. On biomedical named entity recognition: experiments in interlingual transfer for clinical and social media texts. In *European Conference on Information Retrieval*, pages 281–288. Springer.

Zulfat Miftahutdinov, Artur Kadurin, Roman Kudrin, and Elena Tutubalina. 2021. Medical concept normalization in clinical trials with drug and disease representation learning. *Bioinformatics*, 37(21):3856–3864.

Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2023. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188*.

Nathan Peiffer-Smadja, Timothy Miles Rawson, Raheelah Ahmad, Albert Buchard, P Georgiou, F-X Lescure, Gabriel Birgand, and Alison Helen Holmes. 2020. Machine learning for clinical decision support in infectious diseases: a narrative review of current applications. *Clinical Microbiology and Infection*, 26(5):584–595.

Dorian Quelle and Alexandre Bovet. 2024. The perils and promises of fact-checking with large language models. *Frontiers in Artificial Intelligence*, 7:1341697.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.

Andrey Sakhovskiy, Zulfat Miftahutdinov, and Elena Tutubalina. 2021. KFU NLP team at SMM4H 2021 tasks: Cross-lingual and cross-modal BERT-based models for adverse drug effects. In *Proceedings of the Sixth Social Media Mining for Health (#SMM4H) Workshop and Shared Task*, pages 39–43, Mexico City, Mexico. Association for Computational Linguistics.

Andrey Sakhovskiy and Elena Tutubalina. 2022. Multi-modal model with text and drug embeddings for adverse drug reaction classification. *Journal of Biomedical Informatics*, 135:104182.

Marc Cicero Schubert, Wolfgang Wick, and Varun Venkataramani. 2023. Large language model-driven evaluation of medical records using medcheckllm. *medRxiv*, pages 2023–11.

Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*.

Sarvesh Soni and Kirk Roberts. 2021. An evaluation of two commercial deep learning-based information retrieval systems for COVID-19 literature. *J. Am. Medical Informatics Assoc.*, 28(1):132–137.

Mujeen Sung, Minbyul Jeong, Yonghwa Choi, Donghyeon Kim, Jinhyuk Lee, and Jaewoo Kang. 2022. Bern2: an advanced neural biomedical named entity recognition and normalization tool. *Bioinformatics*, 38(20):4837–4839.

Reed T Sutton, David Pincock, Daniel C Baumgart, Daniel C Sadowski, Richard N Fedorak, and Karen I Kroeker. 2020. An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ digital medicine*, 3(1):17.

EV Tutubalina, Z Sh Miftahutdinov, RI Nugmanov, TI Madzhidov, SI Nikolenko, IS Alimova, and AE Tropsha. 2017. Using semantic analysis of texts for the identification of drugs with similar therapeutic effects. *Russian Chemical Bulletin*, 66:2180–2189.

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.

Modest von Korff. 2022. Exhaustive indexing of PubMed records with medical subject headings. In *Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-10)*, pages 8–15, Marseille, France. European Language Resources Association.

Haoran Wang and Kai Shu. 2023. Explainable claim verification via knowledge-grounded reasoning with large language models. *arXiv preprint arXiv:2310.05253*.

Lijing Wang, Timothy Miller, Steven Bethard, and Guergana Savova. 2022. Ensemble-based fine-tuning strategy for temporal relation extraction from the clinical narrative. In *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 103–108, Seattle, WA. Association for Computational Linguistics.

Xin Xie, Zhoubo Li, Xiaohan Wang, Shumin Deng, Feiyu Xiong, Huajun Chen, and Ningyu Zhang. 2022.

Promptkg: A prompt learning framework for knowledge graph representation learning and application. *CoRR, abs/2210.00305*.

Gokul Yenduri, M Ramalingam, G Chemmalar Selvi, Y Supriya, Gautam Srivastava, Praveen Kumar Reddy Maddikunta, G Deepti Raj, Rutvij H Jhaveri, B Prabadevi, Weizheng Wang, et al. 2022. Gpt (generative pre-trained transformer)–a comprehensive review on enabling technologies, potential applications, emerging challenges, and future directions. *arXiv preprint arXiv:2305.10435*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Xuan Zhang and Wei Gao. 2023. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. *arXiv preprint arXiv:2310.00305*.

Yifan Zhang. 2023. Meta prompting for agi systems. *arXiv preprint arXiv:2311.11482*.

## A Appendix 1: prompt examples

### A.1 2-shot prompt prefix

#### A.1.1 Introduction part

"'You are an AI model that checks biomedical records and corrects existing errors, based only on facts. Your goal is to read the medical record text and decide whether there are any errors. If yes, propose the corrected variant and indicate the error sentence number in the text. Correction is just the entire corrected sentence with NO additional explanations or words.

#### A.1.2 Few-shot examples

- Example 1: Text: "0 A 56-year-old man comes to the physician for a follow-up examination. 1 One month ago, he was diagnosed with a focal seizure, and treatment with a drug that blocks voltage-gated sodium channels was begun. 2 Today, he reports that he has not had any abnormal body movements, but he has noticed occasional double vision. 3 His serum sodium is 132 mEq/L, alanine aminotransferase is 49 U/L, and aspartate aminotransferase is 46 U/L. Medications include carbamazepine." Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

- Example 2: Text: "0 A 53-year-old man comes to the physician because of a 1-day history of fever and chills, severe malaise, and cough with yellow-green sputum. 1 He

works as a commercial fisherman on Lake Superior. 2 Current medications include metoprolol and warfarin. 3 His temperature is 38.5 C (101.3 F), pulse is 96/min, respirations are 26/min, and blood pressure is 98/62 mm 4 Hg. 5 Examination shows increased fremitus and bronchial breath sounds over the right middle lung field. 6 After reviewing imaging, the causal pathogen was determined to be Haemophilus influenzae. 7 An x-ray of the chest showed consolidation of the right upper lobe." Error exists: |||Yes||| Correction: |||After reviewing imaging, the causal pathogen was determined to be Streptococcus pneumoniae.||| Error sentence number: |||6|||

### A.1.3  Rules for the model

Output format if an error exists: Error exist: |||Yes||| Correction: |||<Correction text>||| Error sentence number: |||<Sentence number>|||

Output format if no error is present: Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

Please make sure you complete the objective above with the following rules:

- 1. You should focus on errors in named entities like diagnoses, therapies, and biological species names.

- 2. You must not make things up, you should use only your medical knowledge and medical record data.

- 3. Remember, that you will be rewarded for correct corrections, but also fined for the wrong reports.

- 4. You shouldn't check and correct any spelling errors because only semantic errors are important to you.

- 5. For your convenience, you will see the list of named entities from the record and some info about them.

- 6. You will also see the enumerated sentences from the text - if an error is found, please provide the problematic sentence number.

- 7. Please provide NO explanation for your answer, just give me the error status and error corrections, if any, according to the Output format.

"''

## A.2  5-shot prompt prefix

### A.2.1  Introduction part

"'You are an AI model that checks biomedical records and corrects existing errors, based only on facts. Your goal is to read the medical record text and decide whether there are any errors. If yes, you should propose the corrected variant and indicate the error sentence number in the text. Correction is just the entire corrected sentence with NO additional explanations or words.

### A.2.2  Few-shot examples

- Example 1: Text: "0 A 56-year-old man comes to the physician for a follow-up examination. 1 One month ago, he was diagnosed with a focal seizure, and treatment with a drug that blocks voltage-gated sodium channels was begun. 2 Today, he reports that he has not had any abnormal body movements, but he has noticed occasional double vision. 3 His serum sodium is 132 mEq/L, alanine aminotransferase is 49 U/L, and aspartate aminotransferase is 46 U/L. Medications include carbamazepine." Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

- Example 2: Text: "0 A 53-year-old man comes to the physician because of a 1-day history of fever and chills, severe malaise, and cough with yellow-green sputum. 1 He works as a commercial fisherman on Lake Superior. 2 Current medications include metoprolol and warfarin. 3 His temperature is 38.5 C (101.3 F), pulse is 96/min, respirations are 26/min, and blood pressure is 98/62 mm 4 Hg. 5 Examination shows increased fremitus and bronchial breath sounds over the right middle lung field. 6 After reviewing imaging, the causal pathogen was determined to be Haemophilus influenzae. 7 An x-ray of the chest showed consolidation of the right upper lobe." An error exists: |||Yes||| Correction: |||After reviewing imaging, the causal pathogen was determined to be Streptococcus pneumoniae.||| Error sentence number: |||6|||

- Example 3: Text: "1 He complains of anxiety, nausea, abdominal cramping, vomiting, and diarrhea for three days. 2 He denies smoking, drinking alcohol, and using illicit drugs. 3 He appears restless. 4 His temperature is 37 C (98.6 F), pulse is 110/min, and 5 blood

479

pressure is 150/86 mm 6 Hg. 7 Physical examination shows dilated pupils, diaphoresis, and piloerection. 8 His abdominal exam shows diffuse mild tenderness. 9 There is no rebound tenderness or guarding. 10 Suspected overdose, recommend Naloxone administration. 11 His hemoglobin concentration is 14.5 g/dL 12 , leukocyte count is 8,000/mm, and platelet count is 250,000/mm3; serum studies and urinalysis show no abnormalities." An error exists: |||Yes||| Correction: |||Suspected overdose, recommend methadone administration.||| Error sentence number: |||10|||

- **Example 4**: Text: "0 A potassium hydroxide preparation is conducted on a skin scraping of the hypopigmented area. 1 Patient was treated with topical selenium sulfide based on the microscopy findings. 2 Microscopy of the preparation showed long hyphae among clusters of yeast cells." Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

- **Example 5**: Text: "0 A 56-year-old man comes to the physician for a follow-up examination. 1 One month ago, he was diagnosed with a focal seizure, and treatment with a drug that blocks voltage-gated sodium channels was begun. 2 Medications include phenytoin. 3 Today, he reports that he has not had any abnormal body movements, but he has noticed occasional double vision. 4 His serum sodium is 132 mEq/L, alanine aminotransferase is 49 U/L, and aspartate aminotransferase is 46 U/L." Error exists: |||Yes||| Correction: |||Medications include carbamazepine.||| Error sentence number: |||2|||

### A.2.3 Rules for the model

Output format if an error exists: Error exists: |||Yes||| Correction: |||<Correction text>||| Error sentence number: |||<Sentence number>|||

Output format if no error is present: Error exists: |||No||| Correction: |||None||| Error sentence number: |||None|||

Please make sure you complete the objective above with the following rules:

- **1.** You should focus on errors in named entities like diagnoses, therapies, and biological species names.

- **2.** You must not make things up, you should use only your medical knowledge and medical record data.

- **3.** Remember, that you will be rewarded for correct corrections, but also fined for the wrong reports.

- **4.** You shouldn't check and correct any spelling errors because only semantical errors are important to you.

- **5.** For your convenience, you will see the list of named entities from the record and some info about them.

- **6.** You will also see the enumerated sentences from the text - if an error is found, please provide the problematic sentence number.

- **7.** Please provide NO explanation for your answer, just give me the error status and error corrections, if any, according to the Output format.

,"

## A.3 Ensemble prompt prefix

### A.3.1 Introduction part

"'You are an AI model that checks biomedical records and corrects existing errors, based only on facts. Your goal is to read the medical record text and decide whether there are any errors. If yes, propose the corrected variant and indicate the error sentence number in the text. Correction is just the entire corrected sentence with NO additional explanations or words.

### A.3.2 Few-shot examples

- **Example 1**: Text: "0 A 56-year-old man comes to the physician for a follow-up examination. 1 One month ago, he was diagnosed with a focal seizure, and treatment with a drug that blocks voltage-gated sodium channels was begun. 2 Today, he reports that he has not had any abnormal body movements, but he has noticed occasional double vision. 3 His serum sodium is 132 mEq/L, alanine aminotransferase is 49 U/L, and aspartate aminotransferase is 46 U/L. Medications include carbamazepine." Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

- **Example 2**: Text: "0 A 53-year-old man comes to the physician because of a 1-day history of fever and chills, severe malaise, and cough with yellow-green sputum. 1 He works as a commercial fisherman on Lake Superior. 2 Current medications include metoprolol and warfarin. 3 His temperature is 38.5 C (101.3 F), pulse is 96/min, respirations are 26/min, and blood pressure is 98/62 mm 4 Hg. 5 Examination shows increased fremitus and bronchial breath sounds over the right middle lung field. 6 After reviewing imaging, the causal pathogen was determined to be Haemophilus influenzae. 7 An x-ray of the chest showed consolidation of the right upper lobe." Error exist: |||Yes||| Correction: |||After reviewing imaging, the causal pathogen was determined to be Streptococcus pneumoniae.||| Error sentence number: |||6|||

- **Example 3**: Text: "1 He complains of anxiety, nausea, abdominal cramping, vomiting, and diarrhea for three days. 2 He denies smoking, drinking alcohol, and using illicit drugs. 3 He appears restless. 4 His temperature is 37 C (98.6 F), pulse is 110/min, and 5 blood pressure is 150/86 mm 6 Hg. 7 Physical examination shows dilated pupils, diaphoresis, and piloerection. 8 His abdominal exam shows diffuse mild tenderness. 9 There is no rebound tenderness or guarding. 10 Suspected overdose, recommend Naloxone administration. 11 His hemoglobin concentration is 14.5 g/dL 12 , leukocyte count is 8,000/mm, and platelet count is 250,000/mm3; serum studies and urinalysis show no abnormalities." Error exist: |||Yes||| Correction: |||Suspected overdose, recommend methadone administration.||| Error sentence number: |||10|||

- **Example 4**: Text: "0 A potassium hydroxide preparation is conducted on a skin scraping of the hypopigmented area. 1 Patient was treated with topical selenium sulfide based on the microscopy findings. 2 Microscopy of the preparation showed long hyphae among clusters of yeast cells." Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

- **Example 5**: Text: "0 A 56-year-old man comes to the physician for a follow-up examination. 1 One month ago, he was diag-

nosed with a focal seizure, and treatment with a drug that blocks voltage-gated sodium channels was begun. 2 Medications include phenytoin. 3 Today, he reports that he has not had any abnormal body movements, but he has noticed occasional double vision. 4 His serum sodium is 132 mEq/L, alanine aminotransferase is 49 U/L, and aspartate aminotransferase is 46 U/L." Error exist: |||Yes||| Correction: |||Medications include carbamazepine.||| Error sentence number: |||2|||

### A.3.3 Rules for the model

Output format if an error exists: Error exist: |||Yes||| Correction: |||<Correction text>||| Error sentence number: |||<Sentence number>|||

Output format if no error is present: Error exist: |||No||| Correction: |||None||| Error sentence number: |||None|||

Please make sure you complete the objective above with the following rules:

- 1. You should focus on errors in named entities like diagnoses, therapies, and biological species names.

- 2. You must not make things up, you should use only your medical knowledge and medical record data.

- 3. Remember, that you will be rewarded for correct corrections, but also fined for the wrong reports.

- 4. You shouldn't check and correct any spelling errors because only semantical errors are important to you.

- 5. For your convenience, you will see the list of named entities from the record and some info about them.

- 6. You will see the enumerated sentences from the text - if an error is found, please provide also the problematic sentence number.

- 7. You will also see some possible solutions for this text from the other experts, along with the mean expert trust score for each opinion. You could take expert decisions into account, but with respect to the trust score (higher is better).

- 8. Please provide NO explanation for your answer, just give me the error status and error

corrections, if any, according to the Output format.

''

## B  Appendix 2: Resources used

During the discussed approaches evaluation and prediction making, more than 5,600 API requests were made with 10,537,000 tokens transferred, and the total prediction cost was around $93,6.

# Wonder at Chemotimelines 2024: MedTimeline: An End-to-End NLP System for Timeline Extraction from Clinical Narratives

**Liwei Wang, Qiuhao Lu, Rui Li, Sunyang Fu, and Hongfang Liu**
School of Biomedical Informatics,
The University of Texas Health Science Center at Houston,
Houston, TX, USA
{liwei.wang, qiuhao.lu, rui.li.1, sunyang.fu, hongfang.liu}@uth.tmc.edu

## Abstract

Extracting timeline information from clinical narratives is critical for cancer research and practice using electronic health records (EHRs). In this study, we apply MedTimeline, our end-to-end hybrid NLP system combining large language model, deep learning with knowledge engineering, to the ChemoTimeLine challenge subtasks. Our experiment results in 0.83, 0.90, 0.84, and 0.53, 0.63, 0.39, respectively, for subtask1 and subtask2 in breast, melanoma and ovarian cancer.

## 1 Introduction

Patients' medical history plays a crucial role in guiding the decisions made by clinicians. Yet, the vast majority of temporal information, along with the medical events, is embedded in clinical narratives. For instance, details such as the timing of chemotherapy administration for cancer patients, particularly those referred to the current hospital from other healthcare facilities, are often documented within clinical notes during patient consultations with physicians. There is a pressing need to automatically extract timeline information from clinical narratives to facilitate the understanding of disease progression and treatment efficacy and enhance the quality of cancer research and patient care based on electronic health records (EHRs). Large language models (LLMs), trained on a large amount of unstructured text and then applied to a task through instructive prompts (Tam et al., 2023), have recently shown great value in information extraction and garnered significant attention. We developed MedTimeline, an end-to-end hybrid natural language processing (NLP) system, which combines LLMs and deep learning to support knowledge engineering for timeline information extraction. In this ChemoTimeLine challenge, we applied MedTimeline to the two subtasks and had it evaluated based on the tasks-specific data(Yao et al., 2024).

## 2 Related Work

In the 2012 i2b2 clinical temporal relations challenge, Sohn *et al.* constructed an automated system, i.e., MedTime, that leveraged the framework of HeidelTime, for TIMEX3 extraction from clinical text (Sohn et al., 2013). The system extracts temporal information, including date, time, duration, and frequency, along with their normalized values, demonstrating superior performance. In addition, using the THYME corpus (Styler IV et al., 2014), Liu *et al.* developed an attention-based neural network model to extract containment relations within sentences of clinical narratives (Liu et al., 2019), which outperformed the existing state-of-the-art neural network models at the time.

NLP systems derived from challenges are usually limited to functioning within the confines of the tasks they're specifically designed for. Consequently, Wang *et al.* further expanded their NLP work to patient-level event temporal relation extraction based on real EHR data (Wang et al., 2019). Their results revealed that complete data related to patients' journeys was important for accurate identification of diagnosis dates. In addition, domain knowledge, e.g., chemotherapy drug and transplant names of multiple myeloma, and histology cell type of lung cancer were critical for event temporal relation extraction. In addition, this study demonstrated the usability of MedTime and MedTagger, resource-driven open-source UIMA-based frameworks with the capacity to incorporate knowledge engineering (Sohn et al., 2013; Liu et al., 2013; Wen et al., 2019), for EHR-based cancer research.

## 3 Methodology

In this section, we present our solution, MedTimeline, an end-to-end NLP system comprising an event entity (Chemotherapy entity for subtask 2) extractor, a temporal entity extractor (subtask 2), and a patient-level timeline aggregator (subtasks 1

Figure 1: Architecture of MedTimeline

| w/o Synthetic Data | | | w/ Synthetic Data | | |
|---|---|---|---|---|---|
| Relation | Train | Dev | Relation | Train | Dev |
| **Breast Cancer** | | | **Breast Cancer** | | |
| OPEN | 389 | 133 | OPEN | 389 | 133 |
| CONTAINS | 298 | 57 | CONTAINS | 492 | 57 |
| BEGINS-ON | 131 | 27 | BEGINS-ON | 231 | 27 |
| ENDS-ON | 26 | 29 | ENDS-ON | 225 | 29 |
| **Melanoma** | | | **Melanoma** | | |
| OPEN | 35 | 192 | OPEN | 35 | 192 |
| CONTAINS | 37 | 157 | CONTAINS | 37 | 157 |
| BEGINS-ON | 10 | 42 | BEGINS-ON | 205 | 42 |
| ENDS-ON | 1 | 2 | ENDS-ON | 191 | 2 |
| **Ovarian Cancer** | | | **Ovarian Cancer** | | |
| OPEN | 338 | 226 | OPEN | 338 | 226 |
| CONTAINS | 327 | 140 | CONTAINS | 516 | 140 |
| BEGINS-ON | 98 | 34 | BEGINS-ON | 266 | 34 |
| ENDS-ON | 59 | 52 | ENDS-ON | 256 | 52 |

Table 1: Dataset statistics with and without synthetic data.

and 2). The architecture of MedTimeline includes two well-established knowledge engineering NLP pipelines (MedTagger as event extractor and MedTime as temporal expression extractor) from the Open Health NLP (OHNLP) consortium, a context-aware deep learning open-source architecture, and an LLM-empowered data augmentation pipeline (Figure 1). Specifically, the data augmentation pipeline incorporates ChatGPT to generate synthetic data to facilitate the fine-tuning of a pre-trained language model for temporal relation classification within the timeline aggregator.

### 3.1 Event Entity Extractor

MedTimeline leverages MedTagger for event entity extraction. Particularly, the knowledge artifacts of chemotherapy drug names for breast cancer, ovarian cancer and melanoma, are first collected from both the training data set and the online knowledge hub of the American Cancer Society[1], and then made into a MedTimeline rule set that is compatible with MedTagger.

### 3.2 Temporal Entity Extractor

MedTime and MedTagger function as the temporal entity extractors in the MedTimeline to automatically extract temporal information from clinical notes. For MedTime, missing temporal expression

---

[1] https://www.cancer.org/

rules are added to MedTime through the comparison of the results automatically extracted by MedTime (existing rules) with the gold standards of the training set, i.e., subtask1 in this study. For instance, we add "at this time" and "on the day" rules into MedTime. Additionally, we leverage MedTagger to manage complex rules that can not be added to MedTime, in order to extract the temporal information not captured by MedTime. For example, MedTime failed to extract "today" when it was preceded by a number, e.g., "5 today". We made a regular expression rule for this case to enable automatic extraction of "today".

### 3.3 Synthetic Data Augmentation

Training data insufficiency and imbalance are critical issues as they may impact the quality and reliability of predictive models (Lu et al., 2021). To address these issues, MedTimeline synthesizes artificial data to enrich the training data and facilitate model training. Essentially, ChatGPT-4 (i.e., `gpt-4-1106-preview`) prompting is used to generate synthetic data.

In the context of the challenge subtasks, we identified the lack of sufficient data for such a condition as melanoma, and imbalance of the datasets concerning the three temporal relations during the initial data analysis. We instruct ChatGPT-4 to produce artificial data, as shown in Table 1. Specifically, textual segments extracted between chemotherapy events and time expressions demonstrate a unique pattern for each predefined tem-

poral relation as well as each cancer type, e.g., `BEGINS-ON` of melanoma is substantially different from `ENDS-ON` of breast cancer. Following the patterns, we manually design 5 example text pieces for each temporal relation of each cancer type to use as few-shot demonstrations. Notably, we only synthesize textual segments connecting chemo and time instead of the entire clinical note, and their numbers are determined based on preliminary experiments. We use the following prompt:

*You are a helpful assistant in synthetic data generation. Your job is to generate a sentence containing a chemotherapy entity for melanoma (source) and a TIMEX3 entity (target). The relation between them is ENDS-ON. After reading and comprehending the examples, generate 50 data samples. The outputs should be in three columns: source, target and context. Use | as the delimiter and do not add index numbers to the generated samples. Be diverse, representative, and accurate, e.g., the chemo should be for the specific cancer and do not mention the specific cancer in the sentence. Examples: [manually-designed 5-shot demonstrations] Generated Data:*

### 3.4 Relation Extraction

We cast relation extraction for the medical events and temporal expressions as a multi-class text classification problem. Essentially, we extract the textual segment (e.g., *"Chemo started Today."*) that links a chemotherapy event (e.g., *chemo*) with its related time expression (e.g., *today*) from the clinical note. We then categorize the textual segment into one of the predefined temporal relations.

The problem is also an open-world classification problem (Bai et al., 2022), as it requires the model to predict a sample as `OPEN`, which indicates it has an open/unspecified temporal relation or does not have any relation. To create the corresponding training data for this category, we adopt a simple yet effective negative sampling strategy where we extract <*chemo, time*> pairs in the training set whose distance is less than 250 characters[2] and do not belong to any of the predefined temporal relations. We consider such negative samples to be hard and realistic. It is worth noting that since candidate relations are not provided in the test set, we use the same strategy for candidate search during inference.

Formally, given a clinical note $S$ containing a list of chemotherapy events $\mathcal{C} = \{c_1, c_2, \ldots, c_{|\mathcal{C}|}\}$ and a list of time expressions $\mathcal{T} = \{t_1, t_2, \ldots, t_{|\mathcal{T}|}\}$,

we search for candidate pairs using the aforementioned strategy and extract the text between them as input $\mathcal{D} = \{x_i, x_2, \ldots, x_k, \ldots, x_{|\mathcal{D}|}\}$ where $x_k$ is the text between $c_i \in \mathcal{C}$ and $t_j \in \mathcal{T}$. The objective is to predict the corresponding label $y_k \in \mathcal{E}$ where $\mathcal{E} = \{\text{CONTAINS-1}, \text{BEGINS-ON}, \text{ENDS-ON}, \text{OPEN}\}$.

In particular, we use the bio-lm[3] pre-trained language model (Lewis et al., 2020) to encode the text and feed the representation for the `[CLS]` token in the last layer into a linear layer for classification. The model is optimized with cross-entropy loss:

$$\mathcal{L} = -\sum_{l=1}^{4} y_l \log \hat{y}_l \tag{1}$$

where $y_l$ is the ground-truth label and $\hat{y}_l$ refers to the output prediction probabilities.

### 3.5 Time Expression Normalization

We adapt MedTime[4] to convert temporal expression from clinical notes into standardized TIMEX3 format. Types of MedTime output include standard dates and time intervals. For time entities which are directly mapped into standard dates such as *2013-11-12* and *2012-W06*, the MedTime output is used as the standardized TIMEX3 date. For time entities which are mapped into time intervals, the standardized TIMEX3 date is calculated by subtracting time intervals from the principal date.

### 3.6 Patient-level Timeline Aggregation

If the relation of the pair is classified as `OPEN` by bio-lm, we do not assign any specific temporal relation for the pair. We then employ the aforementioned temporal expression normalization method to convert the temporal entities into standardized TIMEX3 format. At last, we aggregate all <*chemo, time*> pairs whose relation are not `OPEN` to construct the patient-level timeline.

## 4 Experiments

### 4.1 Results

In this section, we first show the statistics of the dataset with and without synthetic data in Table 1. We then present the temporal relation classification results across different models and cancers in Table 2. Finally, we show the patient-level timeline extraction results on the dev and test sets, as shown

---

[2]Maximum distance among <*chemo, time*> pairs of a predefined temporal relation in the training set.

[3]We use the best-performing variant RoBERTa-large-PM-M3-Voc across all experiments in this work.

[4]https://github.com/OHNLP/MedTime

| Cancers | Models | CONTAINS | | | BEGINS-ON | | | ENDS-ON | | | OPEN | | | Overall | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | Acc | P | R | F1 |
| Breast | PubMedBERT-tlink | 0.829 | 0.509 | 0.630 | 0.556 | 0.741 | 0.635 | 0.833 | 0.345 | 0.488 | 0.912 | 0.857 | 0.884 | 0.703 | 0.844 | 0.703 | 0.751 |
| | BioClinicalBERT | 0.636 | 0.860 | 0.731 | 0.870 | 0.741 | 0.800 | 0.727 | 0.276 | 0.400 | 0.956 | 0.970 | 0.963 | 0.837 | 0.845 | 0.837 | 0.825 |
| | BioClinicalBERT* | 0.831 | 0.860 | 0.845 | 0.828 | 0.889 | 0.857 | 0.917 | 0.759 | 0.830 | 0.948 | 0.955 | 0.951 | 0.902 | 0.904 | 0.902 | 0.902 |
| | bio-lm | 0.773 | 0.895 | 0.829 | 0.920 | 0.852 | 0.885 | 0.769 | 0.345 | 0.476 | 0.901 | 0.962 | 0.931 | 0.862 | 0.858 | 0.862 | 0.849 |
| | bio-lm* | 0.817 | 0.860 | 0.838 | 0.897 | 0.963 | 0.929 | 0.909 | 0.690 | 0.784 | 0.978 | 0.993 | 0.985 | 0.923 | 0.923 | 0.923 | 0.921 |
| Melanoma | PubMedBERT-tlink | 0.617 | 0.586 | 0.601 | 0.914 | 0.762 | 0.831 | 0.000 | 0.000 | 0.000 | 0.730 | 0.745 | 0.737 | 0.679 | 0.701 | 0.679 | 0.689 |
| | BioClinicalBERT | 0.479 | 0.994 | 0.646 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.985 | 0.344 | 0.510 | 0.565 | 0.672 | 0.565 | 0.507 |
| | BioClinicalBERT* | 0.580 | 0.949 | 0.720 | 0.035 | 0.048 | 0.040 | 0.000 | 0.000 | 0.000 | 0.948 | 0.380 | 0.543 | 0.570 | 0.698 | 0.570 | 0.557 |
| | bio-lm | 0.596 | 0.968 | 0.738 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.957 | 0.688 | 0.800 | 0.723 | 0.705 | 0.723 | 0.686 |
| | bio-lm* | 0.569 | 0.949 | 0.711 | 0.925 | 0.881 | 0.902 | 0.000 | 0.000 | 0.000 | 0.923 | 0.438 | 0.594 | 0.687 | 0.777 | 0.687 | 0.671 |
| Ovarian | PubMedBERT-tlink | 0.807 | 0.507 | 0.623 | 0.392 | 0.588 | 0.471 | 0.435 | 0.192 | 0.267 | 0.942 | 0.929 | 0.935 | 0.688 | 0.800 | 0.688 | 0.727 |
| | BioClinicalBERT | 0.750 | 0.879 | 0.809 | 0.615 | 0.471 | 0.533 | 0.895 | 0.327 | 0.479 | 0.918 | 0.987 | 0.951 | 0.839 | 0.840 | 0.839 | 0.821 |
| | BioClinicalBERT* | 0.774 | 0.907 | 0.836 | 0.800 | 0.588 | 0.678 | 0.920 | 0.442 | 0.597 | 0.929 | 0.978 | 0.953 | 0.865 | 0.870 | 0.865 | 0.855 |
| | bio-lm | 0.703 | 0.879 | 0.781 | 0.667 | 0.588 | 0.625 | 0.905 | 0.365 | 0.521 | 0.951 | 0.951 | 0.951 | 0.834 | 0.848 | 0.834 | 0.824 |
| | bio-lm* | 0.778 | 0.850 | 0.812 | 0.800 | 0.706 | 0.750 | 0.906 | 0.558 | 0.690 | 0.920 | 0.965 | 0.942 | 0.863 | 0.865 | 0.863 | 0.858 |

Table 2: Temporal relation classification performance across different models and cancers with relation-wise and overall scores. * represents fine-tuning with synthetic data.

in Table 3. All experimental results are obtained during the challenge.

We use three pre-trained language models in the clinical domain as baselines, i.e., PubMedBERT-tlink[5], BioClinicalBERT (Alsentzer et al., 2019), and bio-lm (Lewis et al., 2020). Note that we do not fine-tune PubMedBERT-tlink as it is already trained on a similar task and data. For relation classification, we use precision (P), recall (R), and F1-score as the metrics. For patient-level timeline extraction, we use the official script of the challenge where the *relaxed-to-month* F1-score is used as the metric. One key observation is that both BioClinicalBERT and bio-lm demonstrate a significant improvement with synthetic training data, highlighting the effectiveness of data augmentation in this context. All models struggle with `ENDS-ON` for Melanoma even after training data is augmented from 1 to 191. The reason lies in the fact that there are very limited data samples in the dev set, i.e., only 2 samples in the dev set as shown in Table 1.

## 4.2 Error analysis

We compare the patient-level chemo-timeline generated by MedTimeline with the gold standard of dev set to identify errors from our system. The errors mainly originate from two sources, i.e., time normalization and relation classification. The former is caused by wrong anchor time retrieved from MedTime and inaccurate imputation of the incomplete time entity. The latter arises from incomplete and complex text input. Incomplete text input is caused by our strategy of merely extracting the text between the chemo entity and time entity, leading

| Subtask | Split | Breast | Melanoma | Ovarian |
|---|---|---|---|---|
| Subtask 1 | Dev | 0.86 | 0.80 | 0.77 |
| | Test | 0.83 | 0.90 | 0.84 |
| Subtask 2 | Dev | 0.83 | 0.71 | 0.75 |
| | Test | 0.53 | 0.63 | 0.39 |

Table 3: Patient-level timeline evaluation results for Subtasks 1 and 2.

to the missingness of some useful information. For example, the original text in clinical notes is *'She received her 9th and final dose of IL2 at 9/22'*[6], and the timeline in the gold annotation is *["il2", "ends-on", "2012-09-22"]*. However, by extracting *'IL2 at 9/22'* as input, our system wrongly classifies the relation as `BEGINS-ON`. Meanwhile, some text input is too complex for the system to classify the correct relation. For example, given the original text *'Today he feels well. He had been able to control the symptoms of nausea that he has experienced with his TCH chemotherapy'*, our system wrongly classifies the relation as `ENDS-ON` while the relation should be `OPEN`.

## 5 Conclusion

We present MedTimeline, an end-to-end hybrid NLP system generalizable to any medical events for patients' timeline extraction, and evaluate it based on the ChemoTimeLine challenge data. Our system ranks the second place in subtask 1 and the third place in subtask 2. In the future, we will continue to develop the MedTimeline, and tailor it to the scenarios of various medical events.

---

[5]https://huggingface.co/HealthNLP/pubmedbert_tlink

[6]All examples are rephrased in order to avoid data leakage

## References

Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Ke Bai, Guoyin Wang, Jiwei Li, Sunghyun Park, Sungjin Lee, Puyang Xu, Ricardo Henao, and Lawrence Carin. 2022. Open world classification with adaptive negative samples. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4378–4392, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoyanov. 2020. Pretrained language models for biomedical and clinical tasks: Understanding and extending the state-of-the-art. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 146–157, Online. Association for Computational Linguistics.

Hongfang Liu, Suzette J Bielinski, Sunghwan Sohn, Sean Murphy, Kavishwar B Wagholikar, Siddhartha R Jonnalagadda, KE Ravikumar, Stephen T Wu, Iftikhar J Kullo, and Christopher G Chute. 2013. An information extraction framework for cohort identification using electronic health records. *AMIA Summits on Translational Science Proceedings*, 2013:149.

Sijia Liu, Liwei Wang, Vipin Chaudhary, and Hongfang Liu. 2019. Attention neural model for temporal relation extraction. In *Proceedings of the 2nd clinical natural language processing workshop*, pages 134–139.

Qiuhao Lu, Dejing Dou, and Thien Huu Nguyen. 2021. Textual data augmentation for patient outcomes prediction. In *2021 IEEE international conference on bioinformatics and biomedicine (BIBM)*, pages 2817–2821. IEEE.

Sunghwan Sohn, Kavishwar B Wagholikar, Dingcheng Li, Siddhartha R Jonnalagadda, Cui Tao, Ravikumar Komandur Elayavilli, and Hongfang Liu. 2013. Comprehensive temporal information detection from clinical text: medical events, time, and tlink identification. *Journal of the American Medical Informatics Association*, 20(5):836–842.

William F Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C De Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, et al. 2014. Temporal annotation in the clinical domain. *Transactions of the association for computational linguistics*, 2:143–154.

Derek Tam, Anisha Mascarenhas, Shiyue Zhang, Sarah Kwan, Mohit Bansal, and Colin Raffel. 2023. Evaluating the factual consistency of large language models through news summarization. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5220–5255.

Liwei Wang, Jason Wampfler, Angela Dispenzieri, Hua Xu, Ping Yang, and Hongfang Liu. 2019. Achievability to extract specific date information for cancer research. In *AMIA Annual Symposium Proceedings*, volume 2019, page 893. American Medical Informatics Association.

Andrew Wen, Sunyang Fu, Sungrim Moon, Mohamed El Wazir, Andrew Rosenbaum, Vinod C Kaggal, Sijia Liu, Sunghwan Sohn, Hongfang Liu, and Jungwei Fan. 2019. Desiderata for delivering nlp to accelerate healthcare ai advancement and a mayo clinic nlp-as-a-service implementation. *NPJ digital medicine*, 2(1):130.

Jiarui Yao, Harry Hochheiser, WonJin Yoon, Eli Goldner, and Guergana Savova. 2024. Overview of the 2024 shared task on chemotherapy treatment timeline extraction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*. NAACL.

# Edinburgh Clinical NLP at MEDIQA-CORR 2024: Guiding Large Language Models with Hints

**Aryo Pradipta Gema[1][*]**   **Chaeeun Lee[1][*]**   **Pasquale Minervini[1]**   **Luke Daines[2]**
**T. Ian Simpson[1]**   **Beatrice Alex[3,4]**

[1]School of Informatics, University of Edinburgh    [2]Usher Institute, University of Edinburgh
[3]Edinburgh Futures Institute, University of Edinburgh
[4]School of Literatures, Languages and Cultures, University of Edinburgh
{aryo.gema, chaeeun.lee, p.minervini, luke.daines}@ed.ac.uk
{ian.simpson, b.alex}@ed.ac.uk

## Abstract

The MEDIQA-CORR 2024 shared task aims to assess the ability of Large Language Models (LLMs) to identify and correct medical errors in clinical notes. In this study, we evaluate the capability of general LLMs, specifically GPT-3.5 and GPT-4, to identify and correct medical errors with multiple prompting strategies. Recognising the limitation of LLMs in generating accurate corrections only via prompting strategies, we propose incorporating error-span predictions from a smaller, fine-tuned model in two ways: 1) by presenting it as a hint in the prompt and 2) by framing it as multiple-choice questions from which the LLM can choose the best correction. We found that our proposed prompting strategies significantly improve the LLM's ability to generate corrections. Our best-performing solution with 8-shot + CoT + hints ranked sixth in the shared task leaderboard. Additionally, our comprehensive analyses show the impact of the location of the error sentence, the prompted role, and the position of the multiple-choice option on the accuracy of the LLM. This prompts further questions about the readiness of LLM to be implemented in real-world clinical settings.[1]

## 1 Introduction

Medical errors represent a major concern in the healthcare sector, leading to adverse patient outcomes and higher costs for healthcare providers. The detection and correction of such medical errors are critical in enhancing healthcare delivery and outcomes. Recognising the importance of efficient and precise medical documentation, the MEDIQA-CORR 2024 shared task (Ben Abacha et al., 2024a) is initiated to evaluate the potential of using Large Language Models (LLMs) as solutions to locate and correct medical errors within clinical notes.

In our study, we evaluated multiple prompting strategies such as In-context Learning (ICL) and Chain-of-Thought (CoT) to enhance the performance of LLMs, specifically focusing on GPT-3.5 and GPT-4 (OpenAI, 2023). We proposed incorporating a smaller fine-tuned language model, namely BioLinkBERT (Yasunaga et al., 2022), to aid LLMs in locating an error span in a clinical note. We incorporated the predicted error span in two ways: 1) by presenting it as a hint in the prompt to direct the error correction, and 2) by framing it as multiple-choice questions where the LLM can select the most probable correction.

Our findings revealed that the LLMs show noticeable improvements in their generation capability when presented with more ICL examples. Similarly, the CoT prompt also improves the error correction capability of the LLMs. Among the different reasoning styles we experimented with, the LLM performs the best with brief reasoning. Our prompt design, which provides a hint about the typical nature of the errors and a hint from the error span prediction, further improves the LLMs' ability to generate corrections. The combination of 8-shot ICL with Brief CoT reasoning and hints is the best-performing prompting strategy in the two provided validation sets. This pipeline ranked sixth in the shared task leaderboard. In summary, our contributions are as follows:

- A comprehensive analysis of the impact of ICL on the performance of LLMs for medical error correction.

- An extensive exploration of CoT to inject various reasoning styles into the LLM and their impact on the performance.

- Novel approaches to integrate the predictions of a smaller language model into the LLM generation.

- Sensitivity analyses of LLMs, highlighting how minor variations such as the error sentence loca-

---

[*]Equal contribution.
[1]Our code is available at https://github.com/aryopg/mediqa

| Category | Train | | Valid | | Test | |
|---|---|---|---|---|---|---|
| | MS | UW | MS | UW | MS | UW |
| No Error | 970 | 0 | 255 | 80 | - | - |
| Contain Error | 1,219 | 0 | 319 | 80 | - | - |
| **Total** | **2,189** | **0** | **574** | **160** | **597** | **328** |

Table 1: Dataset statistics of each split, categorised by the source and presence of a medical error.

tion, the prompted role, and the multiple-choice positioning can influence generation capabilities.

## 2 Background

### 2.1 Task Description

MEDIQA-CORR 2024 task (Ben Abacha et al., 2024b) comprises three sub-tasks, each addressing a different aspect of medical error correction:

**Binary classification:** Detecting whether the clinical note contains a medical error.

**Span Identification:** Identifying the text span associated with a medical error if it exists.

**Natural Language Generation:** Generating a correction if a medical error exists.

Table 1 shows the statistics for each data split, organised by the source of the data and whether or not it contains a medical error. Each clinical note contains either one or no medical error.

The task uses accuracy for binary classification and span identification. The generated correction is evaluated using an aggregate Natural Language Generation (NLG) score, combining ROUGE-1 (Lin, 2004), BERTScore (Zhang et al., 2020), and BLEURT (Sellam et al., 2020), which is best aligned with human judgement, among other NLG metrics (Ben Abacha et al., 2023).

### 2.2 Related work

LLMs have shown remarkable capabilities in many NLP tasks, including in the clinical domain. Liévin et al. (2022) evaluated LLMs with various prompting strategies, showing LLMs' capability to answer complex medical questions. Falis et al. (2024) uses GPT-3.5 to generate accurate synthetic discharge summaries by prompting it with a list of diagnoses. Gema et al. (2024) also shows GPT-4 in zero-shot setting outperforms other fine-tuned LLMs in a natural language inference task for clinical trial data.

However, despite the increasing use of general LLMs, their performance varies widely depending on the nature of the task. For instance,

fine-tuned smaller encoder-based models (e.g., BioLinkBERT) still maintain the lead in tasks such as medical entity recognition (Kim et al., 2023). Gema et al. (2023) showed that domain-adapted LLaMA (Touvron et al., 2023) outperforms the state-of-the-art models in clinical outcome prediction tasks. Such studies show that fine-tuned models are still preferable, especially in discriminative tasks such as classification and entity recognition.

In this study, we seek to combine the generative capability of LLMs with the discriminative capability of a smaller fine-tuned language model. We compared our novel method with solutions that rely solely on prompting strategies (i.e., ICL and CoT).

## 3 System Overview

We experimented with three strategies:

**End-to-end Prompting Strategy for Error Correction:** This strategy treats all three subtasks as a single prompting task. The LLM simultaneously predicts if the clinical note contains an error, pinpointing its location, and proposing a correction.

**Fine-tuning Error Span Prediction and MCQ-style Error Correction:** This method splits the task into error span prediction and correction. It uses a fine-tuned model for error span prediction, followed by MCQ-style prompts for correction.

**Hybrid Approach:** As shown in Figure 1, This approach uses error span predictions from a fine-tuned model as correction hints injected into the end-to-end prompting strategy. This is our best-performing strategy in both validation and test sets.

The following sections outline the details for the **Error Span Prediction** and **Error Correction**.

### 3.1 Error Span Prediction

We noticed that medical errors appear predominantly in the form of diagnoses or treatments, instead of the patient's factual information. This finding motivated us to fine-tune an encoder model to first detect an error span within the clinical note.

We trained BioLinkBERT and BERT[2] using a question-answering pipeline adapted from the Stanford Question Answering Dataset (SQuAD). We pre-processed the training and validation sets to align them with the SQuAD v1 format, which assumes that there is always an error span in the input. We introduced a template question, "Which part in the given clinical note is clinically incorrect?" in the question column of the SQuAD format. The

---

[2]Both base and large versions of the models

Figure 1: Schema of our best-performing strategy with In-Context Learning (ICL) and Chain-of-Thought (CoT) prompting strategies. The strategy involves fine-tuning BioLinkBERT on the training set for error span prediction. Then, we prompt GPT-3.5 with various reasoning templates to reason pairs of clinical notes and ground truth corrections to gather ICL examples with CoT reasons. Subsequently, this strategy leverages the ICL examples and error span predictions as a hint.

trained model predicts the start and end indices, which indicate the position of the predicted error span in the text.

We trained and evaluated the error span prediction models only on clinical notes that contained errors. We evaluated the models using exact match (EM) and token-based F1 score metrics, using the latter to choose the best checkpoint.

## 3.2 Error Correction

We experimented with GPT-3.5 and GPT-4 for the error correction step. We prompted the LLMs to return the outputs in JSON format for ease of post-processing. In rare cases where the outputs are not JSON-parsable, we default the prediction as if no error was found. We integrated the error span prediction to this error correction step in two ways:

### 3.2.1 Multiple-Choice Question prompt

As shown in Figure 2, this strategy involves two interactions with the LLM: 1) to construct an options

set and 2) to ask a multiple-choice question.

In the first interaction, the model generates potential replacement options for the identified error span. Here, the predicted error span is replaced with a placeholder *"<BLANK>"*, and the LLM is tasked with generating $n$ replacement candidates. During our experiments, we observed a pattern where the model often included the predicted error span or its synonyms in the options. To eliminate this redundancy, we added a directive prompt *"Do not include the <predicted_error_span> or its medical synonyms in your answer"*.

In the second interaction, we query the LLM with an MCQ-style prompt, which presents the full clinical note, with the predicted error span replaced by *"<BLANK>"*, and the options comprised of $n$ LLM-generated options from the first interaction and the predicted error span (totalling $n + 1$ options). The LLM chooses the best correction among these options. Subsequently, we derive the error flag classification based on the LLM's re-

490

Figure 2: Schema of the Multiple-Choice Question prompt strategy.

sponse, 0 if it selects the predicted error span as the correct answer, or 1 if the model selects one of the other choices. We experimented with varying the number of answer choices to two and four options.

### 3.2.2 Hybrid Approach

As illustrated in Figure 1, the pipeline continues with the preparation of the ICL examples after the training for the error span prediction. For solutions that rely only on ICL examples and do not require CoT reasoning, we directly retrieve pairs of clinical notes and their respective ground-truth corrections as ICL examples. In contrast, CoT-based solutions require ICL examples with reasons provided. Inspired by He et al. (2023), we prompted GPT-3.5 (`gpt-3.5-turbo-0613`) to generate a reasoning for the ICL examples. We selected GPT-3.5 particularly because of its generation capability and clinical knowledge (Gema et al., 2024).

We experimented with three CoT reasoning templates: **Brief**, **Long**, and **SOAP**. All reasoning templates require the model to reason the ground-truth correction by identifying the incorrect span and providing the reasoning behind it. However, each format provides a different depth and structure of reasoning. The **Brief CoT** template prompts concise reasoning, the **Long CoT** template requires detailed step-by-step explanations, and the **SOAP CoT** template organises information according to Subjective, Objective, Assessment, and Plan sections before making corrections.

During inference, the solution uses a selected reasoning format with ICL examples to correct clinical notes. The model applies a reasoning strat-

egy to new scenarios based on the reasoned ICL examples which are retrieved using the BM25 algorithm (Robertson et al., 1995), selecting examples similar to the clinical note in question. We also integrate a hint about the typical nature of the errors, focusing the model's attention on specific biomedical entities such as diagnoses and treatments (i.e., *"Pay special attention to biomedical entities such as chief complaints, medical exams, diagnoses, and treatments."*). We denote this as **"Type hint"**. Finally, we leverage the error span prediction by adding it as another hint, denoted as **"Span hint"** (i.e., *"A clinician said that you MAY want to pay attention to the mention of <predicted_error_span>"*).

## 4 Results

Our experiments are structured as answers to sequential research questions. Firstly, we conducted experiments to find the best model for error span prediction, evaluating them on EM and F1 scores. Subsequently, we experimented with various prompting strategies for error correction, evaluating them on the macro-averaged accuracy and aggregate NLG scores across MS and UW datasets. The first error correction experiment starts with an end-to-end prompting approach, relying solely on the LLM capability with ICL and CoT to correct errors. We, then, experimented with integrating the error span prediction model into the error correction process via the MCQ-style prompt. Lastly, we experimented with the hybrid approach, integrating the error span prediction as a hint for the end-to-end prompting approach. We used GPT-3.5 in our error correction experiments on the validation sets[3], choosing the best prompting strategy to be implemented with GPT-4 on the test set.

**RQ1: How well are the smaller LMs performing in the error span detection?**

As shown in Table 2, we experimented with general (i.e., BERT-base and -large) and domain-adapted models (i.e., BioLinkBERT-base and -large) for the error span prediction. We evaluated the models exclusively on a subset of the validation set that contains a medical error as stated in Subsection 3.1.

Among all models, BioLinkBERT-large showed the highest EM and F1 scores on the MS validation set, indicating a superior ability to predict error spans within clinical notes. This suggests that

---

[3]Due to a limited research budget.

| Model | MS | | UW | |
|---|---|---|---|---|
| | **EM** | **F1** | **EM** | **F1** |
| BERT-base | 54.86 | 80.09 | 1.25 | 4.44 |
| BERT-large | 55.17 | 79.30 | 5.00 | 7.92 |
| BioLinkBERT-base | 55.17 | 81.33 | **6.25** | **12.29** |
| BioLinkBERT-large | **58.31** | **82.49** | **6.25** | 8.91 |

Table 2: Performance of fine-tuned error span prediction models. **Bold cell** indicates the highest score for the metric.

| # shots | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| 2 | 0.5089 | 0.3348 | 0.4139 |
| 4 | 0.5242 | 0.4215 | 0.4503 |
| 8 | **0.5268** | **0.4526** | **0.5038** |

Table 3: Performance of GPT-3.5 using different numbers of ICL examples on validation sets. **Bold cell** indicates the highest score for the metric.

the domain-adaptive pretraining that BioLinkBERT has undergone contributes to its performance in medical error detection tasks. However, all models struggle to accurately predict error spans on the UW validation set. Recognising this, we trained BioLinkBERT-large on the MS train dataset and 25% of the UW validation dataset as the error span prediction model for the subsequent experiments.

## RQ2: Can LLMs perform well end-to-end solely with prompting strategies?

Before leveraging the error span prediction, we began our error correction experiment by solely relying on the LLM with prompting strategies to correct errors without any help from the error span prediction. This prompt-only end-to-end approach serves as the baseline for our proposed solutions.

## RQ2.1: Do more ICL examples improve the LLM's performance?

Firstly, we experimented with varying the number of ICL examples on GPT-3.5's performance across MS and UW validation sets. We did not report 0-shot performance as the LLM failed to generate a parseable answer, indicating that the LLM failed to complete the task without any examples. As shown in Table 3, we observe a trend where the performance of the LLM improves in all metrics as the number of shots increases, with the 8-shot setting performing the best. Our subsequent experiments will use the 8-shot ICL setup.

| Type Hint | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| ✗ | **0.5527** | 0.4472 | 0.4467 |
| ✓ | 0.5268 (-0.03) | **0.4526** (+0.01) | **0.5038** (+0.06) |

Table 4: Performance of GPT-3.5 using 8-shot prompt with or without a type hint on validation sets. Values in parentheses indicate the performance difference against the LLM that does not receive a type hint. cyan indicates improvement, red indicates decrease. **Bold cell** indicates the highest score for the metric.

| CoT | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| None | 0.5268 | 0.4526 | 0.5038 |
| Brief | 0.5866 (+0.06) | **0.4989** (+0.05) | **0.5389** (+0.04) |
| Long | **0.6074** (+0.08) | 0.4717 (+0.02) | 0.4930 (-0.01) |
| SOAP | 0.5186 (-0.01) | 0.4058 (-0.05) | 0.4228 (-0.08) |

Table 5: Performance of GPT-3.5 using 8-shot and type hint prompt with various CoT formats on validation sets. Values in parentheses indicate the performance difference against the LLM that does not use CoT reasoning. cyan indicates improvement, red indicates decrease. **Bold cell** indicates the highest score for the metric.

## RQ2.2: Adding a hint about the typical error

In our first experiment, we observed that the LLMs tend to correct non-essential errors (e.g., grammatical and unit errors). Thus, we prompted the LLM with a hint about the typical form of the errors (i.e., *"Pay special attention to biomedical entities such as chief complaints, medical exams, diagnoses, and treatments."*). Table 4 shows the performance comparison between a prompt with and without this hint. When a hint is provided, there is a decrease in the error flag accuracy by 0.03 which may indicate that there are medical errors that are not one of the specified biomedical entities. However, this is compensated by improvements in both sentence ID accuracy and the aggregate NLG score, with the latter seeing a notable increase of 0.06. This indicates that while the hint may slightly hinder the model's binary classification ability, it correctly directs the focus of the LLM in locating the error.

## RQ2.3: Chain-of-Thought with various formats

Table 5 evaluates the effect of different Chain-of-Thought (CoT) formats on GPT-3.5's performance. The absence of CoT (None) serves as a baseline against which the Brief, Long, and SOAP formats are compared. The Brief CoT format leads to improvements across all metrics, particularly in sentence ID accuracy and the aggregate NLG score, underscoring the benefit of concise, targeted rea-

| Prompting Strategy | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| 8-shot + Brief CoT | 0.5866 | 0.4989 | 0.5389 |
| MCQ (2 options) | **0.6131** | **0.6029** | **0.6492** |
| MCQ (4 options) | 0.6087 | 0.5944 | 0.6448 |

Table 6: Performance of GPT-3.5 with the MCQ-style prompt on validation sets. **Bold cell** indicates the highest score for the metric.

| CoT | Span Hint | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|---|
| MCQ (2 opt) | ✓ | **0.6131** | **0.6029** | 0.6492 |
| MCQ (4 opt) | ✓ | 0.6087 | 0.5944 | 0.6448 |
| None | ✗ | 0.5268 | 0.4526 | 0.5038 |
| | ✓ | 0.5671 (+0.04) | 0.5543 (+0.10) | 0.7348 (+0.23) |
| Brief | ✗ | 0.5866 | 0.4989 | 0.5389 |
| | ✓ | 0.5610 (-0.03) | 0.5454 (+0.05) | **0.7385** (+0.20) |
| Long | ✗ | 0.6074 | 0.4717 | 0.4930 |
| | ✓ | 0.6048 (-0.00) | 0.4651 (-0.01) | 0.4822 (-0.01) |
| SOAP | ✗ | 0.5186 | 0.4058 | 0.4228 |
| | ✓ | 0.5237 (+0.01) | 0.4310 (+0.03) | 0.4884 (+0.07) |

Table 7: Performance of GPT-3.5 using 8-shot and type hint prompt with various CoT format and with or without receiving span hint on validation sets. Values in parentheses indicate the performance difference against the solution that does not receive a span hint. cyan indicates improvement, red indicates decrease. **Bold cell** indicates the highest score for the metric.

| Prompting Strategy | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| 8-shot + Hints | 0.5243 | 0.4649 | 0.6274 |
| **8-shot + Brief CoT + Hints** | **0.6681** | 0.5924 | **0.6634** |
| MCQ (2 options) | 0.6573 | **0.5957** | 0.6267 |
| MCQ (4 options) | 0.5935 | 0.5232 | 0.5882 |

Table 8: Results of GPT-4 with either ICL + CoT + hinted prompt or Multiple-Choice-Question prompt on test sets. The models are compared based on the aggregate NLG score.

soning in enhancing model performance. The Long format, while offering the highest accuracy in error flagging, exhibits a decrease in the aggregate score, suggesting that excessive detail may detract from overall correction quality. Conversely, the SOAP format results in declines across all metrics, highlighting that detailed and structured reasoning approaches may not necessarily be beneficial and may even hinder the model's effectiveness.

### RQ3: Can LLMs perform if provided with a span hint?

After the experiments with different prompting setups, we experimented with integrating the error span prediction into the error correction process.

### RQ3.1: Can LLMs perform better with MCQ-style prompts?

As shown in Table 6, MCQ-style prompt using error span prediction improved performance over end-to-end systems. This can be attributed to two reasons. First, the MCQ-style prompt provides options that match the specificity of the predicted error span in the original clinical note, limiting the LLMs' tendency to generate generic corrections. Second, the MCQ-style prompt addresses the LLMs' tendency to be verbose by limiting corrections to a specific error span.

### RQ3.2: Can end-to-end LLMs perform better when provided with a span hint?

In our RQ2 experiments with end-to-end systems, we observed limitations in the LLM's ability to accurately locate errors within the clinical notes. While in RQ3.1, we noticed that integrating error span predictions helped improve the LLM's performance. These insights motivated us to integrate the error span predictions from fine-tuned models to the end-to-end LLM solution. We denoted this solution as the "Hybrid approach", as mentioned in Subsubsection 3.2.2, leveraging the "Span hint" from the error span prediction.

Integrating a span hint into the end-to-end LLM prompt resulted in improvements across all metrics, as shown in Table 7. Notably, span hint significantly improved the aggregate NLG scores of Brief CoT and no-CoT solutions. However, span hint did not improve Long CoT solution, suggesting that the reasoning style may influence the LLM's ability to leverage span hints.

Despite MCQ prompts demonstrating higher accuracy in error sentence identification, "Brief CoT" prompts combined with ICL, type hint, and span hints showed a higher aggregate NLG score, emphasising the different strengths of the two strategies. This indicates that the hybrid approach harnesses the LLM's generative capabilities, while the fine-tuned error span prediction model helps direct these corrections to the appropriate error locations.

### Performance on Test Set

We submitted our four best-performing solutions to be evaluated on the holdout test set. As shown in Table 8, we can observe a similar trend as in the validation set experiments. The 2-options MCQ prompts show strong performance in accurately

identifying the error-containing sentence. The 8-shot + Brief CoT + Hints method performs better, especially in the aggregate NLG score. This suggests that while MCQ prompts effectively direct the model's focus, enabling accurate detection of errors, they may slightly constrain the model's generative capability. Overall, these results highlight the benefit of using concise CoT reasoning in LLMs as well as providing guidance via targeted hints. Our best-performing pipeline, 8-shot + Brief CoT + Hints, ranked sixth in the shared task leaderboard based on the aggregate NLG score.

# 5 Post-hoc Analyses

Commonly reported NLG metrics tend to not be well correlated with human judgement, especially in the clinical domain (Ben Abacha et al., 2023). To understand the limitations of LLMs for clinical note correction, we extend beyond the reported performance metrics by analysing the sensitivity of LLMs to the data and prompt, as well as the common mistakes that LLMs tend to commit.[4]

## 5.1 Sensitivity

It is a well-known fact that the performance of an LLM may differ massively given slight differences in the way we prompt it (Voronov et al., 2024). We analysed factors observed in the data and prompt that may contribute to performance differences.

### 5.1.1 Sensitivity to the position of error sentence in the clinical note

We investigated the sensitivity of the model performance to the position of the error sentence within a given clinical note, dividing them into three cases; if the error sentence is in the first sentence ("*beginning*"), the last sentence ("*end*"), or in between the first and the last sentences ("*middle*").

Figure 3 illustrates the relationship between the NLG metrics and the error sentence position, along with the proportion of the error sentence location. We can observe that ROUGE 1, BERTScore, and BLEURT scores do not vary significantly based on the position of the error sentence. This observation is quantitatively supported by the Kruskal-Wallis H-Test and the post-hoc Dunn's test results shown in Appendix D. The test results reveal that the LLM's ability to generate accurate corrections is not impacted by where the error appears in the input, which is a desirable trait.

---

[4]Post-hoc analyses are conducted on the validation sets.

| Role | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| Clinician assistant | 0.5610 | 0.5454 | 0.7385 |
| No role | 0.5570 (-0.00) | 0.5416 (-0.00) | 0.7504 (+0.01) |
| Assistant | 0.5509 (-0.01) | 0.5442 (-0.00) | 0.7504 (+0.01) |
| Medical student | 0.5539 (-0.01) | 0.5468 (+0.00) | 0.7484 (+0.01) |
| Nurse | 0.5763 (+0.02) | 0.5615 (+0.02) | 0.7424 (+0.00) |
| Clinical note verificator | 0.5554 (+0.01) | 0.5438 (-0.00) | 0.7518 (+0.01) |
| Clinician | 0.5793 (+0.02) | 0.5615 (+0.02) | 0.7615 (+0.02) |

Table 9: Performance of our best-performing solution when prompted with different roles via the system prompt (i.e., *"You are «a role» tasked to ..."*) on the validation sets.

| Generated Option Position | $Acc_{flag}$ | $Acc_{sent\_id}$ | $Score_{agg}$ |
|---|---|---|---|
| A | 0.6131 | 0.6029 | **0.6492** |
| B | **0.6368** | **0.6265** | 0.6380 |

Table 10: Results of the sensitivity analysis of MCQ-style prompt to the position of the LLM-generated option in the 2 options setting on validation sets.

### 5.1.2 Sensitivity to the role described in the system prompt

Owing to their instruction-following ability, LLMs are capable of playing a role as prompted by the user (Wang et al., 2023). In the clinical domain, we tend to prompt an LLM to answer a query as a healthcare professional, such as a clinician. In this analysis, we explored how the role prompted or the lack thereof may affect the performance of the LLM in generating corrections. We modify the system prompt (i.e., *"You are «a role» tasked to ..."*) with various role options. Table 9 details the varying performances of the best-performing 8-shot + Brief CoT + hints solution when prompted with different roles. The LLM performs best when prompted to role-play as a "clinician". This phenomenon, known as *In-Context Impersonation* (Salewski et al., 2024), highlights that role-playing should be examined when developing a prompt-based solution.

### 5.1.3 Sensitivity to the position of the multiple choice options

Table 10 shows the outcome of a sensitivity analysis, based on the relative positioning of the LLM-generated option and the predicted error span within the original text for the systems with MCQ-type prompts. Both binary classification accuracy and error sentence prediction accuracy were improved when the LLM-generated option was positioned as option B, as opposed to option A. On the other hand, the aggregate score for correction

Figure 3: Boxplots of the distribution of ROUGE 1, BERTScore, and BLEURT with respect to the position of the error sentence for MS (left) and UW (right) datasets. "beginning" denotes that the error sentence is at index 0, "end" at the end, while "middle" is in between "beginning" and "end".

reveals a higher score when the LLM-generated option was positioned as option A, achieving a score of 0.6492. This observation of *selection bias* echoes findings by previous studies (Pezeshkpour and Hruschka, 2023; Zheng et al., 2023).

## 5.2 Common LLM mistakes

We qualitatively evaluated the common mistakes found in the generated reasons and corrections.

**Corrections of marginal effects**   LLMs occasionally make minor corrections to clinical notes that, although technically correct, do not significantly affect the correctness. Changes, such as altering "3" to "three" or fixing grammatical mistakes, might enhance readability but are not clinically significant. LLMs also tend to add adjectives, such as "acute" to "pyelonephritis", adding specificity desirable in clinical settings but not always favourably reflected in NLG metrics.

**Near-accurate corrections**   LLMs often suggest near-accurate corrections that lack the required specificity. For example, fixing an error sentence with the generic "antiplatelet therapy" instead of "aspirin" misses the required precision, even though aspirin is an antiplatelet therapy. Likewise, proposing to "Start anticoagulation therapy" instead of the more explicit "dalteparin" lacks specificity. These near-accurate adjustments underscore the difficulty LLMs encounter in achieving the specificity of the ground truth label.

**Mistake due to incomplete context**   LLMs struggle to fix errors in clinical notes when details are

lacking. One example is when the LLM mistakenly suggests changing "pulmonary fibrosis" to "chronic obstructive pulmonary disease". Both conditions share very similar early symptoms that are difficult to differentiate even for clinicians (Chilosi et al., 2012). Another example involves incorrectly adjusting a malnutrition patient's Body Mass Index (BMI) from 30 to 18. Albeit a BMI of 18 signals malnutrition, it deviates from the ground truth label 13. These instances underscore the complexity of the MEDIQA-CORR task, as well as medical error correction in general which is very challenging to do without additional context even for human clinicians.

In summary, the sensitivity and qualitative analyses highlight the current limitations of LLMs in the clinical domain, which prompt further questions about the readiness of LLMs to be implemented in real-world clinical settings.

## 6   Conclusion

This study explores strategies for using LLMs to detect and correct medical error for the MEDIQA-CORR 2024 shared task. In addition to the comprehensive evaluation of prompting strategies based on different reasoning styles, we experiment with integrating error-span predictions from a fine-tuned model. Our best-performing system includes a fine-tuned BioLinkBERT-large for error-span prediction and GPT-4 for error correction. By harnessing LLMs' generative abilities with 8-shot ICL and Brief CoT and presenting predicted error span as a hint in the prompt, our best-performing solu-

495

tion ranked sixth in the shared task leaderboard. Our post-hoc analyses offer insights into the use of LLM in medical error correction, including sensitivity to error location, role-playing bias, and common types of mistakes made by LLMs.

## Limitations

The scope of our study was exclusively confined to GPT-based models, namely GPT-3.5 and GPT-4. The reported findings may differ across different types of LLMs. Furthermore, we independently explored various prompting strategies, such as CoT and MCQ prompt. We did not investigate the effect of integrating MCQ prompt with CoT reasoning. This unexplored combination may offer additional improvements in the LLM's error correction capabilities.

Our post-hoc analyses also reveal a significant limitation of LLMs in clinical settings. Despite the advancements demonstrated through our proposed methodologies, the study underscores that LLMs may not be ready for deployment in real-world clinical environments without human oversight. The analysis highlights the critical need for human supervision, especially given the potential risks associated with inaccuracies in medical documentation and the consequent impacts on patient care. This limitation calls for further research into enhancing the reliability of LLMs as well as the evaluation metrics before considering their implementation in sensitive areas such as healthcare.

## Acknowledgements

## References

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Asma Ben Abacha, Wen-wai Yim, George Michalopoulos, and Thomas Lin. 2023. An investigation of evaluation methods in automatic medical note generation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2575–2588, Toronto, Canada. Association for Computational Linguistics.

Marco Chilosi, Venerino Poletti, and Andrea Rossi. 2012. The pathogenesis of copd and ipf: distinct horns of the same devil? *Respiratory research*, 13:1–9.

Matúš Falis, Aryo Pradipta Gema, Hang Dong, Luke Daines, Siddharth Basetti, Michael Holder, Rose S Penfold, Alexandra Birch, and Beatrice Alex. 2024. Can gpt-3.5 generate and code discharge summaries? *arXiv preprint arXiv:2401.13512*.

Aryo Pradipta Gema, Luke Daines, Pasquale Minervini, and Beatrice Alex. 2023. Parameter-efficient fine-tuning of llama for the clinical domain. *arXiv preprint arXiv:2307.03042*.

Aryo Pradipta Gema, Giwon Hong, Pasquale Minervini, Luke Daines, and Beatrice Alex. 2024. Edinburgh clinical nlp at semeval-2024 task 2: Fine-tune your model unless you have access to gpt-4.

Xuanli He, Yuxiang Wu, Oana-Maria Camburu, Pasquale Minervini, and Pontus Stenetorp. 2023. Using natural language explanations to improve robustness of in-context learning for natural language inference. *arXiv preprint arXiv:2311.07556*.

Hyunjae Kim, Hyeon Hwang, Chaeeun Lee, Minju Seo, Wonjin Yoon, and Jaewoo Kang. 2023. Exploring approaches to answer biomedical questions: From pre-processing to gpt-4 notebook for the bioasq lab at clef 2023. In *CEUR Workshop Proceedings*, volume 3497, pages 132–144. CEUR-WS.

---

[5]https://edinburgh-international-data-facility.ed.ac.uk/

Valentin Liévin, Christoffer Egeberg Hother, and Ole Winther. 2022. Can large language models reason about medical questions? *arXiv preprint arXiv:2207.08143*.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions. *arXiv preprint arXiv:2308.11483*.

Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. *Nist Special Publication Sp*, 109:109.

Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2024. In-context impersonation reveals large language models' strengths and biases. *Advances in Neural Information Processing Systems*, 36.

Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. BLEURT: learning robust metrics for text generation. *CoRR*, abs/2004.04696.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Anton Voronov, Lena Wolf, and Max Ryabinin. 2024. Mind your format: Towards consistent evaluation of in-context learning improvements. *arXiv preprint arXiv:2401.06766*.

Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. LinkBERT: Pretraining language models with document links. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8003–8016, Dublin, Ireland. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In *ICLR*. OpenReview.net.

Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2023. Large language models are not robust multiple choice selectors. In *The Twelfth International Conference on Learning Representations*.

| Parameter | Value |
|---|---|
| Model Name | gpt-3.5-turbo-0613 |
| API Version | 2023-03-15-preview |
| Temperature | 0 |
| Top P | 0 |
| Frequency Penalty | 0 |
| Presence Penalty | 0 |
| Max new token | 256 |

Table 11: GPT-3.5 API call hyperparameters to generate Natural Language Explanations.

## A  Experimental setup

All fine-tuning experiments were run on a single NVIDIA A100-40GB GPUs. We used the HuggingFace's transformer library (Wolf et al., 2020). The validation set was utilised to determine the best checkpoint.

In-context examples were retrieved from the Training set. Additionally, the validation set was used to evaluate and select the optimal prompt design. For the test submission, we also retrieved In-context examples from the MS and UW validation sets.

## B  Hyperparameters

### B.1  GPT-3.5 Hyperparameters for the generation of Natural Language Explanation

We prompted GPT-3.5 (model name: `gpt-3.5-turbo-0613`) with hyperparameters as shown in Table 11. The generation process took approximately 2 hours and cost $2.

### B.2  GPT-4 generation hyperparameters

During inference on the test set, we prompted GPT-4 (model name: `gpt-4-turbo`) as shown in Figure 1 Step 3. We set `temperature=0` to ensure that the model's generation is deterministic. The maximum generation length is 512, allowing longer

CoT reasons. One generation process took approximately 2 hours and cost $35.

## C Prompt Examples

Here, we provide examples of the prompts used in our experiments. The black text within the box represents the prompt input text, the <span style="color:red">red</span> text represents the prediction of the models, and the <span style="color:blue">blue</span> text represents the ground truth.

### C.1 Prompt for In-Context Learning and Chain-of-Thought

#### C.1.1 System Prompt

> **System Prompt for the In-Context Learning and Chain-of-Thought Experiments**
>
> You are a clinician assistant tasked with reviewing clinical texts that have been corrupted by an evil third party. Each clinical text may contain either one perturbed sentence with clinical or factual inaccuracies, or no errors at all. Your objective is to help the clinician in identifying any perturbed sentence and provide a correction. Please respond in JSON format with the following structure:
>
> - 'reason': Describe the reasoning behind identifying a specific sentence as incorrect or all clinical text as correct.
> - 'incorrect_sentence_id': If you identify an incorrect sentence, provide its ID here. If all sentences are correct, respond with -1.
> - 'correction': If an incorrect sentence is identified, provide a corrected sentence or an explanation of the inaccuracy. If all sentences are correct, respond with 'NA'.
>
> When evaluating the text, focus specifically on clinical or factual inaccuracies. This could include incorrect medical information, factual errors related to patient care, or erroneous data interpretations. Your detailed assessment and correction are critical for ensuring the reliability and accuracy of our clinical documentation. You have to be more action-oriented, don't suggest planning, discussion, or something that is not action-oriented.

#### C.1.2 CoT reasons

> **Brief Chain-of-Thought**
>
> Streptococcus pneumoniae is not typically associated with skin lesions with a blue livid center, pale intermediate zone, and a dark red peripheral rim. Mycoplasma pneumoniae is a common cause of atypical pneumonia and can present with skin rashes like the ones described in the text.

> **Long Chain-of-Thought**
>
> The incorrect sentence is sentence 24, which states 'Pathology confirms Streptococcus pneumoniae is the causal organism.' This is inaccurate because Streptococcus pneumoniae is a common cause of community-acquired pneumonia, but in this case, the clinical presentation with skin lesions and the absence of a productive cough suggests a more likely diagnosis of atypical pneumonia caused by Mycoplasma pneumoniae. Mycoplasma pneumoniae is known to cause atypical pneumonia with extrapulmonary manifestations such as skin lesions, which aligns better with the patient's presentation. Therefore, the correction should be 'Pathology confirms Mycoplasma pneumoniae is the causal organism.'

> **SOAP Chain-of-Thought**
>
> Subjective: The patient presents with fever, nausea, nonproductive cough, headaches, generalized fatigue, muscle and joint pain, and increasing shortness of breath.
> Objective: Vital signs show a slightly elevated temperature, normal pulse, respirations, and blood pressure. Skin lesions with specific characteristics are noted on the extremities. Laboratory studies reveal hemoglobin, leukocyte count, electrolyte levels, and other values. An x-ray of the chest is performed.
> Assessment: The patient likely has an infectious process given the symptoms and laboratory findings.
> Plan: Further evaluation and treatment are warranted.
> Inconsistency: The incorrect sentence is in the Assessment section where Streptococcus pneumoniae is mentioned as the causal organism, which is inconsistent with the typical causative agent for atypical pneumonia, which is Mycoplasma pneumoniae.

### C.1.3 Chain-of-Thought Prompt

> **Brief Chain-of-Thought Prompt Example**
>
> «System Prompt»
> «In-Context Learning Examples»
> Clinical text:
>
> 0 A 78-year-old man is brought in to the emergency department by ambulance after his wife noticed that he began slurring his speech and had developed facial asymmetry during dinner approximately 30 minutes ago.
> 1 His past medical history is remarkable for hypertension and diabetes.
> 2 His temperature is 99.1 F (37.3 C), blood pressure is 154/99
> 3 mmHg, pulse is 89/min, respirations are 12/min, and oxygen saturation is 98% on room air.
> 4 CTA of the head is obtained after neurologic exam reveals right upper and lower extremity weakness and an asymmetric smile.
>
> Task: Identify 1 incorrect sentence in the clinical text, or confirm if all sentences are correct. If there is 1 incorrect sentence, how would you fix it? Follow the hints below if you see fit
> Hint:
> - Pay special attention to biomedical entities such as chief complaints, medical exams, diagnoses, and treatments. The mistake often lies within these areas.
> - The clinician said that you MAY want to pay attention to the mention of 'CTA of the head'. If you believe that the mention is incorrect, ONLY SWAP this mention with something more probable. DO NOT MODIFY the sentence in any other way.
> Let's think step by step
> Answer:
> {
> **"reason"**: "The mention of 'CTA of the head' is not the most appropriate imaging study for the patient's presentation of acute stroke symptoms with slurred speech, facial asymmetry, and extremity weakness. A non-contrast CT of the head is typically the initial imaging study to evaluate for acute intracranial haemorrhage or ischemic stroke. Therefore, replacing 'CTA of the head' with 'CT of the head' would align better with the standard diagnostic approach for acute stroke evaluation.",
> **"incorrect_sentence_id"**: "4",
> **"correction"**: "CT of the head is obtained after neurologic exam reveals right upper and lower extremity weakness and an asymmetric smile.",
> }

### C.2 Option Generation Prompt Multiple-Choice Question Prompt

### C.2.1 MCQ 2 options

> **Prompt Used to Generate MCQ Answer Options**
>
> Your job is to review a clinical note that potentially contains a medical error.
>
> In the following clinical note, what should the <BLANK> in the sentence "Suspected of <BLANK>." be replaced with if "primary ciliary dyskinesia" is incorrect? Do not answer with "primary ciliary dyskinesia" or its medical synonyms in your answer. Output your response in JSON format, with keys 'option'.
>
> Clinical note:
>
> A 4-year-old boy is brought to the physician in December for episodic shortness of breath and a nonproductive cough for 3 months. These episodes frequently occur before sleeping, and he occasionally wakes up because of difficulty breathing. His mother also reports that he became short of breath while playing with his friends at daycare on several occasions. He is allergic to peanuts. He is at the 55th percentile for height and weight. Vital signs are within normal limits. Examination shows mild scattered wheezing in the thorax. An x-ray of the chest shows no abnormalities. Suspected of <BLANK>.
>
> Generated answer:
> {
> **"option"**: "asthma"
> }

### C.2.2 MCQ 4 options

## C.3 Inference Prompt Multiple-Choice Question Prompt

### C.3.1 MCQ 2 options

### C.3.2 MCQ 4 options

> **Inference Prompt for Multiple-Choice Question style with 4 options**
>
> Your job is to review a clinical note that potentially contains a medical error.
>
> In the following clinical note, what should the <BLANK> in the sentence "Culture tests indicate <BLANK>." be replaced with for it to be medically informative and accurate? Choose one from the options given below. Output your response in JSON format, with a key 'Answer'.
>
> Clinical note:
>
> A 4-year-old boy is brought to the physician in December for episodic shortness of breath and a nonproductive cough for 3 months. These episodes frequently occur before sleeping, and he occasionally wakes up because of difficulty breathing. His mother also reports that he became short of breath while playing with his friends at daycare on several occasions. He is allergic to peanuts. He is at the 55th percentile for height and weight. Vital signs are within normal limits. Examination shows mild scattered wheezing in the thorax. An x-ray of the chest shows no abnormalities. Suspected of <BLANK>.
>
> Options:
>
> A. asthma
> B. primary ciliary dyskinesia
> C. bronchiolitis
> D. pulmonary embolism
>
> Generated answer: {
> **"Answer"**: "A. asthma"
> }

|   | MS | | | UW | | |
|---|---|---|---|---|---|---|
|   | ROUGE 1 | BERTScore | BLEURT | ROUGE 1 | BERTScore | BLEURT |
| H | 6.0749 | 5.0249 | 7.2848 | 5.6821 | 3.6073 | 2.3457 |
| p | 0.0480 | 0.0811 | 0.0262 | 0.0584 | 0.1647 | 0.3095 |

Table 12: Summary of Kruskal-Wallis H-Test results for sentence position impact on ROUGE 1, BERTScore, and BLEURT metrics. Statistically significant differences ($p < 0.05$) are highlighted in cyan.

|   | MS | | | UW | | |
|---|---|---|---|---|---|---|
|   | ROUGE 1 | BERTScore | BLEURT | ROUGE 1 | BERTScore | BLEURT |
| beginning-middle | 0.1751 | 0.3121 | 0.1389 | 0.3596 | 0.3464 | 0.7118 |
| middle-end | 0.2137 | 0.2479 | 0.1192 | 0.5251 | 1.0000 | 1.0000 |
| beginning-end | 0.3923 | 0.6258 | 0.3586 | 0.0609 | 0.2849 | 0.4757 |

Table 13: Summary of Post-hoc Dunn's Test results for sentence position impact on ROUGE 1, BERTScore, and BLEURT metrics. No significant differences observed.

## D  Statistics of "Sensitivity to the position of error sentence in the clinical note"

The analysis was split into two main tests: the Kruskal-Wallis H-Test to identify overall differences across sentence positions and the Post-hoc Dunn's Test to investigate pairwise differences between sentence positions.

The Kruskal-Wallis H-Test was applied to compare the distributions of scores for ROUGE 1, BERTScore, and BLEURT across three sentence positions (beginning, middle, end) within clinical notes from the validation sets of MS and UW. As shown in Table 13, statistically significant differences were found in the MS dataset for ROUGE 1 and BLEURT metrics, suggesting sensitivity to sentence positioning.

Following the Kruskal-Wallis H-Test, a Post-hoc Dunn's Test was performed to conduct pairwise comparisons between sentence positions for each evaluation metric. The Post-hoc Dunn's Test revealed no statistically significant differences between any pairwise comparisons of sentence positions for all evaluated metrics, suggesting that while overall differences exist, specific pairwise comparisons did not reach statistical significance.

# UMass-BioNLP at MEDIQA-M3G 2024: DermPrompt - A Systematic Exploration of Prompt Engineering with GPT-4V for Dermatological Diagnosis

**Parth Vashisht**[*], **Abhilasha Lodha**[*], **Mukta Maddipatla**[*],
**Zonghai Yao, Avijit Mitra, Zhichao Yang, Junda Wang, Sunjae Kwon, Hong Yu**

**CICS, University of Massachusetts, Amherst, MA, USA**
{pvashisht, alodha, mmaddipatla, zonghaiyao}@umass.edu

## Abstract

This paper presents our team's participation in the MEDIQA-ClinicalNLP 2024 shared task B. We present a novel approach to diagnosing clinical dermatology cases by integrating large multimodal models, specifically leveraging the capabilities of GPT-4V under a retriever and a re-ranker framework. Our investigation reveals that GPT-4V, when used as a retrieval agent, can accurately retrieve the correct skin condition 85% of the time using dermatological images and brief patient histories. Additionally, we empirically show that Naive Chain-of-Thought (CoT) works well for retrieval while Medical Guidelines Grounded CoT is required for accurate dermatological diagnosis. Further, we introduce a Multi-Agent Conversation (MAC) framework and show it's superior performance and potential over the best CoT strategy. The experiments suggest that using naive CoT for retrieval and multi-agent conversation for critique-based diagnosis, GPT-4V can lead to an early and accurate diagnosis of dermatological conditions. The implications of this work extend to improving diagnostic workflows, supporting dermatological education, and enhancing patient care by providing a scalable, accessible, and accurate diagnostic tool. [1]

## 1 Introduction

Diagnosing skin conditions demands a complex blend of visual inspection, patient history examination, and deep clinical acumen, a skill set that dermatologists spend extensive years acquiring (Mangion et al., 2023). Despite the critical nature of these skills, many regions worldwide face a stark scarcity of dermatological expertise (Benner et al., 2009). Even in areas with adequate services, the demand for such specialized knowledge frequently surpasses its availability. The recent global

health crisis has also expedited the shift towards remote clinical diagnostics and treatments, further highlighting the challenges in diagnosing skin diseases (Behar et al., 2020). These challenges include the scarcity of dermatological expertise and the need to accommodate asynchronous patient interactions, including e-visits, emails, and messaging platforms, to ensure continuity and quality of care.

In response to these challenges, recent advancements in Artificial Intelligence (AI), particularly through the development of large language models (LLMs), offer promising solutions to significantly support dermatologists by enhancing clinical diagnosis and treatment processes (McDuff et al., 2023; Singhal et al., 2023b; Tu et al., 2024). Moreover, AI facilitates asynchronous patient services, offering a cost-effective and convenient alternative to traditional methods. Previous works have primarily utilized deep learning for tasks such as skin lesion classification (Udriștoiu et al., 2020; Esteva et al., 2017; Brinker et al., 2019), and dermatopathology (Hekler et al., 2019; Jiang et al., 2020) focusing predominantly on dermoscopic images (Cruz-Roa et al., 2013). These efforts, however, have relied on image-only models, indicating a need for broader applications.

Our research aims to extend the capabilities of AI in dermatology by diagnosing skin diseases and devising appropriate treatment plans based on patients' dermatological images, queries, and medical histories. This approach mirrors the diagnostic process of dermatologists, who rely on high-quality images and comprehensive patient histories to make informed decisions. Although previous studies have explored fine-tuning models on multimodal data (e.g., SkinGPT (Zhou et al., 2023), and MedBLIP (Chen et al., 2023a), our task is particularly challenging due to data availability and image quality limitations, reflecting real-world constraints where high-quality data is either scarce or expensive to obtain.

---

[*] Equal Contribution
[1] The code is released at Github

Studies by OpenAI (Nori et al., 2023a) and Microsoft (Nori et al., 2023b) have demonstrated that generalist foundation models, such as GPT-4, can surpass specifically fine-tuned medical models on various medical benchmarks by employing specialized prompting strategies. Building on these insights, our research leverages both the textual and visual capabilities of GPT-4, targeting the specific task of dermatology diagnosis and treatment.

The diagnostic process for skin lesions or conditions requires meticulous evaluation and is informed by methodologies such as dermatoscopy (Panagoulias et al., 2024), which enables dermatologists to observe skin abnormalities in greater detail. Dermatologists usually follow a common guideline for assessing skin lesions, emphasizing the importance of visual descriptors like shape, size, color, texture, and pattern in differential diagnosis. Inspired by these practices, we have integrated advanced Chain-of-Thought (CoT) techniques with visual features to create medical guidelines tailored for GPT-4V, enhancing its diagnostic precision. This enables the model to emulate dermatologists' diagnostic process.

Furthermore, our research integrated a Multi-Agent Conversation (MAC) Framework (Tao et al., 2024a; Wu et al., 2023; Li et al., 2023b), which involves multiple AI agents that generate additional context and critiques for various candidate skin conditions. These agents collaborate, debate, and consolidate their findings to determine the most accurate skin disease diagnosis from the candidates identified from the retrieval step. This, therefore, introduces a level of dynamic interaction and comprehensive analysis that mirrors the complex decision-making process in clinical dermatology.

Hence, our contributions are twofold:

- We deploy GPT-4V within a novel retrieval and re-ranking framework, critically evaluating the effectiveness of various prompting strategies. These include both naïve prompts and those meticulously crafted based on detailed medical guidelines (CoT), across different stages of our setup. This exploration aims to highlight the adaptability and precision of GPT-4V in simulating the diagnostic reasoning of dermatologists.
- We explore the Multi-Agent Conversation (MAC) Framework in the context of clinical dermatology, examining its potential to en-

rich the diagnostic process. Through this discussion, we identify and delineate the framework's strengths and limitations, offering insights into its applicability and performance in accurately diagnosing skin diseases.

## 2 Related Work

The interdisciplinary fusion of Artificial Intelligence (AI) and dermatology has spawned a myriad of approaches to enhance the diagnosis of skin conditions. Historically, these approaches have often treated diagnosis as a classification task, with literature extensively documenting the use of convolutional neural networks (CNNs) and other deep learning architectures like ResNets for lesion classification from dermoscopic images, which are typically limited to dermatological clinics due to image acquisition constraints (Ba et al., 2022), (Gouda and Amudha, 2020).

Recent advancements have moved beyond traditional clinic-bound methods, exploring the utility of clinical images for broader classifications, such as skin cancer and onychomycosis (Sharma et al., 2022). While these efforts have made significant contributions to disease diagnosis, they have not fully addressed the generative and comprehensive nature of clinical diagnosis, which encompasses treatment planning and patient interaction beyond mere classification.

The evolution of large language models (LLMs) has significantly widened the scope of AI applications in healthcare. LLMs like PubMedBERT (Gu et al., 2020) and BioGPT (Luo et al., 2022) have been fine-tuned on extensive corpora of medical literature, achieving state-of-the-art performance in tasks ranging from biomedical reasoning to question-answering. In the realm of domain-adapted LLMs, models like Meditron and Med-PALM have demonstrated remarkable capabilities in language understanding and generation, setting new benchmarks across biomedical datasets (Chen et al., 2023b; Singhal et al., 2022).

With the advent of multimodal models, integrating visual and textual data has further refined AI applications in medical domains. Vision-language models such as Med-CLIP, Med-BLIP, and Llava-Med have exhibited promising results in image-text retrieval, zero-shot classification, and even multimodal conversations, respectively (Chen et al., 2023a; Li et al., 2023a). Specifically, in the context of dermatology, the Skin-GPT4 model (Zhou

Figure 1: Overview of the AI-assisted dermatology diagnosis pipeline, from initial patient input through to the GPT-4V generated final diagnosis and treatment plan.

et al., 2023) represents a pioneering effort in creating a multimodal setup tailored for skin disease identification and patient interaction.

Studies have shown that generalist foundation models like GPT-4, with their expansive knowledge bases and specialized prompting techniques, outperform domain-specific models such as Med-PALM on various medical benchmarks (Nori et al., 2023a). GPT-4's application in dermatology, particularly in melanoma identification and medical exam question answering, underscores its potential as an assistive tool for educational and diagnostic purposes (Miao et al., 2024; Mishra et al., 2024; Yang et al., 2023).

Our research builds upon these foundations, employing GPT-4's multimodal capabilities (GPT-4V) to enhance dermatological diagnostic processes. By integrating Chain-of-Thought (CoT) techniques and a Multi-Agent Conversation (MAC) Framework (Tao et al., 2024a; Wu et al., 2023; Li et al., 2023b), we aim to emulate the complex decision-making process of dermatologists, enriching the GPT-4V's ability to generate diagnostic and treatment plans from multimodal data. This work not only taps into the multimodal analytical strength of GPT-4V but also seeks to optimize the model's performance in a domain where the nuances of patient history and visual inspection are paramount.

Thus, our contribution to the field involves the innovative use of GPT-4V within a retrieval and re-ranking framework, leveraging both naïve and medically informed CoT prompting strategies.

## 3 Methodology

Our methodology delineates the comprehensive approach we adopted to address the task of mul-

timodal medical answer generation. This process involves two primary stages: retrieving potential diagnoses and the ranking of these to identify the most probable skin condition and treatment plan.

Task Description: The objective of our research is to develop a system capable of diagnosing a possible skin condition and recommending a corresponding treatment plan based on a patient's medical query and associated image. To accomplish this, we propose a two-step pipeline consisting of a retrieval module followed by a ranker module. Specifically:

- Retrieval Module: This component extracts a list of possible skin conditions from the given image and medical query.

- Ranker Module: This module's task is to select the most accurate skin condition diagnosis from the list generated by the retrieval module.

Our overall methodology is mentioned in Figure 1.

### 3.1 Retrieval Module

The Retrieval Module is the initial phase of our diagnostic approach, reflecting the dual aspects a dermatologist considers when evaluating a medical condition: visual inspection and patient history. Inspired by recent works such as MedGE-NIE (Frisoni et al., 2024), along with others (Yu et al., 2022; Zhang et al., 2023; Su et al., 2022), we leverage LLMs as strong context generators instead of traditional retrieval methods, such as keyword-based methods (e.g., BM25 Robertson et al. (2009)), vector-similarity-based methods (e.g., ColBERT Khattab and Zaharia (2020)), and

some internet tools (e.g., Google API). In this retrieval step, we treat LLMs as a knowledge base (Singhal et al., 2023a) to generate potentially valuable information for subsequent steps. Our module employs two distinct strategies:

### 3.1.1 Context-Independent Retrieval (Image-Only)

Recognizing scenarios where comprehensive medical context (patient's medical history and medical queries) may not be readily available, we engage in context-independent retrieval. This approach leverages GPT-4V to identify possible skin conditions based solely on image data. We compared this model's performance against a widely-used online AI tool, First Derm [2]

### 3.1.2 Context-Dependent Retrieval (Image + Context)

The inclusion of medical context is pivotal for accurate diagnosis. Particularly, incorporating details about systemic conditions and patient history can significantly influence differential diagnosis, a critical aspect of clinical dermatology. To this end, we utilize CoT prompting, a technique that simplifies complex problems into manageable objectives, enabling the model to address the larger task.

Within the Context-Dependent Retrieval, we experiment with two strategies:

- Naive CoT: Here, GPT-4V is instructed to methodically analyze all relevant information from the images and medical query before generating a list of potential skin conditions. This process mimics the step-by-step procedural thinking a dermatologist might employ.
- Expert Guidelines Grounded CoT: This approach involves crafting prompts based on the Clinical Guidelines that dermatologists follow, encapsulating a generic framework for skin disease diagnosis. Such frameworks typically comprise patient history, visual inspection, and differential diagnosis. Our Expert-CoT strategy emphasizes key visual characteristics like the lesion's shape, color, size, location, and texture. By integrating this data, the module produces a detailed list of differential diagnoses needed to enhance the model's diagnostic precision further. [3]

## 3.2 Re-Ranker Module

After the retrieval module identifies potential skin conditions, the re-ranker module is critical in our diagnostic pipeline. Its primary objective is to meticulously refine the preliminary list, pinpointing the diagnoses with the highest probability of accuracy. To achieve this, we experiment with four re-ranking strategies:

1. Naive Chain of Thought (CoT)
2. Expert Guidelines Grounded CoT with Context
3. Expert Guidelines Grounded CoT without Context
4. Multi-Agent Conversation Framework

These structured approaches enable a systematic evaluation of the candidate's conditions, ensuring that the decision-making process mirrors the analytical and deductive reasoning of a dermatologist. The specific prompts utilized for these three CoT techniques are presented in Tables 14, 15, and 16. The MAC framework is explained with an example in the Appendix A.

### 3.2.1 Naive Chain of Thought (CoT)

In the Naive CoT approach, GPT-4V is initially instructed to analyze the patient's medical query and the associated images to extract relevant information. Subsequently, each candidate skin condition retrieved from section 3.1 is assigned a score ranging from 1 to 10, where 1 signifies the least probable and 10 denotes the most probable condition. The model identifies the most probable disease based on the scores and analysis. The prompt is mentioned in Table 14.

### 3.2.2 Expert Guidelines Grounded CoT with Context

This method employs a sophisticated strategy by utilizing prompts meticulously designed around the Clinical Guidelines followed by dermatologists. These guidelines encapsulate a comprehensive visual assessment of the affected area, scrutinizing distinct characteristics such as shape, size, color, location, and texture (listed in Table 15). Within this framework, GPT-4V is initially directed to conduct an analysis of the patient's condition, incorporating insights drawn from their medical query. Subsequently, the patient's images undergo visual examination using the defined guidelines, using which

---

[2]https://firstderm.com/
[3]The specific prompts utilized for both the Naive CoT and

the Expert Guidelines Grounded CoT strategies are detailed in Table 13

relevant features are extracted. The final step involves considering a list of possible skin conditions (retrieved from section 3.1) and systematically ruling them out based on the gathered insights and visual inspections to identify the most probable skin condition from the set of candidates.

### 3.2.3 Expert Guidelines Grounded CoT without Context

This approach omits the user query, focusing exclusively on the visual examination of dermatological conditions as per established guidelines. Utilizing GPT-4V, an initial step involves the generation of a detailed visual description, drawing upon ten specified visual features essential for dermatological assessment (as outlined in Table 16). Subsequently, each candidate's skin condition is described visually, emphasizing distinguishing features aligned with the visual guidelines. A comparative analysis is then conducted between the visual descriptions of the candidates and the initial image description, and a score ranging from 1 to 10 is assigned based on the level of match (1 being the lowest match and 10 the highest). The most probable candidate, determined by the highest score in the comparative analysis, is selected as the diagnosis.

### 3.2.4 Multi-Agent Coversation Setup

Inspired by the recent various applications of the Multi-Agent Conversation framework in the general and medical domains (Wu et al., 2023; Tao et al., 2024b), we also implement a multi-agent conversation framework for our re-ranker module (Figure 2). This framework involves multiple AI agents, each specializing in a different aspect of dermatology diagnosis. These agents collaborate, debate, and consolidate their findings to identify the most accurate diagnosis, mirroring the collaborative approach often seen in medical panels. This multi-agent setup not only enriches the model's diagnostic capabilities but also introduces a level of dynamic interaction and comprehensive analysis that mirrors the complex decision-making process in clinical dermatology. Moreover, acknowledging performance gain and consistency improvement obtained using critique-based refinement in large language models, we incorporate feedback generation as an objective of multi-agent debate followed by refinement.

The main components of our multi-agent setup are defined in Table 1.

**Process Flow**

- Assignment and Analysis
  - The Coordinator assigns a distinct probable disease to each Diagnostic Specialist based on the case study and list of probable diseases.
  - Each specialist analyzes the case study, provides evidence supporting their assigned disease and critiques the applicability of other diseases.
- Compilation and Presentation of Findings
  - After receiving inputs from all specialists, the Coordinator compiles the evidence and critiques.
  - The compiled information is presented to the Admin for evaluation.
- Evaluation and Revision
  - The Admin reviews the evidence and critiques, identifying areas where additional clarity or strengthening is needed.
  - If necessary, the Admin requests revisions from specific specialists to enhance their evidence based on critiques.
- Final Diagnosis
  - With the revised evidence, the Admin conducts a final review to determine the most accurate diagnosis.
  - The process concludes once the Admin confirms the diagnosis.

### 3.3 Aligner

The Aligner Module represents the final step in our diagnostic process, focusing on optimizing the model's output to ensure it aligns with clinical standards and expectations. This involves adjusting the prompt to refine the model's language and structure, aiming to emulate the concise, informative style characteristic of professional medical advice. The optimization process is guided by analyzing real doctor responses in the dataset, identifying key elements such as terminology, format, and the inclusion of essential diagnostic and treatment information. The goal is to produce a diagnosis and treatment plan that not only accurately identifies the patient's condition but also provides actionable, understandable advice. This module highlights our commitment to bridging the gap between AI-generated content and the practical needs of clinical practice, ensuring that the output is not

Figure 2: Multi-Agent Conversation (MAC) Setup

*Note: $E_1$ corresponds to the evidence supporting probable_disease_1 generated by Diagnostic_Specialist_1. $C_2$, $C_3$, and $C_4$ are the critiques for probable_disease_2, 3 & 4 generated by Diagnostic_specialist_1. Based on the Re-definement Instructions from the Admin, the diagnostic specialist returns Redefined evidence($R\_E\_1$ )*

| Roles | Tasks |
|---|---|
| *Coordinator* | Orchestrate sequence of consultations. Assign diseases to specialists and manage communications. |
| *Admin* | Evaluate evidence and critiques for accuracy. Request evidence enhancements and finalize the diagnosis. |
| *Agent* | Analyze case study, advocate for one disease. Provide evidence and critique alternative diagnoses. |

Table 1: Tasks for each role in a multi-agent setup

only technically accurate but also clinically relevant and usable in real-world medical contexts.

Recent work has introduced aligners to assist LLMs in generating harmless outputs (Ji et al., 2024), a concept previously applied in the style transfer domain to map model outputs to desired forms (e.g., formality style transfer (Rao and Tetreault, 2018; Yao and Yu, 2021)). When using third-party APIs like GPT, where updating the model's weights is not an option, recent methods have explored the use of Automatic Prompt Optimization (APO) to improve prompts, assuming access to training data and an LLM API (Prasad et al., 2022; Pryzant et al., 2023). Recent studies

have also applied APO in the clinical domain to assist doctors in generating better note-generation prompts (Yao et al., 2023). Inspired by these efforts, we use human responses from training data as APO's training input, allowing the LLM to derive appropriate aligner prompts to facilitate the final step of style transfer. The final prompt generated by APO can be found in Table 17.

### 3.4 Evaluation

The evaluation of our pipeline is dependent on the accuracy metric. Accuracy is defined individually

for each component.

$$Accuracy = \frac{\text{Number of retrieved GT}}{\text{Total number of data points}} \quad (1)$$

*Number of retrieved GT*: Total number of examples for which ground truth skin condition was present in the retrieved list of candidate skin conditions.
*Total number of data points*: Total number of examples for which ground truth skin condition is known. We skip all those examples for which ground truth is not known. In the validation we have 47 examples for which ground truth is known and a total of 56 examples.

Acknowledging the fact that a same skin condition can have multiple names, we implemented GPT-Eval as an evaluator to identify if two skin conditions are similar or not. Our evaluation strategy employs a rule-based approach to assess the similarity between two skin conditions, "A" and "B", according to four predefined rules (refer Table 18). These rules incorporate name identity, synonymity, common root condition, and shared effects and causes to determine similarity systematically. This method addresses the complex nature of dermatological conditions by providing a structured framework that considers linguistic, clinical, and etiological aspects of skin diseases.

## 4 Results

### 4.1 Retrieval Module

| Retrieval Strategies | Methods | Accuracy |
|---|---|---|
| *Context Independent* | First Derm | 0.468085 |
| | GPT-based | 0.595744 |
| *Context Dependent* | Naive CoT | 0.851063 |
| | Expert CoT | 0.744680 |

Table 2: Comparison of Retrieval Strategies and their Accuracy

The accuracy scores, as reported in Table 2, reveal significant insights into the efficacy of each strategy employed within our Retrieval Module.
Firstly, we observed that Context-Independent Retrieval, which relies exclusively on image data, resulted in lower accuracy when compared to Context-Dependent strategies. This indicates that the absence of medical context limits GPT-4V's ability to identify potential skin diseases accurately. Conversely, Context-Dependent Retrieval exhibited superior results. By incorporating medical

queries along with images, this method provides a richer context to GPT-4V, leading to more precise retrieval of potential skin conditions. It appears that the additional contextual data plays a pivotal role in enhancing the model's diagnostic capabilities. When comparing the two strategies within the Context-Dependent Retrieval, Naive CoT outperformed Expert CoT. This may initially seem counterintuitive, given that Expert CoT is grounded in medical guidelines, which one would expect to yield better results. However, our analysis suggests that the Naive CoT strategy's ability to generate a broader range of potential candidates contributed to its higher accuracy. In contrast, the Expert CoT strategy, which employs differential diagnosis principles, likely eliminated some candidates during the retrieval phase, potentially leading to decreased accuracy.

From these observations, we hypothesize that differential diagnosis, while not as effective in the initial retrieval phase, may be better suited to the re-ranking phase of our diagnostic pipeline. The re-ranking phase requires a systematic evaluation to differentiate between closely related skin conditions, aligning with the differential diagnosis's intrinsic nature. Therefore, the nuanced approach of systematically eliminating similar conditions could prove beneficial in the subsequent stage, where precision is paramount.

### 4.2 Re-Ranker Module

| Methods | Top-2 Accuracy | Top-1 Accuracy |
|---|---|---|
| Naive CoT | 0.553191 | 0.425531 |
| Medical Guidelines (image+context) | 0.617021 | 0.531915 |
| Medical Guidelines (image only) | 0.553191 | 0.446808 |

Table 3: Comparison of Re-Ranker Strategies and their Accuracy

The re-ranker module is important in refining the initial list of potential diagnoses obtained from the retrieval module. The metrics used to evaluate the performance of our re-ranker module are:

- **Top-2 Accuracy:** This metric reflects the model's ability to include the correct diagnosis within its top two predictions from the candidate conditions identified in the retrieval phase.

- **Top-1 Accuracy:** This is the precision with which the model identifies the correct diagnosis as its first and final choice from all possible conditions.

As illustrated in Table 3, the evaluation of our re-ranker strategies reveals several insights.
The Naive CoT and Medical Guidelines (image only) strategies exhibit comparable performance, with both Top-2 and Top-1 accuracies closely aligned. This suggests that even without the medical context, the model can leverage visual cues to a degree of effectiveness.
A notable increase in accuracy is observed with the use of Medical Guidelines alongside context (image + patient's query). Incorporating the patient's medical history and associated query, in conjunction with differential diagnosis techniques as outlined in Table 15, enhances the model's discriminatory power. This aligns with our hypothesis that the systematic approach of differential diagnosis—filtering through similar skin conditions—proves more efficacious in the re-ranking phase.

### 4.3 MAC

| Methods | Accuracy |
|---------|----------|
| MG-GR   | 0.53333  |
| MAC     | 0.73333  |

Table 4: Re-Ranking - top 1 Accuracy using MAC. Here MG-GR is the (Medical Guidelines Grounded Re-Ranker.

The multi-agent conversation setup significantly outperforms the traditional top-1 re-ranking strategy, exhibiting a substantial improvement of nearly 20 percentage points. This enhancement was observed across 15 distinct examples where the number of potential solutions retrieved varied between three and five. We propose that the key mechanism driving this enhanced accuracy is the system's critique-based conversational framework. Within this framework, each participating agent is subject to a rigorous process of critique and feedback from other agents. This collaborative interaction encourages continuous reassessment and refinement of each agent's initial diagnoses and the evidence they present. Consequently, this iterative process likely contributes to more precise and reliable diagnostic outcomes, as each agent integrates insights gained

from the critiques to adjust and improve their reasoning and conclusions.

### 4.4 Aligner Module

|            | DeltaBleu |
|------------|-----------|
| Before APO | 0.944723  |
| After APO  | 2.737657  |

Table 5: DeltaBleu scores before and after Automatic Prompt Optimization (APO)

An important evaluation metric for the competition is the deltableu score. The "DeltaBLEU" score is a variation of the BLEU (Bilingual Evaluation Understudy) score, which is a widely used metric for evaluating the quality of text that has been machine-translated from one language to another. The BLEU score measures the correspondence between a machine's output and that of a human, providing a quantitative assessment of translation accuracy. We leverage Automatic Prompt Optimization (APO) to learn a set of rules that bootstraps our prediction and align the responses. Table 5 shows the bleu score improvement by leveraging the rules learned by APO. The learned rules are mentioned in the Table 17.

- Before Alignment: "Based on the visual descriptions, it seems like the most probable condition is Chronic Eczema. I recommend applying topical steroids and moisturizers regularly for treatment."
- After Alignment: "Consider Eczema, which should manifest similarly on both sides. Treat it with regular use of moisturizers and topical steroids."
- Ground Truth: "Should be happening on both sides. Think of Eczema."

## 5 Discussion and Conclusion

In our study, we systematically explored the merits of various prompting strategies within an information retrieval-based dermatology diagnostic framework. By evaluating these strategies through the lens of accuracy metrics, we found that a naive Chain of Thought (CoT) strategy effectively simulates a retrieval module typical of information retrieval systems. This approach is adept at returning a sufficient number of candidate diagnoses, setting a foundational stage for further analysis. Our findings underscore the importance of including patient

history and contextual information in clinical dermatology to enhance diagnostic accuracy.

For the nuanced task of re-ranking diagnostic candidates, our research indicates that a more refined CoT strategy is necessary. Specifically, prompts that incorporate expert guidelines prove critical in conducting differential diagnoses, yielding superior performance in top-1 and top-2 diagnostic outcomes.

Furthermore, we introduce the novel Multi-Agent Critique (MAC) framework, which incorporates agent-based critique and feedback, and has the potential to perform differential diagnosis and refine it's output using feedback.

## 6 Limitations and Future Work

Our current pipeline does not fully comply with stringent data protection regulations, such as the Health Insurance Portability and Accountability Act (HIPAA). Despite Azure's availability of a HIPAA-compliant hosting option, our framework has not been fully aligned with these regulatory standards. The imperative to protect patient data privacy and ensure security is paramount in clinical applications. Our findings suggest that deploying a local model might offer a more privacy-centric approach. However, achieving satisfactory performance with local deployment necessitates further research and development. This limitation underscores the critical need to balance privacy considerations with technological efficacy, especially in the sensitive context of healthcare.

An additional dimension of our study pertains to the inherent variability in the performance of the prompting strategies, attributed to the high temperature setting utilized during GPT-4's open-ended generation tasks. This element of randomness introduces inconsistencies in the model's responses. We hypothesize that employing over-sampling techniques from GPT-4, coupled with self-consistency prompting, could mitigate these inconsistencies and enhance the overall effectiveness of the diagnostic process.

Furthermore, the MAC framework's practical application presents challenges, notably in the seamless integration of inter-agent communication. Our observations point to instances where the system failed due to unexpected behaviors during these interactions, highlighting the complexities of implementing such a framework effectively. Additionally, the MAC study has been conducted on a small set of 15 samples for which the number of retrieved candidates are in the range of 3 to 5 with the MAC system failing for number of candidates greater than 5. This was because of limited context length window with GPT-4 model. Additionally, each call to the agent is financially prohibitive as the number of candidates increase since more rounds of conversations are needed. Such challenges underscore the need for further research and development to refine and optimize the MAC framework for clinical diagnostic applications.

Additionally, given the challenging nature of the dataset with unclean/missing context for a lot of examples, correct bench marking cannot be assumed, but this study can serve as a potential lower bound of GPT-4V's performance on the complex task of clinical dermatology.

## References

Wei Ba, Huan Wu, Wei W. Chen, Shu H. Wang, Zi Y. Zhang, Xuan J. Wei, Wen J. Wang, Lei Yang, Dong M. Zhou, Yi X. Zhuang, Qin Zhong, Zhi G. Song, and Cheng X. Li. 2022. Convolutional neural network assistance significantly improves dermatologists' diagnosis of cutaneous tumours using clinical images. *European Journal of Cancer*, 169:156–165.

Joachim A Behar, Chengyu Liu, Kevin Kotzen, Kenta Tsutsui, Valentina DA Corino, Janmajay Singh, Marco AF Pimentel, Philip Warrick, Sebastian Zaunseder, Fernando Andreotti, et al. 2020. Remote health diagnosis and monitoring in the time of covid-19. *Physiological measurement*, 41(10):10TR01.

Patricia E Benner, Christine A Tanner, and Catherine A Chesla. 2009. *Expertise in nursing practice: Caring, clinical judgment, and ethics*. Springer Publishing Company.

Titus Josef Brinker, Achim Hekler, and et al. 2019. A convolutional neural network trained with dermoscopic images performed on par with 145 dermatologists in a clinical melanoma image classification task. *European journal of cancer*, 111:148–154.

Qiuhui Chen, Xinyue Hu, Zirui Wang, and Yi Hong. 2023a. Medblip: Bootstrapping language-image pretraining from 3d medical images and texts.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023b. Meditron-70b: Scaling medical pretraining for large language models.

Angel Alfonso Cruz-Roa, John Edison Arevalo Ovalle, Anant Madabhushi, and Fabio Augusto González Osorio. 2013. A deep learning architecture for image representation, visual interpretability and automated basal-cell carcinoma cancer detection. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2013*, pages 403–410, Berlin, Heidelberg. Springer Berlin Heidelberg.

Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin M. Ko, Susan M. Swetter, Helen M. Blau, and Sebastian Thrun. 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542:115–118.

Giacomo Frisoni, Alessio Cocchieri, Alex Presepi, Gianluca Moro, and Zaiqiao Meng. 2024. To generate or to retrieve? on the effectiveness of artificial contexts for medical open-domain question answering. *arXiv preprint arXiv:2403.01924*.

Niharika Gouda and J Amudha. 2020. Skin cancer classification using resnet. In *2020 IEEE 5th International Conference on Computing Communication and Automation (ICCCA)*, pages 536–541.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.

Achim Hekler, Jochen Sven Utikal, Alexander H. Enk, Carola Berking, Joachim Klode, Dirk Schadendorf, Philipp Jansen, Cindy Franklin, Tim Holland-Letz, Dieter Krahl, Christof von Kalle, Stefan Fröhling, and Titus Josef Brinker. 2019. Pathologist-level classification of histopathological melanoma images with deep neural networks. *European Journal of Cancer*, 115:79–83.

Jiaming Ji, Boyuan Chen, Hantao Lou, Donghai Hong, Borong Zhang, Xuehai Pan, Juntao Dai, and Yaodong Yang. 2024. Aligner: Achieving efficient alignment through weak-to-strong correction. *arXiv preprint arXiv:2402.02416*.

Y.Q. Jiang, J.H. Xiong, H.Y. Li, X.H. Yang, W.T. Yu, M. Gao, X. Zhao, Y.P. Ma, W. Zhang, Y.F. Guan, H. Gu, and J.F. Sun. 2020. Recognizing basal cell carcinoma on smartphone-captured digital histopathology images with a deep neural network. *British Journal of Dermatology*, 182(3):754–762.

Omar Khattab and Matei Zaharia. 2020. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 39–48.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023a. Llava-med: Training a large language-and-vision assistant for biomedicine in one day.

Jian Li, Xi Chen, Weizhi Liu, Li Wang, Yingman Guo, Mingke You, Gang Chen, and Kang Li. 2023b. One is not enough: Multi-agent conversation framework enhances rare disease diagnostic capabilities of large language models.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics*, 23(6). Bbac409.

Sean E Mangion, Tai A Phan, Samuel Zagarella, David Cook, Kirtan Ganda, and Howard I Maibach. 2023. Medical school dermatology education: a scoping review. *Clinical and Experimental Dermatology*, 48(6):648–659.

Daniel McDuff, Mike Schaekermann, Tao Tu, Anil Palepu, Amy Wang, Jake Garrison, Karan Singhal, Yash Sharma, Shekoofeh Azizi, Kavita Kulkarni, et al. 2023. Towards accurate differential diagnosis with large language models. *arXiv preprint arXiv:2312.00164*.

Jing Miao, Charat Thongprayoon, Wisit Cheungpasitporn, and Lynn D Cornell. 2024. Performance of GPT-4 Vision on kidney pathology exam questions. *American Journal of Clinical Pathology*, page aqae030.

Prakamya Mishra, Zonghai Yao, Parth Vashisht, Feiyun Ouyang, Beining Wang, Vidhi Dhaval Mody, and Hong Yu. 2024. Synfac-edit: Synthetic imitation edit feedback for factual alignment in clinical summarization.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023a. Capabilities of gpt-4 on medical challenge problems.

Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoifung Poon, Tao Qin, Naoto Usuyama, Chris White, and Eric Horvitz. 2023b. Can generalist foundation models outcompete special-purpose tuning? case study in medicine.

Dimitrios P. Panagoulias, Evridiki Tsoureli-Nikita, Maria Virvou, and George A. Tsihrintzis. 2024. Dermacen analytica: A novel methodology integrating multi-modal large language models with machine learning in tele-dermatology.

Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. 2022. Grips: Gradient-free, edit-based instruction search for prompting large language models. *arXiv preprint arXiv:2203.07281*.

Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with" gradient descent" and beam search. *arXiv preprint arXiv:2305.03495*.

Sudha Rao and Joel Tetreault. 2018. Dear sir or madam, may i introduce the gyafc dataset: Corpus, benchmarks and metrics for formality style transfer. *arXiv preprint arXiv:1803.06535*.

Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

Akhilesh Kumar Sharma, Shamik Tiwari, Gaurav Aggarwal, Nitika Goenka, Anil Kumar, Prasun Chakrabarti, Tulika Chakrabarti, Radomir Gono, Zbigniew Leonowicz, and Michał Jasiński. 2022. Dermatologist-level classification of skin cancer using cascaded ensembling of convolutional neural network and handcrafted features based deep neural network. *IEEE Access*, 10:17920–17932.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2022. Large language models encode clinical knowledge.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023a. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023b. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*.

Dan Su, Mostofa Patwary, Shrimai Prabhumoye, Peng Xu, Ryan Prenger, Mohammad Shoeybi, Pascale Fung, Anima Anandkumar, and Bryan Catanzaro. 2022. Context generation improves open domain question answering. *arXiv preprint arXiv:2210.06349*.

Mingxu Tao, Dongyan Zhao, and Yansong Feng. 2024a. Chain-of-discussion: A multi-model framework for complex evidence-based question answering.

Mingxu Tao, Dongyan Zhao, and Yansong Feng. 2024b. Chain-of-discussion: A multi-model framework for complex evidence-based question answering. *arXiv preprint arXiv:2402.16313*.

Tao Tu, Anil Palepu, Mike Schaekermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, et al. 2024. Towards conversational diagnostic ai. *arXiv preprint arXiv:2401.05654*.

Anca Loredana Udriștoiu, Ariana Elena Stanca, Alice Elena Ghenea, Corina Maria Vasile, Mihaela Popescu, Stefan Udristoiu, Andreea Valentina Iacob, Ștefan Cristian Castravete, Lucian Gheorghe Gruionu, and Gabriel Gruionu. 2020. Skin diseases classification using deep leaning methods. *Current Health Sciences Journal*, 46:136 – 140.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.

Zhichao Yang, Zonghai Yao, Mahbuba Tasmin, Parth Vashisht, Won Seok Jang, Feiyun Ouyang, Beining Wang, Dan Berlowitz, and Hong Yu. 2023. Performance of multimodal gpt-4v on usmle with image: Potential for imaging diagnostic support with explanations. *medRxiv*.

Zonghai Yao, Ahmed Jaafar, Beining Wang, Yue Zhu, Zhichao Yang, and Hong Yu. 2023. Do physicians know how to prompt? the need for automatic prompt optimization help in clinical note generation. *arXiv preprint arXiv:2311.09684*.

Zonghai Yao and Hong Yu. 2021. Improving formality style transfer with context-aware rule injection. *arXiv preprint arXiv:2106.00210*.

Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2022. Generate rather than retrieve: Large language models are strong context generators. *arXiv preprint arXiv:2209.10063*.

Yunxiang Zhang, Muhammad Khalifa, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, and Lu Wang. 2023. Merging generated and retrieved knowledge for open-domain qa. *arXiv preprint arXiv:2310.14393*.

Juexiao Zhou, Xiaonan He, Liyuan Sun, Jiannan Xu, Xiuying Chen, Yuetan Chu, Longxi Zhou, Xingyu Liao, Bin Zhang, and Xin Gao. 2023. Skingpt-4: An interactive dermatology diagnostic system with visual large language model.

## A  Appendix

### A.1  Multi-Agent Conversation (MAC) Case Study

In this section, we present an exemplar case study of a debate facilitated by the Multi Agent Chat (MAC) system. The dialogues exemplified in Table 7 to Table 12 illustrate the dynamic interaction between diagnostic specialists and admin agents within our MAC framework.

The ensuing discussion is prefaced with prompts which have successfully generated the anticipated

outcomes, showcasing the MAC system's adeptness. We provide a series of prompts, as detailed in Table 6 and the expert guidelines to ensure that the communication trajectory remains aligned with the system's strategic objectives.

Through the presented case study, we aim to elucidate the capabilities of the MAC system's in the context of clinical diagnosis and the efficacy of its prompts in steering the group chat among various agents to achieve coherent, goal-oriented dialogue.

**Custom Agent (Diagnostic Specialist) Prompt** is to guide the model in adopting the role of a diagnostic specialist. It aims to facilitate the identification of salient features of a skin condition based on clinical observations that align with the designated skin disease. Utilizing these features and the provided details regarding the disease progression, the model is tasked with generating evidence to substantiate the diagnosis of the specified skin disease.

**Coordinator Prompt** is utilized to guide the model to play the role of a coordinator, orchestrating the conversation among the agents, collecting the generated evidences and critiques, consolidating them and passing them to the Admin agent for further analysis and diagnosis. This agent is essential to ensure smooth transition in between the agents.

**Admin Prompt** is designed to instruct the model to play the role of Admin, the admin is a head doctor, who first meticulously assesses the quality of generated evidences and critiques. Then the Admin is tasked with guiding the relevant specialists to refine their evidence in light of the critiques received, thereby enhancing the robustness of the diagnostic case. This iterative process of evaluation and refinement underscores the Admin's pivotal role in ensuring the accuracy and integrity of the final diagnosis, highlighting the significance of expert oversight in collaborative medical diagnostics.

## A.2 Guidelines and Instructions for Multi-Agent Chat

The instructions and expert guidelines are most crucial for the MAC module, since it is a comprehensive briefing of the objectives that the system aims to fulfill. This prompt delineates not only the sequence of actions requisite for task execution but also the intricacies of inter-agent transitions,

| **Diagnostic Specialist Prompt** |
| --- |
| As 'Rick', you are a medical practitioner specializing in dermatology. You are provided with an image description and assigned a specific skin condition, your role is to tune the image description to match with the disease. <br> Once done you can use the image description to generate a detailed report providing evidence that supports this diagnosis. Afterwards critique each of the other probable diseases by explaining why they do not fit the case study as well as your assigned diagnosis does. Ensure clarity and comprehensiveness in your analysis and critiques. |
| **Coordinator Prompt** |
| As the 'Coordinator', your primary responsibility is to oversee the diagnostic process. You will receive a clinical observation of a skin disease, a case study along with a list of probable diseases. Your task is to assign each Diagnostic Specialist a unique probable disease to advocate for, based on the provided details. You'll ensure that each specialist receives all necessary information to perform their analysis effectively. Finally, gather the tuned image descriptions, evidences and critiques from the specialists and present them to the Admin for final evaluation. Your role is crucial for maintaining efficient communication and organization among the specialists. |
| **Admin Prompt** |
| As the 'Admin', your objective is to evaluate the evidence and critiques provided by the Diagnostic Specialists majorly based and aligned to the image description since solely depending on the case study can be tricky to determine the most probable disease for a given case study. Initiate your process by assessing the quality of each critique. Seek consensus among the critiques to strengthen the evidence for a particular diagnosis. You may need to instruct Diagnostic Specialists to refine their evidence based on your findings. Through a structured discussion with the Coordinator and the Diagnostic Specialists, lead the team towards agreeing on a final, most suitable diagnosis for the case study. |

Table 6: Prompts for Multi Agent Chat.

thereby charting the entirety of procedural flow. Furthermore, it encompasses a set of critical guidelines mandating adherence to principles of clarity and precision, alongside the seamless exchange of information among pertinent agents. These directives are imperative to avoid miscommunication and ensure that all interactions remain aligned with the task's end goals. The Task Prompt is furnished to the GroupChatManager, serving as the catalyst for activating dialogues among specialized agents within the framework. This structured approach is pivotal in harmonizing the collective efforts of diverse agents, thus optimizing the overall functionality and efficacy of the MAC system.

## A.3 Multi-Agent Chat - Example

The entire chat is accessible in our GitHub repository for reference. In this section, we present selected excerpts from the Multi-Agent Chat to illustrate the flow and demonstrate the system's

capabilities:

**Set-up of the Chat:**
**Patient Query (with the context of disease progression**): "The skin condition, as shown in the images, presents widespread erythematous patches with violaceous hues across the leg. The patient has multiple crusted plaques and erosions, with sizes varying from a few millimeters to several centimeters. Some lesions have a serpiginous border, suggesting an active edge. The skin's texture looks lichenified in some places, indicating chronicity, and scaling is evident across various regions, signaling some level of dryness and exfoliation. Some patches have merged, forming a larger area of affected skin. Signs of excoriations are present, most likely due to itching, and scattered pustules can also be observed."
**Probable Diseases**: "Prurigo nodularis, Chronic eczema, Psoriasis, Lichen simplex chronicus, Allergic or irritant contact dermatitis"
**Clinical Observation of the Skin Condition**: "The skin condition, as shown in the images, presents widespread erythematous patches with violaceous hues across the leg. The patient has multiple crusted plaques and erosions, with sizes varying from a few millimeters to several centimeters. Some lesions have a serpiginous border, suggesting an active edge. The skin's texture looks lichenified in some places, indicating chronicity, and scaling is evident across various regions, signaling some level of dryness and exfoliation. Some patches have merged, forming a larger area of affected skin. Signs of excoriations are present, most likely due to itching, and scattered pustules can also be observed."

**Ground Truth**: Chronic Eczema

**Excerpt 1 - GroupChat Initialization:** The group chat is initiated by the admin. The task, meticulously crafted for our use case, is provided to the chat manager who then follows the outlined steps.

**Excerpt 2 - Evidence & Critiques:** After the coordinator assigns probable diseases to each diagnostic specialist, they are sequentially called to generate supportive evidence and critique other possibilities.

**Excerpt 3 - Consolidated Evidence & Critiques:** Once the diagnostic specialists have processed their assigned diagnoses, the coordinator gathers and consolidates the evidence and critiques

for each disease.
**Excerpt 4 - Admin Refinement Instructions:** This consolidated evidence and critiques are reviewed by the admin, who assesses them and may request further information. The admin provides instructions for specialists to refine their evidence, aiming for a more accurate diagnosis.
**Excerpt 5 - Enhanced Evidence:** Based on the refinement instructions from the Admin, the designated agent is tasked with enhancing their evidence to better support their diagnosis.
**Excerpt 6 - Final Diagnosis:** Based on the refined evidence, the admin determines the most relevant final diagnosis.

## A.4 Dermatology Guidelines

When a dermatologist evaluates a skin condition, they typically follow a systematic approach that involves several areas.

- Patient History: Look at the "User Query" to extract relevant context that will help in accurate diagnosis of skin conditions.
- Visual Inspection: The initial step involves a thorough visual examination of the affected area.

For visual inspection, the dermatologist looks at the following features and creates a list of possible skin conditions that match the visual features.

1. Size: What is the size of the skin lesions? Is it small or large?
2. Shape: What is the shape of the lesions?
3. Color: What is the color of the skin lesions?
4. Location: Where is the skin lesion or rash located?
5. Distribution Pattern: What is the distribution pattern, is it localized or widespread?
6. Existence of symmetry: Are the lesions symmetric?
7. Borders: Do the edges of the lesion appear sharp, irregular, or blurred?
8. Elevation: Is the lesion is flat, raised, or depressed below the skin surface?
9. Texture: Does the surface quality looks smooth/scaly/rough/soft/hard.

Pattern Recognition: Dermatologists are trained in recognizing patterns that certain skin conditions commonly present. These patterns, combined with the other collected information, help in forming a

| Excerpt - 1: |
| --- |
| **Chat Initialisation with the Admin:**<br>Admin (to chat_manager):<br>The objective is to facilitate a systematic discussion to diagnose a patient's skin condition, based on the clinical observation of the skin condition (The skin condition, as shown in the images, presents widespread erythematous patches with violaceous hues across the leg. The patient has multiple crusted plaques and erosions, with sizes varying from a few millimeters to several centimeters. Some lesions have a serpiginous border, suggesting an active edge. The skin's texture looks lichenified in some places, indicating chronicity, and scaling is evident across various regions, signaling some level of dryness and exfoliation. Some patches have merged, forming a larger area of affected skin. Signs of excoriations are present, most likely due to itching, and scattered pustules can also be observed.), case study (Please help take a look, what kind of skin disease is this? Suffering from the disease for more than 10 years. It is recurrent and is very itchy! It happens wherever I scratch in some places.) and a list of probable diseases (['prurigo nodularis', 'chronic eczema', 'psoriasis', 'lichen simplex chronicus', 'allergic or irritant contact dermatitis']). Give equal emphasis to the clinical observation as well as the symptoms and the disease progression to ensure accuracy. Assume that no further information or biopsy will be conducted.<br>"""<br>**Process Overview:**<br>"""<br>"""<br>**Detailed Steps:**<br>"""<br>"""<br>**Rules for Effective Process Management:**<br>"""<br>**Once "TERMINATE" has been returned no further conversation must happen** |

Table 7: Excerpt 1 : Group chat initialisation

preliminary diagnosis.

Differential Diagnoses: Based on the evaluation, the dermatologist will consider a list of possible conditions (differential diagnoses) and rule them out one by one, based on the evidence and test results.

**Excerpt - 2:**

Rick (to chat_manager):
As Diagnostic Specialist Rick,
Assigned Diagnosis: Prurigo Nodularis.
Supporting Evidence for Prurigo Nodularis: The indications of chronic scratching or rubbing like lichenification suggest that the rash could be Prurigo Nodularis. The patient's description of the condition as being very itchy and recurrent over a span of more than a decade also aligns with this diagnosis. In addition, the presence of scattered pustules can also be seen in cases of Prurigo Nodularis.
Critiques for Other Diseases:
1. **Chronic Eczema**: This condition also presents an itchy rash that can become lichenified from chronic scratching, but usually has a more defined pattern of eruption that is not described here.
2. **Psoriasis**: While this condition presents chronic plaques, they usually exhibit a characteristic silvery scale. No such description is provided here.
3. **Lichen Simplex Chronicus**: While this condition is characterized by lichenification, it generally affects a specific region rather than having a widespread distribution as described here.
4. **Allergic or Irritant Contact Dermatitis**: These conditions generally present rapidly after contact with an offending substance, which does not fit the chronic, decade-long presentation described here.

**CALL COORDINATOR** to validate completion.

Table 8: Excerpt 2: Evidence & Critiques by Diagnostic Specialist Rick.

**Excerpt - 3:**

Coordinator (to chat_manager):
Thank you for your findings, Michael.
As the Coordinator, I acknowledge the completion of Michael's analysis regarding 'allergic or irritant contact dermatitis'.
**Compiling Findings**
I'll now compile and categorize the generated evidences and critiques for each potential disease:
1. **Prurigo Nodularis**
- Supporting Evidence: Chronic scratching or rubbing and recurrent itchiness over a decade. Presence of scattered pustules.
- Consolidated Critiques: More defined pattern of eruption for Chronic Eczema is missing. No silver scales as in Psoriasis. More localized presentation expected as in Lichen Simplex Chronicus. No acute flare-ups post exposure as in Contact Dermatitis.
2. **Chronic Eczema**
- Supporting Evidence: Widespread and merging erythematous patches, lichenified skin at some regions and itchiness wherever scratched.
- Consolidated Critiques: Prurigo Nodularis typically presents hard, itchy lumps. Silvery scales typical to Psoriasis missing. Lichen Simplex Chronicus is usually localized.
3. **Psoriasis**
- Supporting Evidence: Widespread plaques, lichenified skin and scaling.
- Consolidated Critiques: In Prurigo Nodularis itching is more generalized. No typical eczema characteristics like weeping, oozing. Lichen Simplex Chronicus usually is localized.
4. **Lichen Simplex Chronicus**
- Supporting Evidence: Lichenification of skin due to chronic rubbing, intense itchiness, especially in localized areas, plus the serpiginous border.
- Consolidated Critiques: Larger areas of affected skin unlike Prurigo Nodularis. No oozing or crusting unlike Chronic Eczema. Silvery scale of Psoriasis missing. Contact Dermatitis usually presents acute symptoms post exposure.
5. **Allergic or Irritant Contact Dermatitis**
- Supporting Evidence: Chronic dryness, exfoliation, and itching response to certain irritants.
- Consolidated Critiques: Doesn't describe nodules typical to Prurigo Nodularis. Not enough weeping and crusting for Chronic Eczema. No silvery scaly plaques like Psoriasis. Lichen Simplex Chronicus is usually localized.

As the Coordinator, I present the compiled evidence to the Admin for a final evaluation on this patient's skin condition.

Table 9: Excerpt 3: Consolidated Evidences & Critiques

**Excerpt - 4:**

Admin (to chat_manager):
As the Admin,
Thank you, Coordinator, for compiling the findings.
Firstly, it's noticeable that the given evidences for "Allergic or Irritant Contact Dermatitis" and "Psoriasis" are relatively weak compared to other diagnoses. I agree with the critiques that the chronicity and absence of typical defining characteristics such as acute flare-ups in the case of dermatitis and silvery scales in the case of psoriasis make these diagnoses less likely.
Focusing on "Prurigo Nodularis", the evidence is plausible but the image description does not mention the typical nodular lumps that are characteristic of this disease.
"Lichen Simplex Chronicus" has supportive evidence that strongly aligns with chronic itching and lichenification, yet the widespread presence of the disease contradicts the usual localized occurrence of this condition.
"Chronic Eczema" also aligns well with the majority of described symptoms, including itching, chronicity, and lichenification. The absence of strong objections in the critiques and aforementioned aligning symptoms gives strength to this diagnosis.
Considering all compiled evidence and critiques, it seems that Chronic Eczema might be the most likely diagnosis given the available information. However, some uncertainty remains due to overlapping symptoms with other conditions and the lack of additional diagnostic tests.
Diagnostic Specialist Sam, I would like you to enhance your evidence for 'Chronic Eczema' in light of the critiques provided by other specialists. Please refer to the following critiques and provide more specific details that distinguish Chronic Eczema from other conditions:
1. Prurigo Nodularis: Your analysis could benefit from addressing whether or not the absence of nodules is indeed conclusive evidence against this diagnosis.
2. Lichen Simplex Chronicus: Can you further explain the distinguishing factors between these two conditions? Specifically, consider the details regarding distribution and impact of itch-induced scratching.
3. Identification or ruling out of 'Allergic or Irritant Contact Dermatitis': Please provide more info that can make this differentiation clearer.

Table 10: Excerpt 4: Admin Refinement Instructions

Excerpt - 5:

Sam (to chat_manager):
As Diagnostic Specialist Sam,
Enhanced Evidence for 'Chronic Eczema':
To revisit and strengthen my diagnosis supporting Chronic Eczema, let's address the critiques:
1. Absence of nodules in Prurigo Nodularis: Prurigo Nodularis is characterized by hard, itchy nodules which may be paired with lichenification because of chronic scratching. However, such nodules are not explicitly reported in this clinical presentation. Instead, we note widespread erythematous patches with various sizes and serpiginous borders, a pattern more consistent with chronic eczema.
2. Distinguishing features between Chronic Eczema and Lichen Simplex Chronicus: Though both conditions show lichenification due to chronic scratching, they do have differing behaviors. Lichen Simplex Chronicus usually exhibits itself in one or two specific regions of the body, whereas Chronic Eczema can affect larger, more widespread areas as described in this clinical case. Thus, the widespread distribution here lends more credence to a diagnosis of Chronic Eczema rather than the typically localized Lichen Simplex Chronicus.
3. Differentiating between Chronic Eczema and Allergic or Irritant Contact Dermatitis: Contact Dermatitis generally surfaces as an acute flare-up following exposure to a particular substance and often resolves once the irritant or allergen is avoided, whereas Chronic Eczema's cause is multifactorial - influenced not only by external irritants but also by internal factors, such as the patient's immune response. Furthermore, Chronic Eczema exhibits a distinctive pattern of flares and subsiding inflammation over time. This history of enduring for over a decade and recurrent nature of the skin condition directs more towards Chronic Eczema.

Overall, despite sharing common symptoms like itching and skin alterations with the mentioned conditions, Chronic Eczema appears to fit best given the specifics of the condition's distribution and chronicity.

Table 11: Excerpt 5: Enhanced Evidence

Excerpt - 6:

Admin (to chat_manager):
As the Admin,
Thank you, Sam, for the enhanced evidence and Coordinator for your facilitation.
Taking into account the evidences and critiques from all Diagnostic Specialists and the enhanced evidence provided by Sam, I conclude that in the absence of any further diagnostic tests or additional information, the most accurate diagnosis among the provided probable diseases for this case is Chronic Eczema. The patient's long term history, reported symptoms like itching wherever the patient scratches, the recurrent nature of the condition, and clinical information like widespread erythematous patches all point towards Chronic Eczema.

Please **TERMINATE** the conversation, Coordinator.

Table 12: Excerpt 6: Final Diagnosis

| | |
|---|---|
| **Naive CoT** | Act as an expert dermatologist.<br>You are provided with a dermatology case with images. You are tasked to create a list of possible skin conditions for the given case.<br><br>*Instructions:*<br>1: Look at the images and the medical query and see what relevant information you can extract from the medical query that can be useful in diagnosis.<br>2: Create a possible list of skin conditions. |
| **Medical Guidelines based CoT** | You are provided with a dermatology case with images.<br><br>For this case, you are provided with some images and additional user query. You are asked to give a diagnosis for this scenario. Act as a dermatologist. Refer the guidelines below and follow the guidelines to generate the diagnosis.<br><br>*Guidelines:*<br>When a dermatologist evaluates a skin condition, they typically follow a systematic approach that involves several key steps.<br><br>*Visual Inspection:* The initial step involves a thorough visual examination of the affected area.<br>The dermatologist looks at the:<br>   1. Size<br>   2. Shape<br>   3. Color: The color (red, brown, black, blue, white) and whether it's uniform.<br>   4. Location of the lesion or rash.<br>   5. Distribution Pattern (localized/widespread)<br>   6. Existence of symmetry (yes or no)<br>   7. Borders: The edges of the lesion—are they sharp, irregular, or blurred?<br>   8. Elevation: Whether the lesion is flat, raised, or depressed below the skin surface.<br>   9. Texture: The surface quality (smooth, scaly, rough, soft, hard).<br>   10. Pattern Recognition: Dermatologists are trained in recognizing patterns that certain skin conditions commonly present. These patterns, combined with the other collected information, help in forming a preliminary diagnosis.<br>   11. Consideration of Differential Diagnoses: Based on the evaluation, the dermatologist will consider a list of possible conditions (differential diagnoses), and rule them out one by one, based on the evidence and test results.<br>   13. Create a list of possible candidates after the above steps. |

Table 13: CoT based Retrieval

Act as an expert dermatologist.
You are provided with a dermatology case with images and associated medical query. You are tasked to choose the most probable skin condition from the set of candidates.

Medical Query:
*<query>*

Candidates:
*<candidates>*

Instructions:
1: Look at the images and the medical query and see what relevant information you can extract from the medical query that can be useful in diagnosis.
2: Give a score to each candidate skin condition in the range of 1 - 10 with 1 being the least probable and 10 being the most probable.
3: Choose a single most probable disease. If there is a tie in scores for the most probable conditions, pick a single skin condition between those candidates at random and return.

Table 14: Naive CoT for Re-Ranker

Act as an expert dermatologist. You are provided with a dermatology case. For this case, you are provided with some images a user query and list of candidates.

User Query: *<query>*

Candidates: *<candidates>*

Guidelines:
When a dermatologist evaluates a skin condition, they typically follow a systematic approach that involves several areas.
Patient History: Look at the "User Query" to extract relevant context that will help in accurate diagnosis of skin conditions.
Visual Inspection: The initial step involves a thorough visual examination of the affected area.

For visual inspection, the dermatologist looks at the following features and for each, the dermatologist creates a list of possible skin conditions that show such visual features.
1: Size: What is the size of the skin lesions? Is it small or large?
2: Shape: What is the shape of the lesions?
3: Color: What is the color of the skin lesions?
4: Location: Where is the skin lesion or rash located?
5: Distribution Pattern: What is the distribution pattern, is it localized or widespread?
6: Existence of symmetry: Are the lesions symmetric?
7: Borders: Do the edges of the lesion appear sharp, irregular, or blurred?
8: Elevation: Is the lesion is flat, raised, or depressed below the skin surface?
9: Texture: Does the surface quality looks smooth/scaly/rough/soft/hard.

Pattern Recognition: Dermatologists are trained in recognizing patterns that certain skin conditions commonly present. These patterns, combined with the other collected information, help in forming a preliminary diagnosis.

Differential Diagnoses: Based on the evaluation, the dermatologist will consider a list of possible conditions (differential diagnoses) and rule them out one by one, based on the evidence and test results.

Instructions:
step 1: Evaluate the medical images based on Visual Inspection Guidelines.
step 2: Evaluate the medical query as Patient History Guidelines.
step 3: Create a case summary using information extracted at step 1 and step 2.
step 4: For each candidate skin condition present in the list of Candidates, give a score on a scale of 1 - 10 (where 1 is the least probable and 10 is the most probable) that describes how likely is the given skin condition as a diagnosis for the case summary.
step 5: Return the two most probable skin candidates based on scores obtained at step 4.

Table 15: Expert Guidelines Grounded CoT with Context for Re-Ranker

You are provided with a dermatology case. For this case, you are provided with some images and list of possible candidates.

Candidates: *<candidates>*

Visual Inspection Guidelines: The initial step involves a thorough visual examination of the affected area. The dermatologists keep a track of 10 visual features.

1: Size: What is the size of the skin lesions? Is it small or large?
2: Shape: What is the shape of the lesions?
3: Color: What is the color of the skin lesions?
4: Location: Where is the skin lesion or rash located?
5: Distribution Pattern: What is the distribution pattern, is it localized or widespread?
6: Existence of symmetry: Are the lesions symmetric?
7: Borders: Do the edges of the lesion appear sharp, irregular, or blurred?
8: Elevation: Is the lesion is flat, raised, or depressed below the skin surface?
9: Texture: Does the surface quality looks smooth/scaly/rough/soft/hard.
10: Pattern Recognition: Dermatologists are trained in recognizing patterns that certain skin conditions commonly present. These patterns, combined with the other collected information, help in forming a preliminary diagnosis.

Act as a dermatologist and follow the instructions below:

Instructions:
step 1: For the given images, use the guidelines and generate a visual description.
step 2: For each candidate in the "Candidates", generate the visual description that describes the candidate disease. Also mention distinguishing features based on visual guidelines. Include features like shape, colours, lesion type and area of localization to create a visual description for the disease.
step 3: Compare the visual description which was generated for each candidate skin condition at step b with the image description generated at step a. Give a score in the range of 1 to 10 with 1 being the lowest match and 10 being the highest match.
step 4: Choose the most probable candidate which has the highest score with the images based on step 3.

Table 16: Expert Guidelines Grounded CoT without Context for Re-Ranker

**Rules:**

1. Simplify and Be Direct
- Example: "The condition is Chronic Eczema."
- Explanation: Human expert responses tend to be direct and use simpler language. Avoid overly complex explanations and aim for straightforward answers directly addressing the patient's inquiry.

2. Diagnosis Confirmation
- Example: "Your diagnosis is a Myxoid Cyst based on the clear image provided."
- Explanation: Include statements that confirm the diagnosis confidently, as seen in responses like "Chronic Eczema." or "It is myxoid cyst." Use assertive language to convey confidence in your diagnosis.

3. Detail Symptom Correlation
- Example: "The semi-spherical cyst near the end of your thumb, as described, leads to a diagnosis of Myxoid Cyst."
- Explanation: Explicitly connect the diagnosis with observed symptoms or test results when applicable, similar to the detailed descriptions in some valid responses. This helps patients understand why a particular diagnosis is made.

4. Incorporate Treatment Options Clearly
- Example: "For Psoriasis, I recommend oral capsules such as glycyrrhizic acid glycosides, along with transfer factors." - Explanation: When suggesting treatments, mention specific medications or procedures clearly and concisely, as observed in responses with high completeness. If possible, explain the purpose of each treatment briefly.

5. Mention Commonality or Prevalence
- Example: "Chronic Eczema is quite common and effectively manageable with the right treatment."
- Explanation: If applicable, include a brief note on how common the condition is or any relevant statistical information that could reassure the patient or provide context, akin to how some expert responses include prevalence information.

6. Use Patient-Friendly Language
- Example: "Based on the photo you provided, it looks like you have a Myxoid Cyst, which is a fluid-filled lump that's not harmful."
- Explanation: Ensure the language used is patient-friendly, avoiding unnecessary medical jargon that could confuse the patient. When medical terms are unavoidable, consider providing a brief, simple explanation.

7. Personalization and Empathy
- Example: "I understand that dealing with Chronic Eczema can be frustrating. Regular moisturizing and the treatments we've discussed should offer relief." - Explanation: Whenever possible, personalize the response to the patient's situation. Display empathy to make your responses feel more human and less robotic.

Table 17: Automatic Prompt Optimization (APO) Rules

**Rules:**

1: Skin condition A is similar to B if they have same name.

2: Skin condition A is similar to B if B is also known by the name A.

3: Skin condition A is similar to B if both are part of the same root skin condition.
Example: Herpetic Eczema and seborrheic eczema are similar since they have same root, Eczema.

4: Skin condition A is similar to B if they are both have the same effect and share a common cause.

Table 18: Evaluation Guidelines Rules

# KU-DMIS at MEDIQA-CORR 2024: Exploring the Reasoning Capabilities of Small Language Models in Medical Error Correction

**Hyeon Hwang**[1]   **Taewhoo Lee**[1]   **Hyunjae Kim**[1]   **Jaewoo Kang**[1,2]

[1]Korea University   [2]AIGEN Sciences

{hyeon-hwang,taewhoo,hyunjae-kim,kangj}@korea.ac.kr

## Abstract

Recent advancements in large language models (LM) like OpenAI's GPT-4 have shown promise in healthcare, particularly in medical question answering and clinical applications. However, their deployment raises privacy concerns and their size limits use in resource-constrained environments. Smaller open-source LMs have emerged as alternatives, but their reliability in medicine remains underexplored. This study evaluates small LMs in the medical field using the MEDIQA-CORR 2024 task, which assesses the ability of models to identify and correct errors in clinical notes. Initially, zero-shot inference and simple fine-tuning of small models resulted in poor performance. When fine-tuning with chain-of-thought (CoT) reasoning using synthetic data generated by GPT-4, their performance significantly improved. Meerkat-7B, a small LM trained with medical CoT reasoning, demonstrated notable performance gains. Our model outperforms other small non-commercial LMs and some larger models, achieving a 73.36 aggregate score on MEDIQA-CORR 2024.

## 1 Introduction

Large language models (LM) have recently made significant advancements, finding usefulness across diverse applications in healthcare and medicine (Thirunavukarasu et al., 2023; Tian et al., 2024). For instance, OpenAI's GPT-3.5 (Brown et al., 2020) and GPT-4 (Achiam et al., 2024) have demonstrated their capabilities by achieving remarkable accuracy on standardized tests like the United States Medical Licensing Examination (USMLE). They have also shown excellence in real-world clinical applications—from responding to queries to diagnosing complex cases (Kung et al., 2023; Nori et al., 2023a,b; Singhal et al., 2023a,b).

However, deploying proprietary LMs in this sensitive sector presents significant challenges, primarily due to privacy concerns and the need for secure



Figure 1: Overview of our proposed method. (a) In chain-of-though (CoT) dataset generation using GPT-4, we feed GPT-4 with clinical notes, error sentences, and correct sentences to generate CoT explanation that articulates error and correction. (b) In supervised fine-tuning, we fine-tune Meerkat-7B (Kim et al., 2024) with generated dataset to enhance its error detection and correction capabilities.

data handling (Thirunavukarasu et al., 2023; Li and Zhang, 2017; Meskó and Topol, 2023; Bartoletti, 2019). Since these models rely on APIs, it can be hard to use them in hospitals where a significant amount of sensitive personal information is present. Moreover, their vast computational requirements make them impractical for deployment on local servers in hospitals or medical research centers.

For these reasons, smaller open-sourced LMs are emerging as alternatives. For instance, models such as Mistral (Jiang et al., 2023) and BioMistral (Labrak et al., 2024) come with manageable sizes that are more suitable for deployment on local servers, while mitigating the security issues. How-

ever, because these models have significantly fewer parameters (typically 7B) compared to large LMs (more than 100B), there are doubts about whether these models can provide factual responses based on their parametric knowledge. This necessitates rigorous verification before being deployed especially in the medical domain, where reliability is of utmost importance.

In this paper, we evaluate the reliability of small LMs in the medical domain. For this purpose, we utilize the MEDIQA-CORR 2024 shared task (Ben Abacha et al., 2024a), which tasks models with identifying potential errors in clinical notes and correcting them. This task assesses the ability of models to address common medical sense errors, enabling us to verify their reliability and identify hallucinations in small language models.

Our initial experiment found that when small LMs were evaluated in a zero-shot setting or trained using training data through simple supervised fine-tuning, their performance fell short of expectations. Notably, the scores were similar to random guessing in a binary classification task. This result suggests solving complex medical problems is challenging for small models lacking advanced reasoning capabilities.

Thus, we hypothesized that fine-tuning the model with chain-of-thought (CoT) reasoning (Wei et al., 2022) could effectively equip the model with these necessary reasoning capabilities. To implement this, we generated reasoning paths between the inputs and outputs of the training dataset using GPT-4 and then trained the model not only to generate correct answers but also to provide the underlying reasoning for each decision (Figure 1). This approach resulted in noticeable performance improvements, confirming the critical role of CoT reasoning in solving complex medical problems.

Furthermore, we observed that small LMs could benefit from reasoning capabilities aquired from other tasks. Specifically, Meerkat-7B (Kim et al., 2024), trained on an extensive medical CoT reasoning dataset for USMLE-style questions, showed greater performance improvements compared to other small LMs. This significant improvement highlights the importance of reasoning capabilities for small LMs to generate reliable responses.

Using this approach, we achieved an aggregate score of 73.36 for the natural language generation (NLG) evaluation, 63.46 for binary classification accuracy in detecting the presence of an error (error flag accuracy), and 61.51 for accuracy

in identifying the specific sentence containing the error (error sentence accuracy) on the test set. Despite its much smaller size relative to proprietary Large LMs, Meerkat-7B demonstrated competitive performance in the MEDIQA-CORR 2024 shared task, achieving the best score among non-commercial/small LMs. This achievement is particularly significant considering the dominance of GPT-4-based frameworks among other teams.

## 2 Methods

### 2.1 Task Formulation

MEDIQA-CORR 2024 (Ben Abacha et al., 2024a) involves identifying medical errors in clinical notes and correcting them. This task is broken down into the following three sub-tasks: (1) binary classification, determining whether the clinical text contains a medical error or not, (2) span identification, detecting the specific text span associated with the error, and (3) natural language generation (NLG), creating a corrected version of the text in a free-form format. Sub-tasks 2 and 3 are performed only when an error exists in the given note.

In this study, we frame the task around generative models that produce free-form text as output. Let $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\}$ be a dataset, where $N$ is the dataset size, $\mathbf{X} = \{x_n\}_{n=1}^N$, and $\mathbf{Y} = \{y_n\}_{n=1}^N$. The $n$-th clinical note, denoted as $x_n$, is composed of several sentences structured as follows:

$$x_n = \{s_1, s_2, \ldots, s_{T_n}\}, \tag{1}$$

where $T_n$ is the number of sentences within the note, and $s_i$ is the $i$-th sentence ($i \in \{1, \ldots, T_n\}$). The label set $y_n$ consists of an error flag $e \in \{0, 1\}$, error sentence index $k_n$ (if available), and corrected sentence $s^*$ (if available).

We input the entire text $x_n$ to the model, and in return, the model outputs its structured response as shown in Table 1. We parsed the model's structured response to construct the output format. If the model predicts 'the note does not contain an error.', we set the error flag to 0 and fill the fields for error sentence index and corrected sentence with '-1' and 'NA', respectively. Conversely, if the model indicates the presence of an error, we set the error flag to 1 and record both the error sentence index and the corrected sentence. The final submission format for the output is structured as follows:

$$\text{output} = \begin{cases} (e, k_n, s^*) & \text{if } e = 1 \\ (e, -1, \text{`NA'}) & \text{otherwise,} \end{cases} \tag{2}$$

```
<INPUT>
You are an expert tasked with providing a logical explanation as to whether there is an error in the given clinical
note. Your job is to analyze the clinical note step-by-step and provide an explanation leading to the conclusion
regarding the presence or absence of an error. You are strongly recommended to follow the output format: At the
end of your response, without modifications, use the phrase "Therefore, the error sentence {ERROR_SENT} should
be corrected to the corrected sentence {CORRECT_SENT}." or "Therefore, the note does not contain an error."

{NOTE}

ASSISTANT:

<OUTPUT>
{CoT Reasoning}

Therefore, the error sentence {ERROR_SENT} should be corrected to {CORRECT_SENT}.
or
Therefore, the note does not contain an error.
```

Table 1: Input and output format of our CoT dataset that was used to fine-tune small language models. The output
format was specifically structured to simplify the parsing process.

## 2.2 Generating Reasoning Chains

We instructed GPT-4 to conduct a thorough analysis
of clinical notes and provide explanations as to
whether the given note potentially contains an error
or not. Specifically, we use a clinical note, $x$, and
an error flag, $e$, to prompt the model. When the note
contains an error, we also provide both the error
sentence $\hat{s}$ and the corrected sentence $s^*$; otherwise,
we only provide the error flag as follows:

$$r = \begin{cases} \text{gpt}_{\text{err}}(x, \hat{s}, s^*) & \text{if } e = 1 \\ \text{gpt}_{\text{no}}(x) & \text{otherwise,} \end{cases} \quad (3)$$

where $r$ is the generated reasoning chain, and $\text{gpt}_{\text{err}}$
and $\text{gpt}_{\text{no}}$ are the OpenAI API functions with the
pre-defined input prompts. Figure 2 details the
input prompts for error and non-error examples. In
our initial experiments, we observed that when we
did not provide label information to the model and
instead asked it to determine the presence of errors
and correct them, the model often gave incorrect
predictions; therefore, we provided gold-standard
labels to increase the recall rate of the reasoning
data. An example of the reasoning chain generated
by GPT-4 can be seen in Figure 3.

We generated five different reasoning paths for
each example to supplement the limited amount of
data. After filtering out samples that did not follow
the specified output format, we obtained 9,712 and
3,207 examples from the training set and validation
set, respectively. This generated dataset was piv-
otal in training our model, as it helped enhance the
model's reasoning capabilities and as well as per-
formance in correcting errors in clinical notes. The
fine-tuning process enabled the model to generate

explanations as coherent and contextually appropri-
ate as those produced by GPT-4.

## 2.3 Model

As our backbone model, we used Meerkat-7B (Kim
et al., 2024)[1] because it is specifically designed to
handle complex medical queries through advanced
multi-step reasoning. Built on Mistral-7B (Jiang
et al., 2023), Meerkat-7B has been trained on a
high-quality medical instruction-tuning dataset in-
cluding extensive synthetic USMLE-style ques-
tions from 18 medical textbooks and corresponding
CoT reasoning paths. The questions and CoT rea-
soning paths are generated by GPT-4, thereby en-
dowing the model with distilled medical knowledge
and reasoning capabilities from GPT-4. Leverag-
ing these characteristics, Meerkat-7B has achieved
state-of-the-art performance across various med-
ical question-answering benchmarks that require
complex reasoning.

## 2.4 Training and Inference

We adopted supervised fine-tuning to fine-tune a
language model using our reasoning dataset. For a
given clinical note, the model was trained to gen-
erate a reasoning path $r$ first, and then structured
output as shown in Table 1.

During inference, we employed a self-
consistency method (Wang et al., 2023) to mitigate
potential instability in the outputs generated by
a single model. This method, often used as an
ensemble technique, helps aggregate predictions

---
[1]https://huggingface.co/dmis-lab/
meerkat-7b-v1.0

528

```
[Error Example]
...

We will let you know if there's an error and
pinpoint its location if found. Your job is
to analyze the clinical note step-by-step and
provide  an  explanation  leading  to  the
conclusion regarding the presence or absence
of an error. Below is the input problem:

# Clinical Note
{NOTE}

# Conclusion
Error sentence: {ERROR_SENT}
Corrected sentence: {CORRECT_SENT}


...

# Your explanation:
```

```
[Non-error Example]
...

We will let you know if there's an error and
pinpoint its location if found. Your job is
to analyze the clinical note step-by-step and
provide  an  explanation  leading  to  the
conclusion regarding the presence or absence
of an error. Below is the input problem:

# Clinical Note
{NOTE}

# Conclusion
No error in the given note.


...

# Your explanation:
```

Figure 2: The input prompts for generating CoT reasoning paths from error (left) and non-error (right) examples using GPT-4. These prompts guide GPT-4 through a detailed analysis of a clinical note to determine and explain the presence and the absence of errors within step-by-step reasoning.

from generative language models. The model generated 30 separate outputs for each input and then these outputs are aggregated to determine the most reliable result. If 'Therefore, the note does not contain an error.' is the predominant output, it is interpreted that the clinical note contains no errors. Conversely, if a specific corrected sentence emerges as the most consistent across the outputs, that sentence is selected as the final correction. This strategy reduces the impact of potentially erroneous outputs by leveraging the consensus from multiple outputs.

## 3 Experimental Settings

In all our experiments, we utilized eight 80GB NVIDIA A100 GPUs. When fine-tuning, we used a learning rate of 1e-6 and a batch size of 128.[2] For generating the CoT dataset, we used GPT-4 Turbo (gpt-4-1106-preview) through the OpenAI API.

### 3.1 Dataset

For our experiments, we utilized the official dataset (Ben Abacha et al., 2024b) provided by the MEDIQA-CORR 2024 shared task. Table 2 details the number of samples in each split. We used the training set for initial model tuning and selecting the best model and hyperparameters based on validation performance. For the final submis-

---

[2]We tested a range of learning rates, {1e-7, 5e-7, 1e-6, 5e-6, 1e-5}, and picked the best one based on performance on the MS validation set.

| Dataset | Training | Validation | Test |
|---------|----------|------------|------|
| MS      | 2,189    | 574        | 597  |
| UW      | -        | 160        | 328  |

Table 2: Statistics of the MEDIQA-CORR 2024 dataset. The training and validation sets were provided for model development, whereas the test split was specifically designated for the official evaluation during the challenge.

sion of the test set, the model was trained using a combination of the training and validation sets.

### 3.2 Metrics

For binary classification, we used error flag accuracy to evaluate whether the model accurately determines if a clinical text contains a medical error. We used error sentence detection accuracy for span identification to evaluate whether the model accurately outputs the index of the error sentence.

For NLG, we utilized the following evaluation metrics: ROUGE (Lin, 2004), which measures the overlap of ngrams between the generated text and the reference; BERTScore (Zhang et al., 2020), which evaluates semantic similarity using BERT embeddings; and BLEURT (Sellam et al., 2020), which assesses text generation quality based on a learned metric. Additionally, we used an AggregateScore, calculated as the arithmetic mean of ROUGE-1, BERTScore, and BLEURT. Note that these NLG evaluation metrics are computed when the model corrects an error sentence in the clinical

note that contains an error.

# 4 Results

## 4.1 Effect of Medical Reasoning on Clinical Note Correction

To verify the impact of fine-tuning with medical reasoning on clinical note correction, we evaluated three small LMs—Mistral-7B (Jiang et al., 2023), BioMistral-7B (Labrak et al., 2024), and Meerkat-7B (Kim et al., 2024)—using three methods: zero-shot CoT, fine-tuning with CoT reasoning, and fine-tuning without CoT reasoning.

Table 3 demonstrates that zero-shot CoT models exhibited poor accuracy and NLG evaluation results compared to models fine-tuned with CoT reasoning. Specifically, Mistral-7B performed worse than a random guess in the binary classification task, and BioMistral-7B largely failed to adhere to the output formats suggested in the prompts. Meerkat-7B demonstrated relatively strong performance, but there was considerable room for improvement. When fine-tuning Meerkat-7B with CoT reasoning, the performance improved by 33.51% in AggregateScore (AS), indicating that the model requires fine-tuning to adapt effectively to the target task.

In fine-tuning settings, models trained on the CoT dataset notably outperformed those trained without CoT reasoning in all metrics. Specifically, Meerkat-7B showed substantial improvements when trained with CoT reasoning: error flag accuracy increased by 9.23%, error sentence detection by 10.28%, AggregateScore by 3.36%. The result highlights the crucial role of medical reasoning in enhancing the reliability and performance of small LMs for medical domain problems.

Meerkat-7B, which was extensively trained on question-answering CoT data to enhance its complex reasoning capabilities, significantly outperformed other small language models in terms of accuracy metrics and NLG evaluation results when fine-tuned with CoT. Specifically, Meerkat-7B exceeded both Mistral-7B and BioMistral-7B in error flag accuracy, with improvements of 5.75% and 8.71% respectively. It also scored higher on NLG aggregate scores, outperforming Mistral-7B by 3.08% and BioMistral-7B by 6.79%. These results are attributed to the transfer of complex medical reasoning skills, acquired from other tasks, to the task of clinical note correction.

## 4.2 Official Evaluation

Based on the observations in the previous sections, we selected Meerkat-7B as our backbone model for the final submission, affirming its effectiveness for tasks requiring complex medical reasoning. Table 4 shows the official test results in the MEDIQA-CORR 2024.[3] Among the fourteen final submissions, seven teams employed large models, predominantly GPT-4, and five teams used smaller models. Large LMs demonstrated superior performance in both accuracy and NLG evaluation metrics. However, the results indicate that Meerkat-7B achieves competitive outcomes compared to them. Despite having significantly fewer parameters, our model secured fourth place overall and was the top performer among open-source and smaller LMs.

Based on the official results, our model shows substantial error flag accuracy and error sentence detection accuracy compared to other models. Still, a 63.46% accuracy rate in binary classification suggests room for improvement. To enhance our performance in binary classification, we could consider adopting an encoder model (such as BERT (Devlin et al., 2019)) focused specifically on this classification task, rather than relying solely on a general decoder model. A Two-step approach using an encoder model and decoder model in each step may help address both binary classification and correction of error sentences.

Conversely, our model achieved strong results in NLG evaluation, indicating a robust capability to generate accurate corrected sentences within the context of identified errors. This highlights its effectiveness in detailed text generation and correction tasks within the clinical domain. Given these strengths, we can expect more promising usability of our model in tasks where error existences are known, enhancing its practical application in error correction scenarios.

## 4.3 Case Study

We present a case study comparing reasoning from different approaches, using an example from the validation dataset. Figure 4 provides example outputs from three approaches: zero-shot CoT from each Mistral-7B and Meerkat-7B, and fine-tuned Meerkat-7B with CoT reasoning.

The clinical note describes a rash that started behind the ears and spread to the trunk and ex-

---

[3]We obtained basic information about the submissions from the official result sheet, but please note that we do not have precise details about the submissions.

CLINICAL NOTE:

A 44-year-old man comes to the physician because of fatigue and increased straining during defecation for 3 months. During this time, he has lost 5 kg (12 lb) despite no change in appetite. He has a family history of colon cancer in his maternal uncle and maternal grandfather. His mother died of ovarian cancer at the age of 46. Physical examination shows conjunctival pallor. His hemoglobin concentration is 11.2 g/dL, hematocrit is 34%, and mean corpuscular volume is 76 μm3. Colonoscopy shows an exophytic mass in the ascending colon. Pathologic examination of the resected mass shows a poorly differentiated adenocarcinoma. Genetic analysis shows a mutation in the MSH2 gene. Patient was diagnosed with Turcot syndrome.

CoT REASONING:

The clinical note describes a 44-year-old man with symptoms of fatigue, increased straining during defecation, weight loss, and family history of colorectal and ovarian cancers.

. . .

Turcot syndrome, on the other hand, is a genetic syndrome that is associated with brain tumors and colorectal cancer. However, there is no mention of brain tumors in the clinical note, which would be more indicative of Turcot syndrome. Therefore, given the presence of a MSH2 mutation, family history of cancers associated with Lynch syndrome, and the absence of any mention of brain tumors, the diagnosis of Turcot syndrome seems less likely.

Considering these findings and the association between MSH2 mutation and Lynch syndrome, rather than Turcot syndrome, the corrected sentence is "Patient was diagnosed with Lynch syndrome."

ERROR SENTENCE:

Patient was diagnosed with Turcot syndrome.

CORRECTED SENTENCE:

Patient was diagnosed with Lynch syndrome.

Figure 3: Example of CoT reasoning generated by GPT-4. The CoT reasoning shows a detailed explanation process in which GPT-4 uses the provided clinical note, error sentence, and corrected sentence to construct a logical reasoning path leading to the appropriate correction.

tremities, accompanied by mild sore throat, red itchy eyes, and headache. It concludes with a diagnosis of measles. However, the rash pattern and postauricular and suboccipital lymphadenopathy align more closely with rubella, which typically presents with milder symptoms and lymph node swelling. Measles would more likely involve a cough and more severe conjunctivitis, which are not mentioned.

In a zero-shot CoT setting, Mistral-7B did not detect any error in the note due to insufficient reasoning, while Meerkat-7B accurately identified the error sentence 'The patient has measles,' noting the lack of adequate evidence to conclude that the patient has measles, through step-by-step reasoning. However, the model failed to correct the sentence, indicating that it is not fully adapted to the task. In contrast, the fine-tuned Meerkat-7B with CoT reasoning successfully corrected the clinical note. It suggested that rubella is more consistent with the patient's symptoms by providing appropriate supporting reasoning. This case study demonstrates that although Meerkat-7B exhibits relatively decent medical reasoning in error detection within clinical notes compared to other baselines, fine-tuning is necessary to tailor the model for the target task.

## 5 Related Works

### 5.1 Commonsense Detection

Commonsense detection refers to the ability of an AI system, to use basic knowledge about the world that is typically obvious to humans, to understand and respond appropriately in various situations. It has traditionally been explored within general domains, such as SemEval-2020 Task 4 on Commonsense Validation and Explanation (Wang et al., 2020) and the CREAK dataset (Onoe et al., 2021). Unlike these general applications, MEDIQA-CORR 2024 Shared Task (Ben Abacha et al., 2024a) is specifically focused on the medical domain, where the implications of errors are particularly critical. Medical texts require a high degree of expertise and knowledge to not only detect errors but also correct them appropriately. This focus emphasizes the need for AI systems that perform reliably and accurately in healthcare, where factuality directly affects patient care.

531

| | Accuracy Results | | NLG Eval Results | | | |
|---|---|---|---|---|---|---|
| **Model** | **EF** | **ES** | **R1** | **BS** | **BL** | **AS** |
| *Zero-shot CoT* | | | | | | |
| BioMistral-7B* | - | - | - | - | - | - |
| Mistral-7B | 48.95 | 35.89 | 17.81 | 25.97 | 36.56 | 26.78 |
| Meerkat-7B | **54.18** | **45.99** | **25.83** | **33.06** | **40.88** | **33.26** |
| *Fine-tuning w/o CoT reasoning* | | | | | | |
| BioMistral-7B | **52.61** | 47.74 | 49.09 | 57.27 | 50.77 | 52.38 |
| Mistral-7B | 48.78 | 46.86 | 61.01 | 66.63 | 58.83 | 62.16 |
| Meerkat-7B | 52.09 | **50.17** | **61.60** | **68.26** | **60.38** | **63.41** |
| *Fine-tuning w/ CoT reasoning* | | | | | | |
| BioMistral-7B | 52.61 | 51.22 | 56.55 | 65.82 | 57.58 | 59.98 |
| Mistral-7B | 55.57 | 54.70 | 61.07 | 68.93 | 61.08 | 63.69 |
| Meerkat-7B | **61.32** | **60.45** | **64.98** | **71.30** | **64.03** | **66.77** |

Table 3: Performance of small language models on the MS validation set, evaluated through three methods: zero-shot CoT, fine-tuning without CoT reasoning, and fine-tuning with CoT reasoning. Metrics include error flag accuracy (EF), error sentence detection accuracy (ES), ROUGE-1 (R1), BERTScore (BS), BLUERT (BL), and AggregateScore (AS). We did not evaluate BioMistral-7B in the zero-shot CoT method (marked with an asterisk(*)) because this model does not generate responses in the required format, making parsing impossible. Due to superior performance compared to other models, we chose Meerkat-7B as our base model for the final submission.

## 5.2 Biomedical Language Models

With the success of transformer-based models on various NLP tasks, ongoing research has focused on applying them to the medical domain. Different transformer architectures have been trained with large amounts of biomedical text to encapsulate domain-specific context, including encoder-decoder-based (Yuan et al., 2022; Phan et al., 2021), encoder-based (Lee et al., 2020; Gu et al., 2021), and decoder-based (Luo et al., 2022) architectures. More recently, models equipped with billions of parameters have opened the era of Large LMs, showing superior performance and generalizability compared to smaller models. In line with this trend, recent works (Singhal et al., 2023a) have deployed various training strategies that enable Large LMs to excel at highly complex biomedical tasks, such as MedQA (Jin et al., 2021).

## 5.3 Reasoning Distillation

LMs have shown to generate CoT reasoning steps that can benefit end task performance, but only when equipped with at least 100 billion parameters (Wei et al., 2022). To this end, recent works have focused on distilling reasoning chains derived from larger models to smaller models (Li et al., 2022; Magister et al., 2023). SOCRATIC CoT (Shridhar et al., 2023) suggests a two-step approach, where a *problem decomposer* model interacts with a *sub-problem solver* model to reach the final solution.

## 6 Conclusion

In this study, we explored the capabilities of small open-sourced language models in medical error correction and the effect of CoT reasoning on this problem. Our findings confirm that CoT reasoning capabilities are highly encouraged for the task of clinical note correction, especially for small LMs. Particularly, Meerkat-7B, initially trained to solve complex medical questions using an extensive CoT dataset, demonstrates superior performance compared to other open-sourced small LMs. Despite having far fewer parameters than proprietary large LMs, Meerkat-7B achieves competitive performance in clinical note correction. This underscores the potential of well-designed smaller models to handle demanding medical AI tasks effectively. In future research, there should be ongoing efforts to continuously improve small LMs to enhance the reliability and safety of automated systems in healthcare, paving the way for more

| Rank | Base Model | Model Size | Accuracy Results | | NLG Eval Results |
|:---:|:---|:---:|:---:|:---:|:---:|
| | | | EF | ES | AS |
| 1 | GPT-4 | Large | **86.49** | **83.57** | **78.91** |
| 2 | GPT-4 & Claude Opus | Large | 62.16 | 60.86 | 78.66 |
| 3 | GPT-4 | Large | 52.22 | 52.00 | 78.06 |
| 4 | Meerkat-7B (**Ours**) | Small | 63.46 | <u>61.56</u> | <u>73.36</u> |
| 5 | Palmyra | Small | 56.00 | 52.00 | 73.30 |
| 6 | OpenAI (Not Specified) | Large | 66.92 | 61.08 | 71.09 |
| 7 | GPT-4 | Large | 69.41 | 61.95 | 65.81 |
| 8 | OpenAI (Not Specified) | Large | 68.00 | 64.00 | 58.75 |
| 9 | GPT-4 | Large | 67.41 | 60.97 | 58.10 |
| 10 | GPT-4 | Large | 67.78 | 59.03 | 55.87 |
| 11 | GPT-4 | Large | 56.65 | 49.08 | 48.09 |
| 12 | BioMistral-7B | Small | 50.16 | 37.84 | 45.01 |
| 13 | BioMistral-7B | Small | 53.95 | 36.32 | 44.83 |
| 14 | BART & SVM | Small | <u>73.73</u> | 60.00 | 44.56 |

Table 4: Official evaluation on MEDIQA-CORR 2024, featuring metrics such as error flag accuracy (EF), error sentence detection accuracy (ES), and aggregate score (AS). The table lists each base model used and roughly categorizes them into 'Large' or 'Small' based on their parameter size. Ranks are determined based on the aggregate score (AS). The best performance in each metric is highlighted in **bold**, while the best performance among small models is <u>underlined</u>. Our Meerkat-7B-based model achieved an aggregate score of 73.37, outperforming all small models and several large model-based systems.

accurate and trustworthy medical assistants.

## Limitation

One limitation of our current approach is that the model's integration of external knowledge sources is not fully developed (e.g., knowledge base- or retrieval-augmented generation). While Meerkat-7B exhibits high-quality reasoning capabilities, it has not yet been optimized to incorporate external knowledge. Integrating this model with a retriever and utilizing biomedical knowledge sources could significantly improve its ability to solve complex cases while reducing the likelihood of generating hallucinations. Future work will explore adapting Meerkat-7B to harness external biomedical corpora, potentially increasing its accuracy and reliability.

## Acknowledgements

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Ivana Bartoletti. 2019. Ai in healthcare: Ethical and privacy challenges. In *Artificial Intelligence in Medicine*, pages 7–10, Cham. Springer International Publishing.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec

**CLINICAL NOTE:**

A 4-year-old boy is brought to the physician because of a progressive rash for 2 days. The rash started behind the ears and now involves the trunk and extremities. Over the past 4 days, he has had mild sore throat, red, itchy eyes, and headache. He was born at term and has been healthy except for recurrent episodes of tonsillitis and occasional asthma attacks. Two weeks ago, he was treated for tonsillitis that resolved with penicillin therapy. He immigrated with his family from Brazil 3 weeks ago. His immunization status is unknown. The patient appears weak and lethargic. His temperature is 38 C (100.4 F), pulse is 100/min, and blood pressure is 100/60 mm Hg. Examination shows postauricular and suboccipital lymphadenopathy. There is a non-confluent, pink maculopapular rash over the torso and extremities. His palms and soles show no abnormalities. The patient has measles.

**Mistral-7B zero-shot CoT**

The clinical note is well-written and contains all the necessary information regarding the patient's history, presentation, and examination. Therefore, the note does not contain an error.

**Meerkat-7B zero-shot CoT**

To determine if there is an error in the clinical note, let's analyze the information provided step-by-step:

. . .

In the context of an immigrant child with an unknown immunization status, it would be prudent to confirm the diagnosis with laboratory tests such as serology or PCR to detect the presence of the measles virus.
Therefore, the error sentence "The patient has measles." should be corrected to "The patient is suspected to have measles based on clinical presentation. Laboratory tests should be conducted to confirm the diagnosis."

**Meerkat-7B fine-tuned with CoT**

. . .

Given that the rash is non-confluent and pink, and considering the patient's recent immigration from a region where other viral exanthems, such as rubella, might be more common, the clinical presentation is not entirely typical of measles. Rubella, also known as German measles, typically presents with a non-confluent rash that starts on the face and spreads to the rest of the body, which is consistent with the patient's rash description. Additionally, the presence of postauricular and suboccipital lymphadenopathy can be seen in rubella.Considering the symptoms, the progression of the rash, the patient's immunization status, and the geographic origin of the patient, the diagnosis of measles is not as strongly supported as the diagnosis of rubella.
Therefore, the error sentence "The patient has measles." should be corrected to "The patient has rubella.".
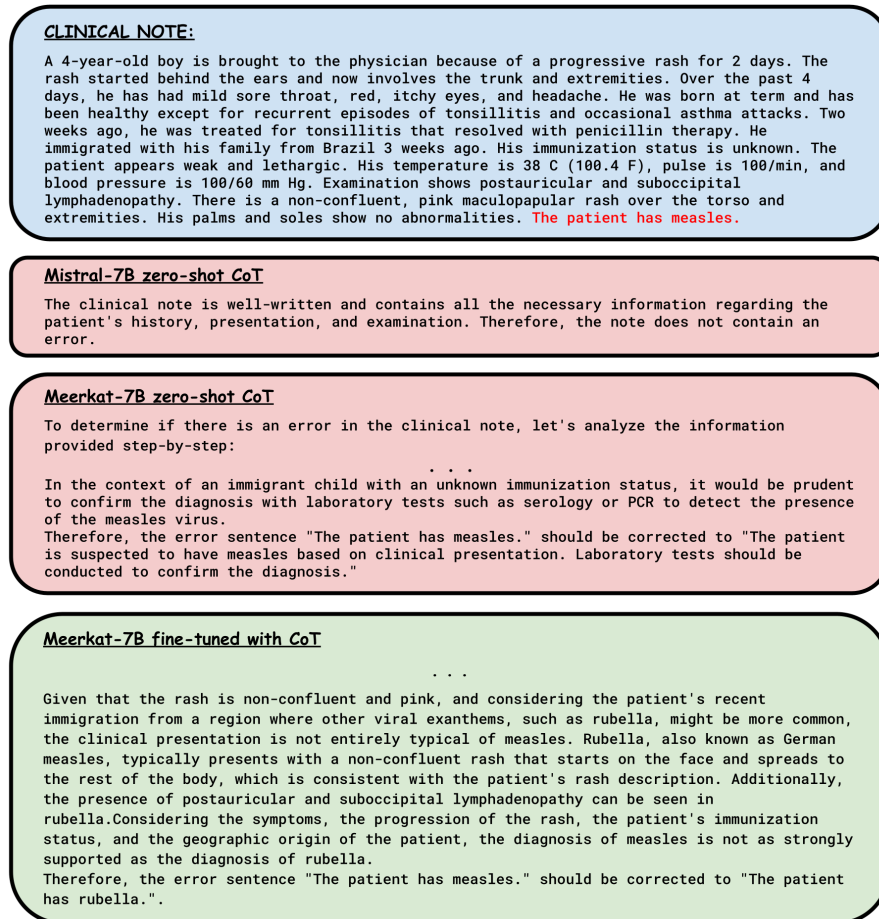
Figure 4: Example outputs from three distinct approaches to clinical note correction. The clinical note contains an error sentence highlighted in red text. The approaches include zero-shot CoT from both Mistral-7B and Meerkat-7B, and fine-tuned Meerkat-7B with CoT reasoning. A green rounded rectangle indicates an accurate correction of the error, while a red rounded rectangle signifies an incorrect response. Meerkat-7B zero-shot CoT detected the error sentence accurately but failed to correct the error sentence due to not fully adapting to the correction task.

Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans. Comput. Healthcare*, 3(1).

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

Hyunjae Kim, Hyeon Hwang, Jiwoo Lee, Sihyeon Park, Dain Kim, Taewhoo Lee, Chanwoong Yoon, Jiwoong Sohn, Donghee Choi, and Jaewoo Kang. 2024. Small language models learn enhanced reasoning skills from medical textbooks. *Preprint*, arXiv:2404.00376.

Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. 2023. Performance

of chatgpt on usmle: potential for ai-assisted medical education using large language models. *PLoS digital health*, 2(2):e0000198.

Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. *Preprint*, arXiv:2402.10373.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, Wenhu Chen, and Xifeng Yan. 2022. Explanations from large language models make small reasoners better. *Preprint*, arXiv:2210.06726.

Xiuquan Li and Tao Zhang. 2017. An exploration on artificial intelligence application: From security, privacy and ethic perspective. In *2017 IEEE 2nd International Conference on Cloud Computing and Big Data Analysis (ICCCBDA)*, pages 416–420.

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.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 23(6):bbac409.

Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Severyn. 2023. Teaching small language models to reason. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1773–1781, Toronto, Canada. Association for Computational Linguistics.

Bertalan Meskó and Eric J Topol. 2023. The imperative for regulatory oversight of large language models (or generative ai) in healthcare. *NPJ digital medicine*, 6(1):120.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023a. Capabilities of gpt-4 on medical challenge problems. *Preprint*, arXiv:2303.13375.

Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoifung Poon, Tao Qin, Naoto Usuyama, Chris White, and Eric Horvitz. 2023b. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *Preprint*, arXiv:2311.16452.

Yasumasa Onoe, Michael J. Q. Zhang, Eunsol Choi, and Greg Durrett. 2021. CREAK: A dataset for commonsense reasoning over entity knowledge. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.

Long N. Phan, James T. Anibal, Hieu Tran, Shaurya Chanana, Erol Bahadroglu, Alec Peltekian, and Grégoire Altan-Bonnet. 2021. Scifive: a text-to-text transformer model for biomedical literature. *Preprint*, arXiv:2106.03598.

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.

Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2023. Distilling reasoning capabilities into smaller language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7059–7073, Toronto, Canada. 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, et al. 2023a. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.

Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. 2023b. Towards expert-level medical question answering with large language models. *Preprint*, arXiv:2305.09617.

Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930–1940.

Shubo Tian, Qiao Jin, Lana Yeganova, Po-Ting Lai, Qingqing Zhu, Xiuying Chen, Yifan Yang, Qingyu Chen, Won Kim, Donald C Comeau, Rezarta Islamaj, Aadit Kapoor, Xin Gao, and Zhiyong Lu. 2024. Opportunities and challenges for ChatGPT and large language models in biomedicine and health. *Briefings in Bioinformatics*, 25(1):bbad493.

Cunxiang Wang, Shuailong Liang, Yili Jin, Yilong Wang, Xiaodan Zhu, and Yue Zhang. 2020. SemEval-2020 task 4: Commonsense validation and explanation. In *Proceedings of the Fourteenth Workshop*

*on Semantic Evaluation*, pages 307–321, Barcelona (online). International Committee for Computational Linguistics.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.

Hongyi Yuan, Zheng Yuan, Ruyi Gan, Jiaxing Zhang, Yutao Xie, and Sheng Yu. 2022. BioBART: Pretraining and evaluation of a biomedical generative language model. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 97–109, Dublin, Ireland. 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.

# CLD-MEC at MEDIQA- CORR 2024 Task: GPT-4 Multi-Stage Clinical Chain of Thought Prompting for Medical Errors Detection and Correction

**Renad M. Alzghoul[1], Abdulrahman Tabaza[1], Aya Abdelhaq[1], Ahmad Altamimi[1]**
[1]Princess Sumaya University for Technology, Amman, Jordan
{ren20228156, abd20200209, aya20228163, a.altamimi}@std.psut.edu.jo

## Abstract

This paper demonstrates CLD-MEC team submission to the MEDIQA-CORR 2024 shared task for identifying and correcting medical errors from clinical notes. We developed a framework to track two main types of medical errors: diagnostics and medical management-related errors. The tracking framework is implied utilizing a GPT-4 multi-stage prompting-based pipeline that ends with the three downstream tasks: classification of medical error existence (Task 1), identification of error location (Task 2), and correction error (Task 3). Throughout the pipeline, we employed clinical Chain of Thought (CoT) and Chain-of-Verification (CoVe) techniques to mitigate the hallucination and enforce the clinical context learning. The model performance is acceptable, given it is based on zero-shot learning. In addition, we developed a RAG system injected with clinical practice guidelines as an external knowledge datastore. Our RAG is based on the Bio_ClinicalBERT as a vector embedding model. However, our RAG system failed to get the desired results. We proposed recommendations to be investigated in future research work to overcome the limitations of our approach.

## 1 Introduction

Medical errors identification and handling in clinical practice is paramount to ensuring optimized patient safety and efficient healthcare delivery. Yet, tracking medical errors is complex to achieve. The umbrella of medical error instances is wide, along with the different phases of the patient journey from assessment to diagnosis, followed by medical management. In addition to the various medical error types that occur throughout these phases. The most common types are diagnostics errors and clinical management errors. Understanding the patient journey-related phases and types of medical errors is crucial in modeling a tracking system framework. Clinical notes of Electronic Health Records (EHR) are considered as documented references of the entire patient's medical journey, from the first point of care to post-medical care plan follow-up. Developing technologies that work on processing clinical notes and notifying healthcare providers with real-time medical error signaling and correction will move healthcare to its next level with a new leveraged paradigm in patient safety. MEDIQA-CORR 2024 shared tasks covers three tasks related to medical error detection and correction from clinical notes (Ben Abacha et al., 2024a). This competition is to establish state-of-the-art techniques (SOTA) to formulate a reliable clinical task of this kind. In this paper, we demonstrate our participation in addressing these three challenges. The complexity of designing a tool that addresses all incidents of medical errors from clinical notes comes from the variety of clinical notes' architecture/context. This variation occurs across:

- Different patient care stages (assessment, diagnosis, medical management plan, follow-up) The type of medical error incidents varies based on the phase of a patient's journey. For instance, diagnostic errors are most likely during the assessment or the diagnosis stage. While at the medical management and plan phase, clinical management errors are the most common.

- The level of documented details related to the history of Present Illness (HPI), Past Medical History (PMH), medication history, clinical findings, diagnosis, medical management plans, and follow-up.

Considering the stage of patient care and the level of clinical note details when mapping case scenarios of medical error incidents within the solution framework is a functional step in building a sustainable medical error tracking system. For instance,

detailed documentation of clinical findings and diagnostic tests may facilitate the identification of diagnostic errors during the assessment or diagnosis stage. On the other hand, comprehensive documentation of treatment plans and medication history may aid in identifying medical management errors during the medical management and planning phase.

Our approach to handling this challenge involves designing a framework that consistently addresses the most common two types of medical errors: diagnostics and medical management errors. To detect these types, we propose a tracking algorithm based on classifying the context of clinical notes to map them with their related medical error case scenarios, ending with detecting medical error incidents of one of these two types. To formulate this in the framework, each clinical note should be screened for three case scenarios of medical error incidents. We categorize clinical notes into two levels of contextual architecture. Level 1 (L1) addresses the first case scenario of medical error. While level two (L2) helps us track the second and third case scenarios of medical error instances. The details of these case scenarios and clinical note levels are demonstrated in section 5.

We implied our tacking approach in a GPT-4 multi-stage prompting-based pipeline that ends with the three downstream tasks: classification of medical error existence (Task 1), identification of error location (Task 2), and correction of error (Task 3). The pipeline is composed of four main stages, as illustrated in Figure 1. Throughout the pipeline, we applied clinical Chain of Thought (CoT) and Chain-of-Verification (CoVe) to mitigate the hallucination and enforce the model to reference its reasoning rational response according to its Evidence-Based Medicine (EBM) clinical practice guidelines attributed knowledge of GPT-4 acquired during training.

As a side work, we developed a RAG system injected with clinical practice guidelines as an external knowledge datastore.

## 2 Background and Related Work

Large language models (LLMs) have proven their potential in various domains, including finance, marketing, and education. Healthcare is a wide area with many horizons (medical education, translational medicine, clinical practice, domain-specific clinical specialty), and the efficiency of LLMs

varies within each horizon area. In some instances, pretrained language models (PLMs) show an efficient performance on specific basic NLP clinical tasks such as Named Entity Recognition (NER), classification, and relationship extraction (RE). However, efficient performance is yet to be reliable and implemented on generative advanced clinical NLP tasks, including clinical text generation, medical question answering, and clinical text summarization. Thus, there is a considerable area for optimizing and leveraging the state-of-the-art in the area of applied generative clinical NLP.

Hallucination and out-of-source generation are some of the main limitations that LLMS and PLMs face, especially with up-to-date and niche-focused domain-related tasks.

RAG and CoT are leading techniques that have been shown to mitigate the limitations mentioned above (Towhidul Islam Tonmoy et al., 2024). RAG framework works on optimizing the output of LLMs by appending LLMs with an external up-to-date knowledge/ data store to be attributed/injected in the generative process through a retrieval and query process (Shuster et al., 2021). RAG contextualizes the model to be more aligned with domain-specific downstream tasks, ending with a more accurate, customized, and specific evidence-grounded response with its data source to be more valid. It encompasses 3 main components: the retrieval, the generation, and the augmentation techniques. Pretrained LLMs performance is comparable with LLMs with RAG. RAG can overcome the need to retrain/finetune LLMs on up-to-date or domain-specific information. Instead, it augments the knowledge with LLMs without retraining the model and results with applicable performance (Gupta et al., 2024).

RAG with LLMS has shown its potential to drastically advance LLMs' usability and reliability. In the healthcare domain, integrating RAG with LLMs has been applied with notable enhancements in the generated responses of LLMs to make them more accurate, informative, and reliable. LLMs output aligned remarkably with the augmented RAG case-specific medical knowledge. (Zakka et al., 2024; Ge et al., 2023) incorporated RAG into LLMs with a medical knowledge database for medical guidelines and treatment recommendations. These LLMs with RAG outperform standard LLMs significantly on the level of accuracy, user satisfaction and consistency. Another study illuminates the impact of appending clinical trials related

to medical knowledge to LLM with RAG on an exceeding performance of this framework compared to experts in clinical trial screening (Unlu et al., 2024).

From the perspective of our shared task, medical error correction is one of its downstream tasks. Which needs for techniques to support and enhance formulating this task. Factual Error Detection and Correction with Evidence Retrieved (FLEEK) (Bayat et al., 2023), is an innovative solution that overcomes hallucinations. It performs two tasks: fact verification and fact revision. It splits an input passage into sentences and uses a sequential pipeline to verify each sentence and correct it so it reduces hallucinations with unstructured knowledge, such as web-based and structured knowledge graphs. Facts are defined as units of information that describe entities, relations, or events and are represented using a semi-structured triple format. FLEEK's performance is evaluated using benchmarks and preliminary experiments using manually created evaluation data. (Dhuliawala et al., 2023) introduced Chain-of-Verification (CoVe), a method to reduce hallucinations in large language models by breaking down verifications into more straightforward questions and self-correcting them. Factored CoVe helps alleviate copying hallucinations and provides performance gains over original responses.

For accurate clinical diagnosis, (Savage et al., 2024) explored LLMs in medicine to imitate the Clinical Reasoning Rationale (CRR) as a COT approach to perform differential diagnosis steps during the medical diagnosis process. They created a diagnostic prompting method that allows LLMs to construct diagnosis while accurately mimicking clinical reasoning using CoT prompts. This led to GPT-4 being prompted to imitate the thought processes of clinicians, giving doctors a comprehensible justification for assessing the precision of LLM replies. These techniques are utilised in order to enhance our model's ability to identify medical errors in the clinical context and consider potential corrections.

RRED (Min et al., 2022) is a deep learning framework designed to detect errors in radiology reports. The system creates artificial existing errors using an error generator and supervised learning techniques.

The method addresses error detection in radiology reports using a deep learning framework with a rich contextual and medical understanding. The error generator generates realistic errors from existing radiology reports, creating synthesized datasets for training the error detector. The error detector employs a BERT-based architecture to detect errors based on a semantic understanding of radiology reports.

## 3 Dataset

### 3.1 MEDIQA-CORR 2024

MEDIQA-CORR 2024 proposed three shared tasks related to medical error detection and correction from clinical notes. Table 1 illustrates the characteristics of the data concerning both the input and output parameters for each task.

Three datasets were provided in this challenge: training, validation, and testing (Ben Abacha et al., 2024b).

#### 3.1.1 Training Dataset

The initial dataset is derived from the University of Washington (MS) Training Set. It comprises 2,189 clinical texts, all of which either have one error or none at all (1 denoting that the text has an error and 0 denoting that there are no errors). This data set includes the original clinical note, the error sentence, the corrected sentence, and the corrected text as a whole. These clinical notes document the patient's related conditions throughout different patient care phases.

#### 3.1.2 Validation Dataset

Two validation datasets: MS validation set contains 574 clinical texts, and the University of Washington (UW) validation set includes 160 clinical texts. Clinical notes are unlabeled in this dataset for validation purposes. However, The labeled notes of the dataset were accessible.

#### 3.1.3 Testing Dataset

This dataset was provided by MEDIQA-CORR 2024 and contains 574 clinical texts, which include only the sentences without flagging the error or correcting it. It will serve as a means of testing our model to determine its performance for tasks 1, 2, and 3.

### 3.2 Clinical Guidelines Dataset

For clinical knowledge enhancement, we utilized the Clinical Guidelines corpus dataset (Chen et al., 2023), comprising 47,000 clinical practice guidelines sourced from 17 reputable online medical references. We utilized this dataset as a data store to

Table 1: Example of the input and the structured output format for tasks 1, 2, and 3.

| Input | 0 A 9-year-old girl is brought to the pediatrician by her mother who reports that the girl has been complaining of genital itching over the past few days.<br>1 She states she has noticed her daughter scratching her buttocks and anus for the past week; however, now she is scratching her groin quite profusely as well.<br>2 The mother notices that symptoms seem to be worse at night.<br>3 The girl is otherwise healthy, is up to date on her vaccinations, and feels well.<br>4 She was recently treated with amoxicillin for a middle ear infection.<br>5 The child also had a recent bought of diarrhea that was profuse and watery that seems to be improving.<br>6 Her temperature is 98.5 F (36.9 C), blood pressure is 111/70<br>7 mmHg, pulse is<br>8 83/min, respirations are 16/min, and oxygen saturation is 98% on room air.<br>9 Physical exam is notable for excoriations over the girl's anus and near her vagina.<br>10 Suspected of infection with Giardia lamblia. |
|---|---|

| Task | Output | |
|---|---|---|
| **1: Error Flag** | 1 | |
| **2: Error Sentence** | 10 Suspected of infection with Giardia lamblia. | |
| **3: Corrected Sentence, Corrected Text** | **Corrected Sentence:** Suspected of infection with Enterobius vermicularis. | **Corrected Text:** A 9-year-old girl is brought to the pediatrician by her mother, who reports that the girl has been complaining of genital itching over the past few days.<br>....<br>Suspected of infection with Enterobius vermicularis. |

be augmented with LLM by the RAG framework to enhance the clinical practice domain of knowledge.

## 4 Description of the Three Tasks

In this research, we worked on three tasks; the description of each task is as follows:

**Binary Classification (Detecting Medical Errors)**: In this task, we had to determine whether the text contained a medical error or not. This task involved binary classification (0/1) regarding the existence of the medical error in the text.

**Span Identification (Locating Errors within Text)**: In the second task, if there is a medical error in the given clinical text, the model should identify the precise text span linked to it. The exact location or the section where the error is found in the medical text.

**Natural Language Generation (Correction of Errors)**: In the last task, the model must provide a free text correction where the medical error is detected in the clinical text. This task aims to improve the quality and accuracy of mistake resolution in clinical situations by producing a human-like context to augment the automated correction process.

## 5 Methods

We build a GPT-4 prompting-based pipeline that processes the text of clinical notes and endeavors to detect medical errors and correct them if they exist. The model pipeline addresses medical error incidents using three types of clinical note context case scenarios. These case scenarios are classi-
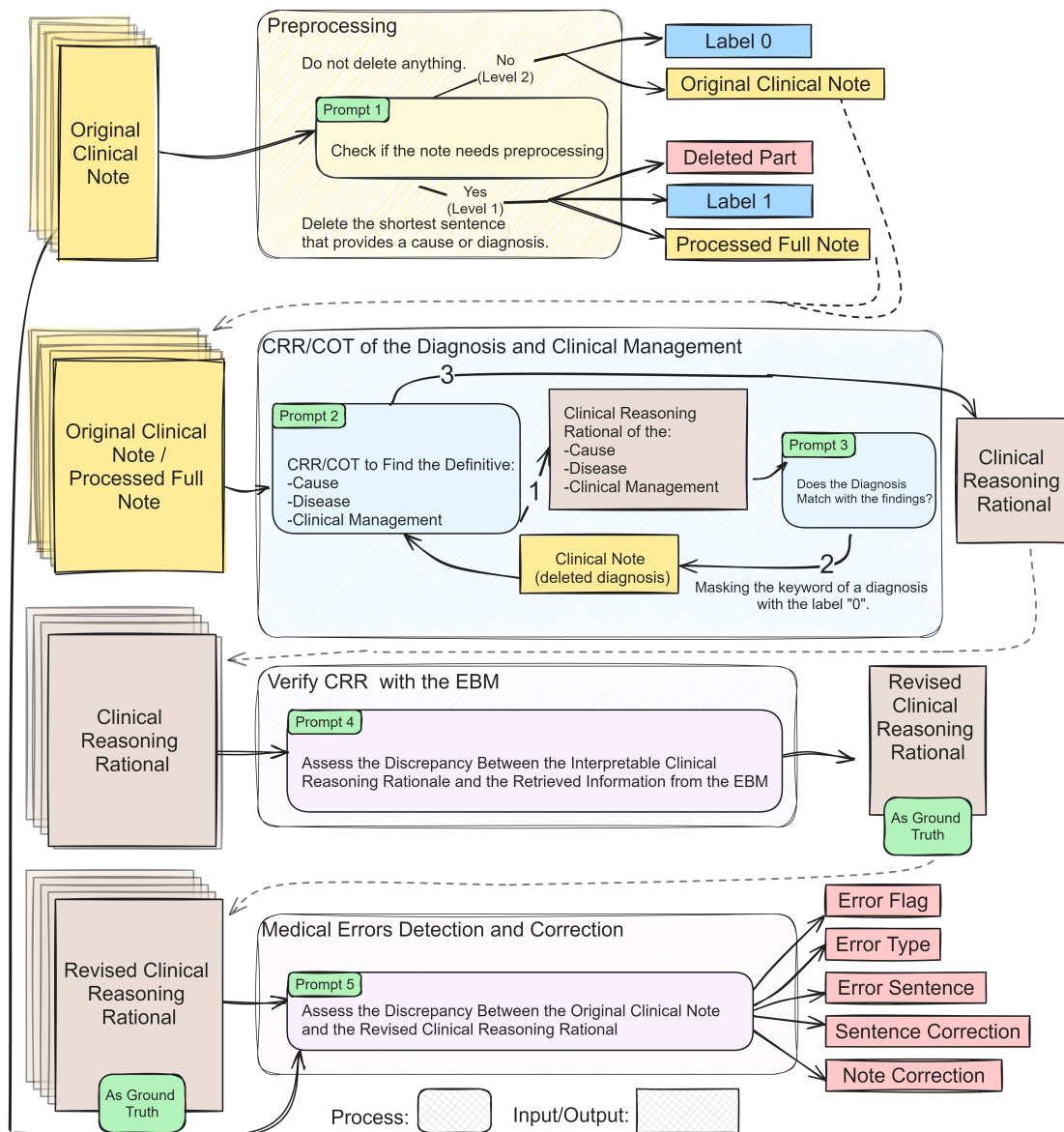
Figure 1: GPT-4 Multi-stage prompting pipeline. The pipeline comprises four phases, each with certain prompts and processes. The data flow process throughout the pipeline with detailed inputs and outputs is depicted in the figure. CRR: Clinical Reasoning Rationale; COT: Chain of Thought; EBM: Evidence-Based Medicine.

fied according to the types of medical errors we aim to track, including diagnostic and clinical management errors. **Diagnostic errors** occur when a patient's medical condition is attributed to an incorrect cause (pathogen, poison, etc.) or misdiagnosis. While **medical management types of errors** involve an incorrect medication, lab test, procedure, or medical image for a patient's medical condition. These errors occur due to an incorrect interpretation of findings to the correct cause and diagnosis or an incorrect interpretation of the cause and diagnosis to the most appropriate medical management. These consequences of incorrect interpretations that lead to medical errors are addressed

through our solving approach, which is systematically employed throughout the pipeline to track incidents of these types of medical errors.

In a brief overview, the pipeline involves removing the cause or diagnosis from the clinical note and subsequently generating a CoT process. This process is tailored to find the most probable cause, findings, and clinical management of a patient's clinical case. Then, assess for discrepancy between the generated CoT and original clinical note, thus indicating if an error exists and correcting it. The output of each prompt is formulated to be unified across all dataset. The pipeline is composed of four main processing phases. In the following subsec-

tions 5.1, 5.2, 5.3, 5.4, we elucidate these phases with the details of the techniques used within each phase.

Clinical notes can be articulated in various architectures, depending on the timeline and stages of the patient's medical journey. We categorized clinical notes into two levels of contextual architecture. Level one (L1) addresses the first case scenario of medical error, while level two (L2) helps us track the second and third case scenarios of medical error instances. This facilitates formulating a clear link between the context of clinical note and the most suspected medial error case scenario to be signaled as the following:

**L1: Documented HPI or PMH without clinical management actions related to the cause or the diagnosis.** At this level, the clinical management would be directly moved to detect the suspected cause and diagnosis of the medical condition upon the existing clinical findings without thinking of any further needed medical management actions to confirm the diagnosis or to manage the current patient's medical status.

**L2: Documented HPI or PMH with clinical management actions related to the cause or the diagnosis**. At this level, the patient is already known with the diagnosis and related cause. A wider margin of settings. This level will track the second and third case scenarios of medical error instances. Medical errors at this level mostly will be related to signal the second case scenario. The third case scenario is related to the previous step of the patient medical journey, at the diagnosis step. Specifically when the diagnosis and clinical findings are not directly connected to each other in most common clinical contexts, yet the note contains a clinical management action related to the incorrect patient's diagnosis. Refer to Figure 2 of Appendix B that demonstrates an example of each case scenario.

## 5.1 Preprocessing

In order to orient the model to the downstream tasks of detecting medical errors and correcting them, the model is firstly promoted to detect the shortest part of a sentence that declares the cause or the diagnosis of a medical condition in the clinical note in order to be masked. The rationale behind this masking is to force the generated CoT (in the next phase) without being biased by the already declared cause or diagnosis in the clinical note. This process is applied to L1 clinical notes to help track the first case scenario of medical error incidents. With L1

clinical notes, the model will return the "Deleted Part" and the "Processed Full Note". While L2 type clinical notes should be returned as they are. As shown in Figure 1. "Processed Full Note" and L2 type of clinical notes are designed to be passed to the next stage. The prompt, "Prompt1" of this stage is shown in Table 3 of Appendix A.

## 5.2 Clinical Reasoning Rationale/CoT of the Diagnosis and Clinical Management

This stage is tailored to return the cause, diagnosis, and medical management of each clinical note returned from the previous step. These returns are based on CRR. CRR is a CoT technique applied in the clinical context. Two prompts are used at this stage. The first one, "Prompt 2", as shown in Figure 1. We engineered the prompt to do step by step deduction to create a differential diagnosis from which to find the most likely cause and diagnosis of medical condition in a clinical note. The answer is constrained to the documented clinical findings of the clinical note, directing the model to be more definitive to the most probable correct cause and diagnosis without expanding the probability of other differential diagnoses based on further clinical investigations actions beyond what the note handles. Subsequently, upon the most likely cause and diagnosis of a medical condition, the model deduces the most correct clinical management (treatment, clinical care plan, intervention, procedure...etc.) using a step-by-step process. "Prompt 2" is demonstrated in Table 3 of Appendix A. The "Clinical Reasoning Rationale" output from "Prompt 2" is then employed to be used as a reference for the next prompt, "Prompt 3". "Prompt 3" serves as a checkpoint of the third case scenario of medical error instances. If the documented diagnosis is based on clinical findings that are not directly related to each other in the common clinical context, it indicates a diagnostic type of medical error that needs to be processed. The model at this prompt is designed to process any discrepancy between the clinical finding and the diagnosis by marking the keyword of the incorrect diagnosis with the label "0". Then, the processed note will be passed again to "Prompt 2" to find the correct cause, diagnosis, and clinical management. The final "Clinical Reasoning Rationale" output at this phase will be passed to the next phase, "Verify CRR with the EBM". "Prompt 3" is demonstrated in Table 3 of Appendix A.

## 5.3 Verify CRR with the EBM

The functionality of this stage is to verify the baseline interpretable "Clinical Reasoning Rationale" output by instructing the model to generate questions that target and retrieve each information in the CRR note, then correct any discrepancy, following (Bayat et al., 2023; Dhuliawala et al., 2023). The prompt of this phase, "Prompt 4", demonstrated in Table 4 of Appendix A, helps in forcing the model to reference EBM clinical practice guidelines attributed knowledge of GPT-4 acquired during training. The output, verified (CRR), will be taken as ground truth knowledge for the next stage, "Medical Errors Detection and Correction".

## 5.4 Medical Errors Detection and Correction

This final stage, "Prompt 5", as shown in Table 5 of Appendix A, is designed to be the cut-point step for the three tasks. The verified CRR is taken as ground truth knowledge for clinical notes. It should include the correct cause, diagnosis, and clinical management for the note. The model is instructed to compare the verified CRR with the original clinical note for cause, diagnosis, and clinical management discrepancies, as shown in Figure 1. Clinical notes with the contextual architecture of L1, discrepancies related to the cause or diagnosis should be cached. For L2, discrepancies related to the cause or diagnosis should be cached as well, along with clinical management discrepancies.

Clinical management is a wide aspect, including interventions related to treatment, ordering certain lab tests and images, transfer, and procedure. The CRR includes all the necessary clinical management actions related to the clinical note case, while the original note might include one of them. This case scenario might drive the model to detect it as a discrepancy, correcting it with the appropriate completed clinical management plan. Additionally, if a diagnostic error exists, it should be corrected, ending with two medical errors identified. For L2, since the clinical note contains only one error. The model is instructed to prioritize correction for diagnostics errors (cause and disease). Then, to clinical management-related errors. From this phase, we should have the "Error Flag" for task 1, the "Error Sentence" for task 2, the "Sentence Correction" for task 3, and the "Note Correction" as a full note.

## 5.5 RAG

To enhance the accuracy and relevance of the generated clinical response while throughout the processes related to CRR/CoT and CoVe, our approach conducted an experiment using a RAG framework. Our RAG system is based on the parameters outlined in Table 2.

Clinical Guidelines corpus dataset is utilized as an external knowledge database. Our RAG system should integrate this knowledge into the prompt output through the query and retrieval process. The process of generating the query is based on the instructions stated in "Prompt 4". Where the RAG system is utilized as the ground truth for this stage.

## 6 Experiments and Results

Here, we report the experimental findings demonstrating our model's effectiveness on the shared tasks. Our model was performed on the three tasks utilizing a zero-shot learning approach and a GPT-4 prompting-based pipeline with CoT and CoVe methods and structured output.

Since our approach mainly focuses on zero-shot learning. Thus, we only used the training and validation for prompt optimization until we reached a reasonable output in tracking the three case scenarios. For a comprehensive show of our approach's functionality in tracking the three case scenarios, please refer to Appendix B. It provides an example experiment of tracking each one of the case scenarios, illustrating its input/output at each prompt through the entire process.

The Function Calling feature of OpenAI API and JSON mode is utilized to get the aimed structured format and ensure consistent output throughout the dataset.

The results of our approach performance on the testing dataset show that the accuracy of the first and second tasks is 0.566 and 0.49, respectively, without using external knowledge sources, fine-tuning methods, or group learning. While for task 3 and the main results, performance metrics yielded the following scores: ROUGE-1-F of 0.427, BERTScore of 0.48, BLEURT of 0.53, their Aggregate-Score (Mean of ROUGE-1-F, BERTScore, BLEURT-20) of 0.48, and their Composite Scores of 0.34. These metrics assess the model's ability to produce contextually appropriate corrections for clinical errors identified in clinical text.

Table 2: RAG parameters.

| Parameter | Value |
|---|---|
| Chunk_Size | 500 |
| Chunk_Overlap | 32 |
| Embedding Model | Bio_ClinicalBERT |
| Embedding Dimension | 768 |
| Model Pooling Strategy | Mean |
| Vector Index | Faiss Hierarchical Navigable Small Worlds Index with Neighboring Vectors of 32 |
| Chunking Strategy | Recursive Character Text Splitter from langchain |

## 7 Discussion and Future Scope

The performance of our model is somehow acceptable but unreliable when applied to the testing dataset. The algorithmic approach of following three case scenarios of medical events based on clinical note contextual architecture might undertrack other medical error incidents case scenarios. For the LLM we have used, GPT-4, we have aimed to be built based on a RAG framework incorporated with clinical practice guidelines. Our hypothesis was to optimize the output of un pretrained generative model in the arena of clinical practice (niche-focused) to get reliable, inferential and ground truth knowledge without hallucinations. RAG framework is the best to be employed with a massive LLM such as GPT, BART, or T5. For limited hardware resources, we took GPT-4, as an open LLM model. Our RAG system failed in retrieving relevant queries. It was supposed to be connected to our pipeline at phase 4, "Verify CRR-CoT with the EBM", but for irrelevant retrieved chunks, we continued the work without it. This shortcoming performance could be one of the following:

- The vector embedding model we used, Bio_ClinicalBER, is not one of the vector embedding models that are already designed for the RAG frameworks (trained with a retrieval objective). The choice of Bio_ClinicalBERT was to test a clinical embedding model rather than general used ones. In addition to a limited time, we could not test the SOTA models with our RAG such as ColBERT (Khattab and Zaharia, 2020).

- The chunking strategy we used is a naive technique, which might be the cause of the poor informative chunks.

- The complicated structure of the utilized Clinical Guidelines dataset as an external datastore for our RAG system.

Along the pipeline, we used CRR-CoT and CoVe as prompt optimization techniques. Future work should investigate the performance of public LLMs, to unlock their known capabilities for these downstream tasks within hardware accessible facilities. In addition to exploring advanced chunking strategies such as semantic chunking and finetuning a domain-specific model such as Bio_ClinicalBERT for retrieval.

## 8 Limitations

Our work is limited by two points. The first one is our inability to produce a reliable RAG system due to time constraints. We could not explore how our approach would perform with a successful RAG system implementation; specifically, we utilized rich, niche-focused external knowledge to boost the reliability and applicability of the generated output. Secondly, with our limited computing and financial capacity, we would not be able to experiment with other massive LLMS, whether they are general, clinically fine-tuned, or pre-trained ones.

## Acknowledgments

## References

Farima Fatahi Bayat, Kun Qian, Benjamin Han, Yisi Sang, Anton Belyi, Samira Khorshidi, Fei Wu, Ihab F Ilyas, and Yunyao Li. 2023. Fleek : Factual error detection and correction with evidence retrieved. *arXiv preprint arXiv:2310.17119.*

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview

of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Zeming Chen, Alejandro Hernández, Cano Angelika, Romanou Antoine, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, and Vinitra Swamy. 2023. MEDITRON -70B: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*.

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-Verification reduces hallucination in large language models. *arXiv preprint arXiv:2309.11495*.

Jin Ge, Steve Sun, Joseph Owens, Victor Galvez, Oksana Gologorskaya, Jennifer C Lai, Mark J Pletcher, and Ki Lai. 2023. Development of a liver disease-specific large language model chat interface using retrieval augmented generation. *medRxiv*.

Aman Gupta, Anup Shirgaonkar, Angels de Luis Balaguer, Bruno Silva, Daniel Holstein, Dawei Li, Jennifer Marsman, Leonardo O Nunes, Mahsa Rouzbahman, and Morris Sharp. 2024. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture. *arXiv preprint arXiv:2401.08406*.

Omar Khattab and Matei Zaharia. 2020. Colbert : Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 39–48.

Dabin Min, Kaeun Kim, Jong Hyuk Lee, Yisak Kim, and Chang Min Park. 2022. Rred: A radiology report error detector based on deep learning framework. In *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 41–52.

Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath Rangan, and Jonathan H Chen. 2024. Diagnostic reasoning prompts reveal the potential for large language model interpretability in medicine. *NPJ Digital Medicine*, 7(1):20.

Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3784–3803.

S M Towhidul Islam Tonmoy, S M Mehedi Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv e-prints*, pages arXiv–2401.

Ozan Unlu, Jiyeon Shin, Charlotte J Mailly, Michael F Oates, Michela R Tucci, Matthew Varughese, Kavishwar Wagholikar, Fei Wang, Benjamin M Scirica, and Alexander J Blood. 2024. Retrieval augmented generation enabled generative pre-trained transformer 4 (GPT-4) performance for clinical trial screening. *medRxiv*, pages 2002–2024.

Cyril Zakka, Rohan Shad, Akash Chaurasia, Alex R Dalal, Jennifer L Kim, Michael Moor, Robyn Fong, Curran Phillips, Kevin Alexander, and Euan Ashley. 2024. Almanac—retrieval-augmented language models for clinical medicine. *NEJM AI*, 1(2):AIoa2300068.

## A    The Pipeline's Prompts Templates

The instruction template of each prompt throughout the pipeline is demonstrated in Tables 3, 4, and 5.

## B    The Output of the Pipeline

The model's performance at each stage is demonstrated with an example for each of the three case scenarios. Figure 2 displays an example for each case scenario, including the original note as an input to the pipeline and the corrected sentence. The outputs throughout the pipeline's phases are illustrated as follows: at the Preprocessing B.1, Clinical Reasoning Rationale/CoT of the Diagnosis and Clinical Management B.2, Verify CRR with the EBM B.3, and Medical Errors Detection and Correction B.4 stage.

### B.1    Preprocessing

Figure 3 illustrates the outputs at stage one, the Preprocessing, for the first case scenario medical errors that are in particular related to clinical notes with level one context. Figures 7 and 11, represent the second and third case scenarios, respectively, both of which are in particular related to clinical notes with level two context.

### B.2    Clinical Reasoning Rationale/CoT of the Diagnosis and Clinical Management

The outputs of this phase utilizing Prompts 2 and 3 are depicted in Figures 4, 8, and 12 for the three case scenarios.

### B.3    Verify CRR with the EBM

At stage 3, the verified CRR output of clinical notes corresponding to the first, second, and third case scenarios is shown in Figures 5, 9, and 13, respectively.

Table 3: The templates of Prompts 1,2, and 3.

| Prompt | Instruction Template |
|---|---|
| 1 | I will give you a clinical note. You have to delete the shortest sentence that shows the cause or diagnosis, following to these conditions: <br> 1) If the clinical note mentions any of clinical management actions (treatment, clinical care plan, or any intervention,....etc.) related to ( management of past medical history, management history of present illness, diagnosis), then do not delete anything. Give this label 0. <br> 2) Else, then delete the sentence that shows the cause and diagnosis. Give this label 1 <br> 3) Print the assigned labels 1 or 0. <br> 4) Print the deleted part if applicable. <br> 5) Print the full final note. |
| 2 | 1) Based on Evidence-Based Medicine, use step-by-step deduction to create a differential diagnosis and then use step by step deduction to identify both of the most likely causing (Pathogen name of the bacteria, worm, virus, fungi,....etc., poison,.... etc.) and diagnosis separately. The answer should also be definitive to one cause and one diagnosis without requiring any further clinical investigating action. <br> 2) Then, step by step, deduce the most correct (treatment, clinical care plan, clinical management, intervention) <br> You are designed to output JSON. <br> The JSON should be structured like this: <br> { <br> "Differential Diagnosis Step by Step": { <br> "Step 1": ..., <br> "Step 2": ..., <br> "Step N": ... <br> }, <br> "Differential Diagnosis": { <br> "Most Likely Cause": ..., <br> "Explanation": ... <br> }, <br> "Treatment Step by Step": { <br> "Step 1": ..., <br> "Step 2": ..., <br> "Step N": ... <br> }, <br> "Definitive Diagnosis": ..., <br> "Treatment": { <br> "Definitive Treatment": ... <br> } <br> } |
| 3 | 1) Use this interpretable clinical reasoning rationale you have produced for this clinical note: cot <br> 2) Based on the interpretable clinical reasoning rationale, If the clinical note mentions a diagnosis or a medical condition that is based on a clinical presentation or findings that are not directly connected to each other in most common clinical contexts, then there should be a medical error in the diagnosis. <br> 3) Delete the diagnosis or a medical condition-related keyword from the clinical note. <br> 4) Print the deleted keyword if applicable. <br> 5) Print the full final note, where the deleted keyword should be masked with this label -> "0" <br> You are designed to output JSON. has to be structured like this: <br> {{ <br> "DeletedKeyword": ..., <br> "FullFinalNote": ... <br> }} |

Table 4: The template of Prompt 4.

| Prompt | Instruction Template |
|---|---|
| 4 | You have to verify your interpretable clinical reasoning rationale of the diagnosis you have produced of its related clinical note. The verification should be done by generating questions that target and retrieve information from the most appropriate clinical practice guidelines. <br> -Make the query address the name of the guideline you want to retrieve that response from. <br> -If you want to check for the diagnosis of clinical findings, make the query address the related clinical findings you want to check for the diagnosis. <br> -Make the directed query address the most likely correct (cause, diagnosis). <br> -Make the directed query address the recommendations part of the guideline related to (diagnosis, clinical management, treatment, drug of choice) <br> -Search from the directed guidlines. <br> -Return the information you gained. <br> -Compare your interpretable clinical reasoning rationale with the retrieved information from the guideline; if there is a discrepancy, show it. <br> -If there is a major discrepancy, take the retrieved information as ground truth and print out the final CoT after being revised. <br> You are designed to output JSON. <br> It has to be structured like this: <br> {{ <br> "VerificationQueries": { <br> "Query 1": ..., <br> "Query 2": ..., <br> "Query 3": ..., <br> "Query N": ... <br> }, <br> "RetrievedInformation": { <br> "Response 1": ..., <br> "Response 3": ..., <br> "Response N": ... <br> }, <br> "Comparison": { <br> "Clinical Findings": ..., <br> "Causes": ..., <br> "Treatment": ... }, <br> "Discrepancy": ... (could be nullable), <br> "FinalCoT": { <br> "Differential Diagnosis Process": { <br> "Step 1": ..., <br> "Step 2": ..., <br> "Step 3": ..., <br> "Step N": ... <br> }, <br> "Definitive Cause": { <br> "Most Likely Pathogen/Cause": ... }, <br> "Definitive Diagnosis": ..., <br> "Treatment Plan": { <br> "Step 1": ..., <br> "Step 2": ..., <br> "Step 3": ..., <br> "Step 4": ..., <br> "Step N": ... <br> } <br> } <br> }} |

Table 5: The template of Prompt 5.

| Prompt | Instruction Template |
|---|---|
| 5 | 1) Use this interpretable clinical reasoning rationale you have produced as a ground truth<br>{verified_cot}<br>2) Compare if the clinical note matches the ground truth to tell if the clinical note has a medical error in (diagnosis (pathogen, poison, disease), clinical management (treatment, clinical care plan, intervention (order certain lab test, transfer, certain image by name, procedure).).<br>3) Identify any discrepancy between the ground truth and the clinical note.<br>4)Then, if there is anything in the clinical note related to either diagnosis or cause that is not available (referenced) in the ground truth reference, then label it as a medical error. And skip the steps related to clinical management.<br>5)Then else, if there is anything in the clinical note related to clinical management after diagnosis that is not available (referenced) in the ground truth reference, specifically in (clinical management-related sections), then label it as a medical error. And skip the steps related to the diagnosis or cause. If there is a medical error, identify its type (diagnosis, cause, or clinical management) and print it, identify the exact related shortest part and print it, and correct it with the shortest possible correction. Do not change the format of the corrected part. Only correct the related keyword. Then, if the error type is related to clinical management-related errors, the corrected sentence should be definitive to the exact needed medication, procedure, image,..... etc., not general. Not as a recommendation. Correct the note directly with the most correct probable needed audit. If the error type is related to diagnosis, cause, or clinical management, consider this error correction to be edited on the final corrected note. The priority to add the correction of diagnosis and cause first to be considered. Consider one correction only, depending on the context. Finally print out the corrected final note.<br>The clinical note you have to correct is split into sentences with an index for each.<br>The correction you return includes the error flag, the error location, and the sentence correction. |

## B.4 Medical Errors Detection and Correction

At stage 4, using Prompt 5. The final outputs include the "error flag" to identify if an error exists, "error location", and "sentence correction" if there is an error within a clinical note. See Figures 6, 10, and 14 corresponding to the first, second, and third case scenarios, respectively.
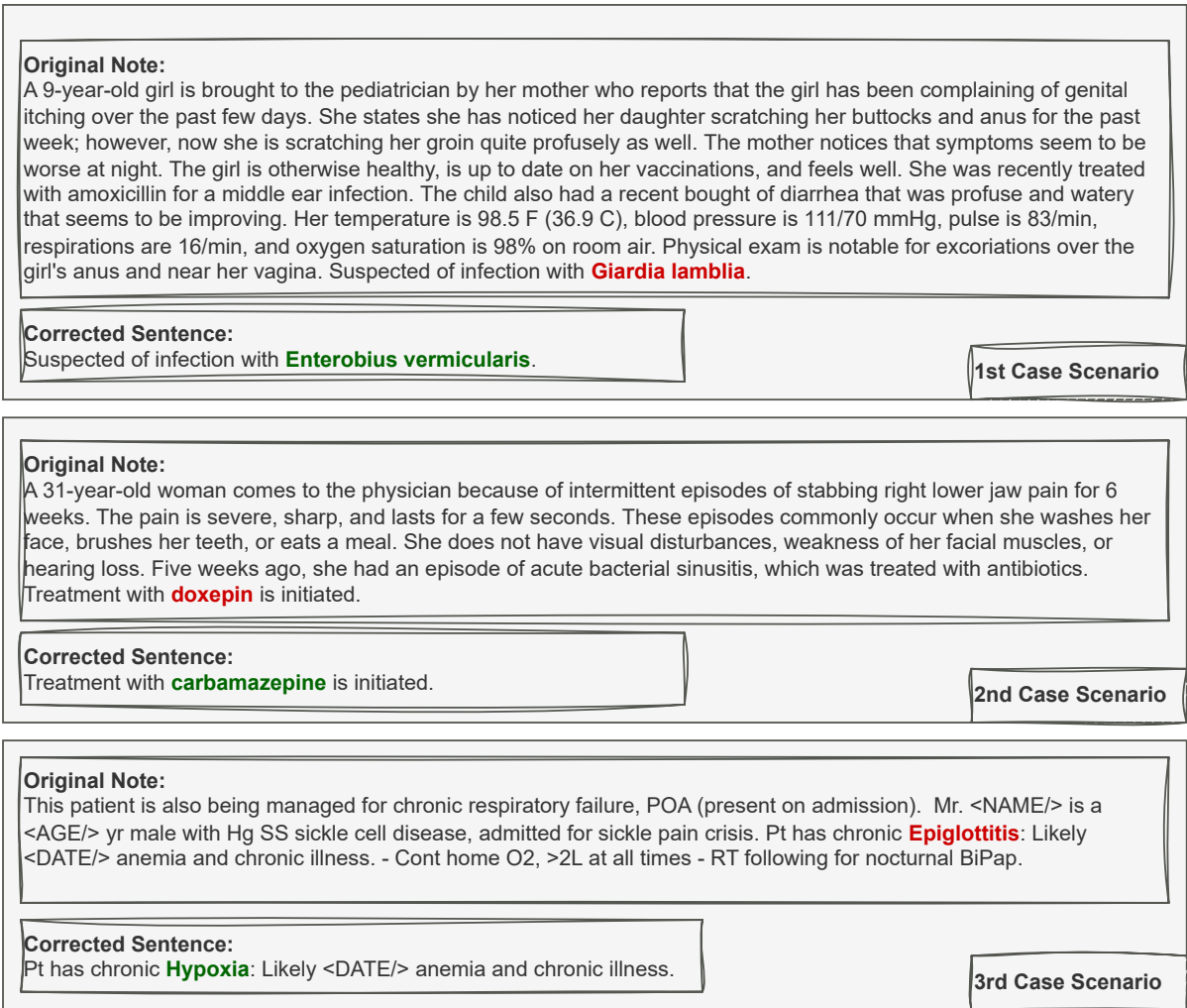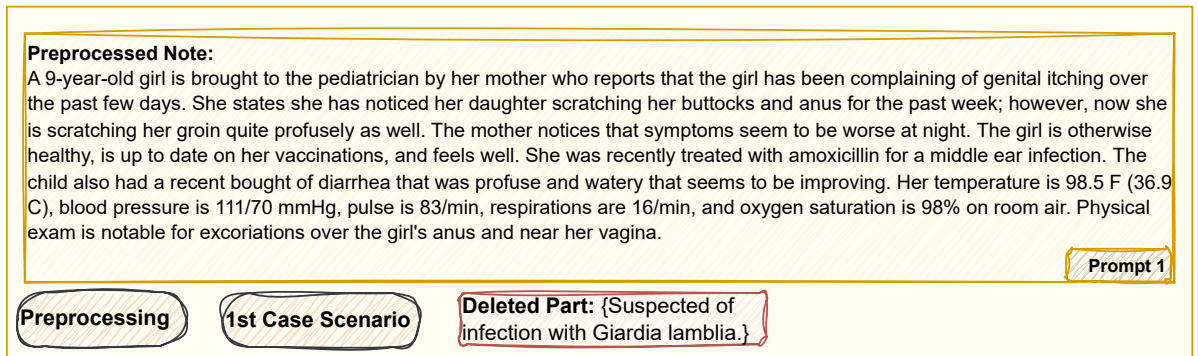
Figure 2: Example for each type of case scenario.



Figure 3: The outputs of clinical notes with the first case scenario, level one notes. At stage 1, the Preprocessing, using Prompt 1.

**CRR Round 1:**
{'**Differential Diagnosis Step by Step**': {'Step 1': 'Begin by considering common causes of genital itching in children, including irritant contact dermatitis, infections (bacterial, viral, fungal, parasitic), and infestations.',
'Step 2': 'Note the symptom of worsening itching at night, which is a characteristic feature of certain parasitic infections.',
'Step 3': 'Acknowledge the recent use of antibiotics (amoxicillin) and a bout of diarrhea, which could predispose to opportunistic infections or alter normal flora, but these are less likely to cause localized itching primarily at night.',
'Step 4': 'Consider the physical exam findings of excoriations over the genital and anal area, which suggest an infestation rather than a generalized dermatitis or systemic illness.',
'Step 5': 'Given the symptom of itching worsening at night and the specific locations of the excoriations, focus on parasitic infestations that exhibit nocturnal activity.'},
'**Differential Diagnosis**': {'Most Likely Cause': 'Enterobius vermicularis (Pinworm) Infestation',
'Explanation': 'The symptoms of nocturnal itching, particularly around the anus, and the presence of excoriations in a child, are classic for pinworm infestation. Pinworms are a common cause of genital itching in children and are known for increased activity and egg-laying at night, which markedly exacerbates symptoms during these hours.'},
'**Treatment Step by Step**': {'Step 1': 'Confirm the diagnosis if possible, typically by using the tape test to collect eggs for microscopic identification. However, given the classic presentation, empirical treatment may be initiated without confirmation.',
'Step 2': 'Prescribe a single dose of albendazole or mebendazole, with a second dose to be taken in two weeks to eradicate any subsequent hatchlings.',
'Step 3': 'Educate the family on the importance of strict hygiene measures to avoid reinfection. This includes frequent handwashing, especially after using the bathroom and before eating, daily showering, changing underwear and bed linens frequently, and avoiding nail-biting.',
'Step 4': 'Advise treatment of all family members living in the household to prevent reinfection or spread, as pinworms are highly contagious.'},
'**Definitive Diagnosis**': 'Enterobiasis (Pinworm Infestation)',
'**Treatment**': {'Definitive Treatment': 'Albendazole or mebendazole, with a repeat dose in 2 weeks, coupled with stringent personal and household hygiene measures.'}}

Prompt 2

**Clinical Note (deleted diagnosis with 0):** A 9-year-old girl is brought to the pediatrician by her mother who reports that the girl has been complaining of genital itching over the past few days. She states she has noticed her daughter scratching her buttocks and anus for the past week; however, now she is scratching her groin quite profusely as well. The mother notices that symptoms seem to be worse at night. The girl is otherwise healthy, is up to date on her vaccinations, and feels well. She was recently treated with amoxicillin for a middle ear infection. The child also had a recent bought of diarrhea that was profuse and watery that seems to be improving. Her temperature is 98.5 F (36.9 C), blood pressure is 111/70 mmHg, pulse is 83/min, respirations are 16/min, and oxygen saturation is 98% on room air. Physical exam is notable for excoriations over the girl's anus and near her vagina. Suspected of infection with **0**.

Prompt 3

**CRR Round 2:**
{'**Differential Diagnosis Step by Step**': {'Step 1': 'Given the symptoms of genital itching, particularly noted to be worse at night, and the presence of excoriations over the anus and near the vagina, parasitic infections are highly suspected.',
'Step 2': 'The recent history of diarrhea could suggest a gastrointestinal infection or infestation, which may be related to the itching if the pathogen has a cutaneous manifestation or provokes such a reaction.',
'Step 3': 'The recent use of amoxicillin for an ear infection introduces the possibility of a secondary yeast infection due to antibiotic use; however, the description of symptoms being worse at night and the specific location of itching are less typical for a yeast infection.',
'Step 4': "Considering the girl's age and the specific symptoms described, common infections of the perianal region in children should be considered, such as pinworm (Enterobius vermicularis) infestation."},
'**Differential Diagnosis**': {'Most Likely Cause': 'Enterobius vermicularis (pinworm) infestation',
'Explanation': "The nocturnal itching described, along with the location of the excoriations, strongly suggests Enterobius vermicularis as the causative agent. Pinworm eggs are deposited around the anus by the adult female worm, causing intense itching, especially at night. The history of improving diarrhea could be coincidental or related to the stress on the child's gastrointestinal system."},
'**Treatment Step by Step**': {'Step 1': "Confirm the diagnosis by identifying pinworm eggs. This is typically done with the 'tape test,' where clear adhesive tape is applied to the perianal area in the morning before washing or using the toilet, then examined microscopically for eggs.",
'Step 2': 'Treat with a specific antiparasitic medication. The first-line treatment for pinworm infestation is either albendazole or mebendazole, given as a single dose and then repeated in two weeks to prevent reinfection.',
'Step 3': 'Educate the family on hygiene measures to prevent reinfection and the spread to other family members. This includes washing hands and nails thoroughly, especially after using the toilet and before eating; changing underwear, bed linen, and towels daily; and avoiding scratching the perianal area.',
'Step 4': 'Consider treating all household members simultaneously to prevent reinfection, given the high contagion rate among individuals living in close quarters.'},
'**Definitive Diagnosis**': 'Enterobius vermicularis (pinworm) infestation',
'**Treatment**': {'Definitive Treatment': 'Treatment consists of administering a single dose of albendazole or mebendazole, repeated in two weeks. Alongside, stringent personal and household hygiene practices should be emphasized to prevent reinfection and spread.'}}

Prompt 2/ 2nd time

CRR of Diagnosis and Clinical Management          1st Case Scenario

Figure 4: The outputs of clinical notes with the first case scenario, level one notes. At phase 2, The CRR of Diagnosis and Clinical Management, using Prompts 2 and 3.

**Verified CRR:**
{'**Verification Queries**': {'Query 1': 'What are the clinical practice guidelines for diagnosing Enterobius vermicularis (pinworm) infestation?',
'Query 2': 'What are the first-line treatments for Enterobius vermicularis (pinworm) infestation according to current clinical guidelines?',
'Query 3': 'What hygiene measures are recommended in clinical practice guidelines to prevent pinworm reinfection and spread?',
'Query 4': "Is the 'tape test' recommended by clinical practice guidelines for the diagnosis of pinworm infestation?",
'Query 5': 'Are there any guidelines that recommend treating all household members when one person is diagnosed with pinworm infestation?'},
'**Retrieved Information**': {'Response 1': "Clinical practice guidelines recommend the diagnosis of Enterobius vermicularis (pinworm) infestation to be confirmed through visualization of worms or by identifying eggs using the 'tape test'.", 'Response 2': 'Guidelines confirm that the first-line treatment for pinworm infestation is either albendazole or mebendazole, given as a single dose and repeated in two weeks.', 'Response 3': 'To prevent reinfection and spread, guidelines recommend stringent hygiene measures including washing hands, changing bed linens and underwear frequently, and maintaining clean living environments.', 'Response 4': "Yes, the 'tape test' is specifically recommended for diagnosing Enterobius vermicularis infestation.", 'Response 5': 'Clinical practice guidelines do recommend considering the treatment of all household members simultaneously to prevent reinfection, reflecting the high contagion rate of pinworm among individuals living in close contact.'},
'**Comparison**': {'Clinical Findings': "The clinical findings and the use of the 'tape test' for diagnosis align with clinical practice guidelines.",
'Causes': 'The correlations between symptoms and Enterobius vermicularis as the cause are consistent with the expected clinical presentation.',
'Treatment': 'The treatment plan including the use of albendazole or mebendazole and the emphasis on hygiene measures, matches the recommendations from the clinical practice guidelines.'},
'**Discrepancy**': None,
'**Final COT**': {'**Differential Diagnosis Process**': {'Step 1': 'Given the symptoms of genital itching, particularly noted to be worse at night, and the presence of excoriations over the anus and near the vagina, parasitic infections are highly suspected.', 'Step 2': 'The recent history of diarrhea could suggest a gastrointestinal infection or infestation, which may be related to the itching if the pathogen has a cutaneous manifestation or provokes such a reaction.', 'Step 3': 'The recent use of amoxicillin for an ear infection introduces the possibility of a secondary yeast infection due to antibiotic use; however, the description of symptoms being worse at night and the specific location of itching are less typical for a yeast infection.', 'Step 4': "Considering the girl's age and the specific symptoms described, common infections of the perianal region in children should be considered, such as pinworm (Enterobius vermicularis) infestation."},
'**Definitive Cause**': {'Most Likely Pathogen/Cause': 'Enterobius vermicularis (pinworm) infestation'},
'**Definitive Diagnosis**': 'Enterobius vermicularis (pinworm) infestation',
'**Treatment Plan**': {'Step 1': "Confirm the diagnosis by identifying pinworm eggs with the 'tape test'.", 'Step 2': 'Treat with a specific antiparasitic medication, either albendazole or mebendazole, given as a single dose and then repeated in two weeks.', 'Step 3': 'Educate the family on hygiene measures to prevent reinfection and the spread to other family members.', 'Step 4': 'Consider treating all household members simultaneously to prevent reinfection.'}}}

Prompt 4

CRR of Diagnosis and Clinical Management     1st Case Scenario

Figure 5: The verified CRR output of clinical notes, corresponding to the first case scenario and level one notes. At stage 3, the CRR of Diagnosis and Clinical Management, using Prompt 4. The Final CoT is the verified CRR. CRR; Clinical Reasoning Rationale

**Final Note:**
{"error_flag":1,
"error_location":10,
"sentence_correction":"Suspected of **Enterobius vermicularis** (pinworm) infestation."}

Prompt 5

Medical Errors Detection and Correction     1st Case Scenario

Figure 6: The last structured output includes: "error flag" to identify if an error exists. And "error location" and "sentence correction" if an error within a clinical note corresponds to the first case scenario and level one notes. At stage 4, the Medical Errors Detection and Correction, using Prompt 5.

**Preprocessed Note:**
A 31-year-old woman comes to the physician because of intermittent episodes of stabbing right lower jaw pain for 6 weeks. The pain is severe, sharp, and lasts for a few seconds. These episodes commonly occur when she washes her face, brushes her teeth, or eats a meal. She does not have visual disturbances, weakness of her facial muscles, or hearing loss. Five weeks ago, she had an episode of acute bacterial sinusitis, which was treated with antibiotics. Treatment with **doxepin** is initiated.

Prompt 1

Preprocessing     2nd Case Scenario     Deleted Part: {}

Figure 7: The outputs of clinical notes with the second case scenario, level two clinical notes. At stage 1, the Preprocessing, using Prompt 1.

**CRR Round 1:**
{'**Differential Diagnosis Step by Step**': {'Step 1': 'Given the location of the pain (right lower jaw) and the trigger factors (washing face, brushing teeth, eating), conditions involving the trigeminal nerve should be considered, as well as dental causes.', 'Step 2': 'The absence of visual disturbances, weakness of facial muscles, or hearing loss rules out conditions that involve other cranial nerves or central nervous system disorders.', 'Step 3': 'The recent history of acute bacterial sinusitis suggests a possible link between sinus inflammation and the jaw pain, but the nature of the pain described is not typical for sinusitis radiating pain.', 'Step 4': 'The nature of the pain (severe, sharp, and lasts for a few seconds) is highly suggestive of Trigeminal neuralgia, a disorder of the trigeminal nerve.', 'Step 5': 'Considering the treatment initiated (doxepin), which is an antidepressant used off-label for chronic pain management, further supports the neuralgia hypothesis over other dental or sinus-related causes.'},
'**Differential Diagnosis**': {'Most Likely Cause': 'Trigeminal Neuralgia',
'Explanation': "The patient's symptoms of intermittent episodes of stabbing pain in the right lower jaw triggered by facial movements or activities, along with her recent sinusitis episode and the absence of other neurological signs, strongly suggest Trigeminal Neuralgia. This condition involves the trigeminal nerve and is characterized by sudden, severe, brief, stabbing recurrent episodes of facial pain in one or more divisions of the nerve."},
'**Treatment Step by Step**': {'Step 1': 'Confirm the diagnosis of Trigeminal Neuralgia through a detailed patient history and physical examination to rule out other causes of facial pain.', 'Step 2': 'Initiate pharmacological treatment with carbamazepine or oxcarbazepine as the first-line treatment for trigeminal neuralgia, adjusting the dosage according to patient response and tolerance.', 'Step 3': 'Monitor for side effects of medication, such as dizziness, nausea, and diplopia.', 'Step 4': 'If the patient does not respond to or cannot tolerate the first-line medications, consider alternative treatments such as gabapentin, pregabalin, or low doses of tricyclic antidepressants like amitriptyline or nortriptyline.', 'Step 5': 'In cases where medication does not provide adequate relief, or for individuals who prefer not to take medication long term, refer the patient to a neurosurgeon for consideration of surgical options such as microvascular decompression or gamma knife radiosurgery.', 'Step 6': 'Provide patient education regarding the nature of the disease, the importance of medication adherence, and possible side effects.'},
'**Definitive Diagnosis**': 'Trigeminal Neuralgia',
'**Treatment**': {'Definitive Treatment': 'First-line pharmacological treatment with carbamazepine or oxcarbazepine, with consideration for surgical intervention in refractory cases.'}}

**Prompt 2**

**Clinical Note (deleted diagnosis with 0):** A 31-year-old woman comes to the physician because of intermittent episodes of stabbing right lower jaw pain for 6 weeks. The pain is severe, sharp, and lasts for a few seconds. These episodes commonly occur when she washes her face, brushes her teeth, or eats a meal. She does not have visual disturbances, weakness of her facial muscles, or hearing loss. Five weeks ago, she had an episode of acute bacterial sinusitis, which was treated with antibiotics. Treatment with doxepin is initiated.

**Prompt 3**

**CRR Round 2:**
{'**Differential Diagnosis Step by Step**': {'Step 1': 'Given the location of pain in the right lower jaw and its triggering by mechanical stimuli such as washing the face and brushing teeth, the differential diagnosis initially includes dental causes like dental abscess, temporomandibular joint disorders, and neuralgias.', 'Step 2': 'The absence of dental pathology signs (e.g., toothache, swelling), temporomandibular joint dysfunction symptoms (e.g., jaw clicking, difficulty chewing), and the description of the pain as intermittent, stabbing, and severe, directs the differential towards neuralgias.', 'Step 3': "Excluding facial muscle weakness, hearing loss, and visual disturbances rules out conditions such as Bell's palsy, acoustic neuroma, and multiple sclerosis.", 'Step 4': 'The history of acute bacterial sinusitis treated with antibiotics suggests a possible relationship; however, the lack of ongoing sinusitis symptoms decreases the likelihood of this being the direct cause.', 'Step 5': 'The pain experienced during common stimuli like eating and brushing teeth specifically implicates triggers typical for trigeminal neuralgia.'},
'**Differential Diagnosis**': {'Most Likely Cause': 'Trigeminal Neuralgia', 'Explanation': 'Given the characteristics of the pain (severe, sharp, and lasts for a few seconds) occurring in the distribution of the trigeminal nerve, and triggered by touch, combined with the exclusion of other facial or dental pathologies, trigeminal neuralgia is the most likely diagnosis.'},
'**Treatment Step by Step**': {'Step 1': 'Initiate pharmacologic treatment with anticonvulsants, considering carbamazepine as the first-line treatment due to its effectiveness in relieving neuralgic pain.', 'Step 2': 'If carbamazepine is not tolerated or contraindicated, other medications such as gabapentin, lamotrigine, or oxcarbazepine may be considered as alternatives.', 'Step 3': 'Monitor for side effects of the medication and adjust dosage accordingly to optimize pain control while minimizing adverse effects.', 'Step 4': 'In cases where medication is ineffective or side effects are intolerable, referral to a neurosurgeon for evaluation for surgical intervention, such as microvascular decompression, may be considered.', 'Step 5': 'Provide patient education on trigeminal neuralgia, including potential triggers to avoid and the importance of adherence to medication for pain management.'},
'**Definitive Diagnosis**': 'Trigeminal Neuralgia',
'**Treatment**': neuralgic conditions, but the first-line treatment for trigeminal neuralgia should be carbamazepine, due to its proven efficacy in this condition. Other options include gabapentin and surgical interventions for refractory cases.'}}

**Prompt 2/ 2nd time**

**CRR of Diagnosis and Clinical Management**  **2nd Case Scenario**

Figure 8: The outputs of clinical notes with the second case scenario, level two notes. At stage 2, The CRR of Diagnosis and Clinical Management, using Prompts 2 and 3.

**Verified CRR:**
{'**VerificationQueries**': {'Query 1': 'What are the clinical practice guidelines for diagnosing trigeminal neuralgia?', 'Query 2': 'What are the first-line and alternative treatments for trigeminal neuralgia according to the latest clinical practice guidelines?', 'Query 3': 'Are there any recommended diagnostic criteria for trigeminal neuralgia that include the absence of other neurological signs?', 'Query 4': 'Does the clinical practice guideline recommend carbamazepine as the first-line medication for trigeminal neuralgia? What are the alternatives if carbamazepine is contraindicated or not tolerated?', 'Query 5': 'What are the indications for surgical intervention in trigeminal neuralgia according to clinical practice guidelines?'},
'**Retrieved Information**': {'Response 1': 'The latest clinical practice guidelines for diagnosing trigeminal neuralgia involve the identification of characteristic pain symptoms, such as sudden, brief, and severe episodes of facial pain along the distribution of the trigeminal nerve, without other neurological deficits.', 'Response 2': 'According to the guidelines, carbamazepine is recommended as the first-line treatment for trigeminal neuralgia, followed by oxcarbazepine. Gabapentin or lamotrigine are considered alternatives when the first-line treatments are not tolerated.', 'Response 3': 'The guidelines specify that diagnosis should consider the absence of other neurological signs to rule out secondary causes of facial pain.', 'Response 4': 'Yes, the guidelines confirm carbamazepine as the first-line therapy for trigeminal neuralgia and list gabapentin, lamotrigine, and oxcarbazepine as alternatives.', 'Response 5': 'Surgical intervention, such as microvascular decompression, is recommended for patients who do not respond to pharmacological treatment or who cannot tolerate the medication.'},
'**Comparison**': {'Clinical Findings': 'The differential diagnosis steps align well with guidelines, focusing on specific pain characteristics and ruling out other pathologies.',
'Causes': "The identified cause, trigeminal neuralgia, is directly supported by the absence of other neurological deficits and the pain's distribution, which matches guideline criteria.",
'Treatment': 'The initial treatment plan aligns with guidelines, recommending carbamazepine as first-line treatment. Alternatives and surgical intervention for refractory cases are also in agreement with recommended practices.'},
'**Discrepancy**': None,
'**Final COT**': {'**Differential Diagnosis Process**': {'Step 1': 'Given the location of pain in the right lower jaw and its triggering by mechanical stimuli such as washing the face and brushing teeth, the differential diagnosis initially includes dental causes like dental abscess, temporomandibular joint disorders, and neuralgias.', 'Step 2': 'The absence of dental pathology signs (e.g., toothache, swelling), temporomandibular joint dysfunction symptoms (e.g., jaw clicking, difficulty chewing), and the description of the pain as intermittent, stabbing, and severe, directs the differential towards neuralgias.', 'Step 3': "Excluding facial muscle weakness, hearing loss, and visual disturbances rules out conditions such as Bell's palsy, acoustic neuroma, and multiple sclerosis.", 'Step 4': 'The history of acute bacterial sinusitis treated with antibiotics suggests a possible relationship; however, the lack of ongoing sinusitis symptoms decreases the likelihood of this being the direct cause.', 'Step 5': 'The pain experienced during common stimuli like eating and brushing teeth specifically implicates triggers typical for trigeminal neuralgia.'},
'**Definitive Cause**': {'Most Likely Pathogen/Cause': 'Trigeminal Neuralgia'}, '
'**Definitive Diagnosis**': 'Trigeminal Neuralgia', '
'**Treatment Plan**': {'Step 1': 'Initiate pharmacologic treatment with anticonvulsants, considering carbamazepine as the first-line treatment due to its effectiveness in relieving neuralgic pain.', 'Step 2': 'If carbamazepine is not tolerated or contraindicated, other medications such as gabapentin, lamotrigine, or oxcarbazepine may be considered as alternatives.', 'Step 3': 'Monitor for side effects of the medication and adjust dosage accordingly to optimize pain control while minimizing adverse effects.', 'Step 4': 'In cases where medication is ineffective or side effects are intolerable, referral to a neurosurgeon for evaluation for surgical intervention, such as microvascular decompression, may be considered.', 'Step 5': 'Provide patient education on trigeminal neuralgia, including potential triggers to avoid and the importance of adherence to medication for pain management.'}}}

**Prompt 4**

**CRR of Diagnosis and Clinical Management**      **2nd Case Scenario**

Figure 9: The verified CRR output of clinical notes corresponds to the second case scenario and level two notes. At stage 3, the CRR of Diagnosis and Clinical Management, using Prompt 4. The Final CoT is the verified CRR. CRR; Clinical Reasoning Rationale

**Final Note:**
{"**error_flag**":1
,"**error_location**":5,
"**sentence_correction**":"Treatment with **carbamazepine** is initiated."}

**Prompt 5**

**Medical Errors Detection and Correction**      **2nd Case Scenario**

Figure 10: The last structured output includes: "error flag" to identify if an error exists. And "error location" and "sentence correction" if an error within a clinical note corresponds to the second case scenario and level two notes. At stage 4, the Medical Errors Detection and Correction, using Prompt 5.

**Preprocessed Note:**
This patient is also being managed for chronic respiratory failure, POA (present on admission). Mr. <NAME/> is a <AGE/> yr male with Hg SS sickle cell disease, admitted for sickle pain crisis. Pt has chronic **Epiglottitis**: Likely <DATE/> anemia and chronic illness. - Cont home O2, >2L at all times - RT following for nocturnal BiPap.

Prompt 1

( Preprocessing )   ( 3rd Case Scenario )   ( Deleted Part: {} )

Figure 11: The outputs of clinical notes with the third case scenario, level two clinical notes. At stage 1, the Preprocessing, using Prompt 1.

---

**CRR Round 1:**
{'**Differential Diagnosis Step by Step**': {'Step 1': "The patient's history of Hg SS sickle cell disease and chronic respiratory failure suggests complications related to these conditions.", 'Step 2': 'The presentation of sickle pain crisis indicates an exacerbation of the sickle cell disease, potentially triggering acute chest syndrome or other vaso-occlusive complications.', 'Step 3': 'Chronic epiglottitis, chronic anemia, and chronic illness might indicate ongoing or recurrent infections, inflammation, or complications from sickle cell disease.'},
'**Differential Diagnosis**': {'Most Likely Cause': 'Complications associated with Hg SS sickle cell disease.',
'Explanation': "Given the patient's chronic conditions and the presentation of a sickle pain crisis, complications of sickle cell disease are likely contributing to the current state. These can include acute chest syndrome, exacerbated by chronic respiratory failure, or other vaso-occlusive events leading to organ damage and pain. Chronic epiglottitis could be a result of recurrent infections or chronic inflammation."}, '
'**Treatment Step by Step**': {'Step 1': "Manage the sickle pain crisis with hydration, oxygen therapy, and pain management tailored to the patient's needs.", 'Step 2': 'Evaluate the need for transfusion or exchange transfusion to address anemia and reduce sickle hemoglobin levels.', 'Step 3': 'Continue home oxygen therapy and monitor the need for nocturnal BiPAP to manage chronic respiratory failure.', 'Step 4': 'Assess the epiglottitis and consider antibiotics if a bacterial infection is suspected, along with steroids to reduce inflammation.', 'Step 5': 'Long-term, explore disease-modifying treatments for sickle cell disease, such as hydroxyurea or stem cell transplantation, to reduce the frequency of vaso-occlusive crises and complications.'},
'**Definitive Diagnosis**': 'Exacerbation of Hg SS sickle cell disease leading to acute sickle pain crisis and potential acute chest syndrome, coupled with chronic respiratory failure.',
'**Treatment**': {'Definitive Treatment': 'The treatment involves managing the acute crises and the chronic conditions concurrently. This includes pain management, possibly transfusions, continued oxygen therapy and possibly BiPAP for respiratory support, antibiotics and steroids if an infection or significant inflammation is present, and considering disease-modifying therapies for long-term management of sickle cell disease.'}}

Prompt 2

**Clinical Note (deleted diagnosis with 0):** This patient is also being managed for chronic respiratory failure, POA (present on admission). Mr. <NAME/> is a <AGE/> yr male with Hg SS sickle cell disease, admitted for sickle pain crisis. Pt has chronic **0**: Likely <DATE/> anemia and chronic illness. - Cont home O2, >2L at all times - RT following for nocturnal BiPap.

Prompt 3

**CRR Round 2:**
{'**Differential Diagnosis Step by Step**': {'Step 1': 'Identify key clinical features: Patient with known Hg SS sickle cell disease presented for pain crisis, has chronic anemia, chronic respiratory failure, requires home oxygen therapy (>2L), and nocturnal BiPAP.', 'Step 2': 'Consideration of common complications in sickle cell disease: Acute pain crisis, acute chest syndrome, pulmonary hypertension, and chronic organ damage.', 'Step 3': 'Analysis of provided information regarding respiratory management suggests chronic respiratory complications, possibly pulmonary hypertension or hypoxemia secondary to chronic sickle cell lung disease.', 'Step 4': 'Given the chronic nature of the respiratory failure requiring oxygen and BiPAP, identify potential causes such as vaso-occlusive episodes leading to pulmonary infarction or chronic acute chest syndromes contributing to progressive lung function decline.'},
'**Differential Diagnosis**': {'Most Likely Cause': 'Pulmonary Hypertension secondary to chronic sickle cell disease.',
'Explanation': 'Chronic anemia and repeated vaso-occlusive crises can lead to pulmonary hypertension, a common complication in patients with chronic sickle cell disease. This would also explain the requirement for continuous oxygen therapy and nocturnal BiPAP to manage resultant chronic respiratory failure.'},
'**Treatment Step by Step**': {'Step 1': 'Continue current supportive measures, including oxygen therapy and nocturnal BiPAP, to manage chronic respiratory failure.', 'Step 2': 'Management of sickle cell disease exacerbations with hydration, pain control, possibly blood transfusions to reduce sickling episodes.', 'Step 3': 'Consider initiating or optimizing specific therapies for pulmonary hypertension, which could include phosphodiesterase-5 inhibitors (e.g., sildenafil) or other pulmonary vasodilators.', 'Step 4': 'Regular follow-up with a comprehensive sickle cell management team, including specialists in hematology, pulmonology, and potentially cardiology.', 'Step 5': 'Evaluation and possible adjustment of current hydroxyurea therapy to reduce the frequency of pain crises and other complications related to sickle cell disease.', 'Step 6': 'Encourage vaccinations and prophylactic antibiotics as appropriate to reduce the risk of infections leading to acute chest syndrome.'},
'**Definitive Diagnosis**': 'Chronic Respiratory Failure secondary to Pulmonary Hypertension in the setting of Sickle Cell Disease.',
'**Treatment**': {'Definitive Treatment': 'Optimization of sickle cell disease management including pain crisis management, chronic respiratory failure support with oxygen and BiPAP, consideration of pulmonary hypertension targeted therapy, hydroxyurea therapy optimization, and preventive measures for complications.'}}

Prompt 2/ 2nd time

( CRR of Diagnosis and Clinical Management )   ( 3rd Case Scenario )

Figure 12: The outputs of clinical notes with the third case scenario, level two notes. At stage 2, The CRR of Diagnosis and Clinical Management, using Prompts 2 and 3.

**Verified CRR:**
{'VerificationQueries': {'Query 1': 'What are the American Thoracic Society guidelines for diagnosing pulmonary hypertension in the context of sickle cell disease?', 'Query 2': 'What are the clinical findings indicative of pulmonary hypertension as per the latest guidelines?', 'Query 3': "What are the National Heart, Lung, and Blood Institute's recommendations for the management of sickle cell disease complications, particularly pulmonary hypertension?", 'Query 4': 'What are the recommended treatments for chronic respiratory failure secondary to pulmonary hypertension in sickle cell disease patients according to the most recent clinical practice guidelines?'},
'Retrieved Information': {'Response 1': 'The American Thoracic Society guidelines recommend echocardiographic screening for pulmonary hypertension in sickle cell disease patients, with confirmatory diagnosis requiring right heart catheterization.', 'Response 2': 'Clinical findings indicative of pulmonary hypertension include dyspnea, chest pain, and signs of right heart failure. Echocardiography may show elevated tricuspid regurgitant jet velocity suggesting pulmonary hypertension.', 'Response 3': 'The National Heart, Lung, and Blood Institute recommends comprehensive care including hydroxyurea therapy to reduce sickle cell-related complications, aggressive treatment of acute chest syndrome, and evaluation for pulmonary hypertension in patients with sickle cell disease.', 'Response 4': 'Recommended treatments for chronic respiratory failure in sickle cell disease include oxygen therapy, treatment of underlying pulmonary hypertension with medications such as phosphodiesterase-5 inhibitors, and management of sickle cell disease exacerbations.',
'Comparison': {'Clinical Findings': 'Echocardiographic screening and right heart catheterization for pulmonary hypertension align with proposed diagnosis considerations.',
'Causes': 'Chronic anemia and vaso-occlusive crises leading to pulmonary hypertension are consistent with suggested complications of sickle cell disease.',
'Treatment': 'Supportive measures, specific therapies for pulmonary hypertension, hydroxyurea optimization, and preventive measures align with guideline recommendations.'},
'Discrepancy': None,
'Final COT': {'Differential Diagnosis Process': {'Step 1': 'Identify key clinical features: Patient with known Hg SS sickle cell disease presented for pain crisis, has chronic anemia, chronic respiratory failure, requires home oxygen therapy (>2L), and nocturnal BiPAP.', 'Step 2': 'Consideration of common complications in sickle cell disease: Acute pain crisis, acute chest syndrome, pulmonary hypertension, and chronic organ damage.', 'Step 3': 'Analysis of provided information regarding respiratory management suggests chronic respiratory complications, possibly pulmonary hypertension or hypoxemia secondary to chronic sickle cell lung disease.', 'Step 4': 'Given the chronic nature of the respiratory failure requiring oxygen and BiPAP, identify potential causes such as vaso-occlusive episodes leading to pulmonary infarction or chronic acute chest syndromes contributing to progressive lung function decline.'},
'Definitive Cause': {'Most Likely Pathogen/Cause': 'Pulmonary Hypertension secondary to chronic sickle cell disease.'},
'Definitive Diagnosis': 'Chronic Respiratory Failure secondary to Pulmonary Hypertension in the setting of Sickle Cell Disease.', 'Treatment Plan': {'Step 1': 'Continue current supportive measures, including oxygen therapy and nocturnal BiPAP, to manage chronic respiratory failure.', 'Step 2': 'Management of sickle cell disease exacerbations with hydration, pain control, possibly blood transfusions to reduce sickling episodes.', 'Step 3': 'Consider initiating or optimizing specific therapies for pulmonary hypertension, which could include phosphodiesterase-5 inhibitors (e.g., sildenafil) or other pulmonary vasodilators.', 'Step 4': 'Regular follow-up with a comprehensive sickle cell management team, including specialists in hematology, pulmonology, and potentially cardiology.', 'Step 5': 'Evaluation and possible adjustment of current hydroxyurea therapy to reduce the frequency of pain crises and other complications related to sickle cell disease.', 'Step 6': 'Encourage vaccinations and prophylactic antibiotics as appropriate to reduce the risk of infections leading to acute chest syndrome.'}}}

Prompt 4

CRR of Diagnosis and Clinical Management      3rd Case Scenario

Figure 13: The verified CRR output of clinical notes corresponds to the third case scenario and level two notes. At stage 3, the CRR of Diagnosis and Clinical Management, using Prompt 4. The Final CoT is the verified CRR. CRR; Clinical Reasoning Rationale

**Final Note:**
{"error_flag":1
,"error_location":3,
**"sentence_correction":**"Pt has chronic **Hypoxia** likely due to pulmonary hypertension secondary to sickle cell disease."}

Prompt 5

Medical Errors Detection and Correction      3rd Case Scenario

Figure 14: The last structured output includes: "error flag" to identify if an error exists. And "error location" and "sentence correction" if an error within a clinical note corresponds to the third case scenario and level two notes. At stage 4, the Medical Errors Detection and Correction, using Prompt 5.

# Overview of the 2024 Shared Task on Chemotherapy Treatment Timeline Extraction

*Jiarui Yao[1], *Harry Hochheiser[2], WonJin Yoon[1], Eli Goldner[1], Guergana Savova[1]

[1]Boston Children's Hospital and Harvard Medical School

{jiarui.yao, wonjin.yoon, eli.goldner, guergana.savova}@childrens.harvard.edu

[2]University of Pittsburgh

harryh@pitt.edu

## Abstract

The 2024 Shared Task on Chemotherapy Treatment Timeline Extraction aims to advance the state of the art of clinical event timeline extraction from the Electronic Health Records (EHRs). Specifically, this edition focuses on chemotherapy event timelines from EHRs of patients with breast, ovarian and skin cancers. These patient-level timelines present a novel challenge which involves tasks such as the extraction of relevant events, time expressions and temporal relations from each document and then summarizing over the documents. De-identified EHRs for 57,530 patients with breast and ovarian cancer spanning 2004-2020, and approximately 15,946 patients with melanoma spanning 2010-2020 were made available to participants after executing a Data Use Agreement. A subset of patients is annotated for gold entities, time expressions, temporal relations and patient-level timelines. The rest is considered unlabeled data. In **Subtask1**, gold chemotherapy event mentions and time expressions are provided (along with the EHR notes). Participants are asked to build the patient-level timelines using gold annotations as input. Thus, the subtask seeks to explore the topics of temporal relations extraction and timeline creation if event and time expression input is perfect. In **Subtask2**, which is the realistic real-world setting, only EHR notes are provided. Thus, the subtask aims at developing an end-to-end system for chemotherapy treatment timeline extraction from patient's EHR notes. There were 18 submissions for Subtask 1 and 9 submissions for Subtask 2. The organizers provided a baseline system. The teams employed a variety of methods including Logistic Regression, TF-IDF, n-grams, transformer models, zero-shot prompting with Large Language Models (LLMs), and instruction tuning. The gap in performance between prompting LLMs and finetuning smaller-sized LMs indicates that for a challenging task such as patient-level chemotherapy timeline extraction, more sophisticated LLMs or prompting techniques are necessary in order to achieve optimal results as finetuing smaller-sized LMs outperforms by a wide margin.

## 1 Introduction

Cancer treatment is rarely simple. Complex protocols involving multiple drugs, given over extended period of times in specified orders, are the norm (Warner et al., 2019). This poses a challenge for clinical researchers. Ideally, real-world studies of the impact of specific protocols would require to identify which patients have been given which protocols. In practice, this task is complicated by a dearth of detailed information: although medication records and clinical notes might indicate the administration of a given chemotherapeutic agent to a patient, they rarely, if ever, name specific protocols. Furthermore, structured medication administration records are insufficient, as clinical notes may contain mentions of medications in the context of reasons for discontinuing treatment, prior treatments given at differing institutions, or reactions to treatment.

Extracting chemotherapy timelines from clinical notes involves a series of challenges. Individual mentions of relevant drug administrations (chemotherapy events) must be extracted and mapped to appropriate medication terminologies. Each event must then be assigned a time extent, based on the date of the note and any temporal modifiers and indicators (e.g. time expressions) identified alongside the medication event (Laparra et al., 2018). Finally, these individual instances must be ordered into a timeline. Each of these tasks involves substantial challenges, several of which have been the focus of previous SemEval challenges (Elhadad et al., 2015; Laparra et al., 2018; Bethard et al., 2017).
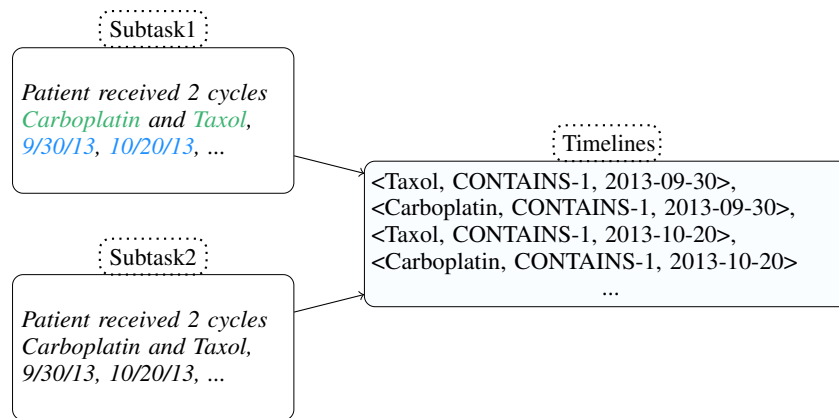
---

* indicates co-first authors.

Figure 1: Illustration of the two subtasks in the 2024 Chemotherapy Treatment Timeline Extraction shared task. The input of Subtask1 is patient notes with gold events (highlighted in green) and time expressions (highlighted in blue). The input of Subtask2 is patient notes only. The output of both subtasks is a list of chemotherapy treatment timelines with normalized time expressions. See details in section 2.

The 2015-2021 SemEval shared tasks (Bethard et al., 2015, 2016, 2017; Laparra et al., 2018, 2021) on temporal relation extraction from the clinical narrative used the THYME and THYME2 corpora (Styler IV et al., 2014; Wright-Bettner et al., 2020), each with a separate focus on one of the following tasks – pairwise temporal relation extraction, time expression normalization, and domain adaptation. The SemEval shared tasks provided the gold event and time expressions so that the teams focus on the temporal relation extraction to advance approach development. The state-of-the-art methodologies and results they established allowed the community to start exploring applications to real world biomedical use cases.

The 2024 Chemotherapy Treatment Timeline Extraction shared task[*] elevates the technical challenges to a new level by presenting participants with two challenges: assembling timelines from individual event mentions and temporal/time expressions provided as input (Subtask1), and building timelines directly from clinical notes, thus the real-world task of end-to-end extraction (Subtask2). Both subtasks go beyond the 2015-2021 Semeval shared tasks, however they build on the community knowledge advanced through them. For the 2024 Chemotherapy Treatment Timeline Extraction shared task the organizers provided a dataset of the Electronic Health Records (EHRs) of more than 73,000 cancer patients from 2004-2020 from University of Pittsburgh Medical Center (UPMC).

In the next sections, we describe the shared task,

its substasks, the dataset, the evaluation methodology, the baseline system, the teams with highlights of their approaches, and finally the results. Details of each team's approach is described in a separate paper by the team.

## 2 Description of the Shared Task and Subtasks

The overall goal of the task was to create patient-level timelines of *chemotherapy treatment events* from all the notes in the EHR available for a given patient. In general, timelines can be represented in different formats. We can describe a patient's treatment timeline in natural language, such as "*2 cycles Carboplatin and Taxol, 9/30/13, 10/20/13*" which is easy to understand by humans, however, it cannot be "understood" directly by machines. Over the years, the research community has developed a parsimonious set of relations to express temporality between two events or between an event and temporal/time expression (Wright-Bettner et al., 2020; Styler IV et al., 2014). We adopt these conventions where an event is any occurrence that can be positioned on a timeline (in our case chemotherapy events) and the set of temporal relations are defined as BEFORE, CONTAINS (with inverse CONTAINS-1 which is the equivalent of CONTAINED-BY), OVERLAP, NOTED-ON, BEGINS-ON, ENDS-ON. We limit events to only chemotherapy treatment events. Therefore, for the shared task we represent the chemotherapy treatment timelines in a computable format as a list of *<chemotherapy, temporal_relation, time_expression>* triplets. Thus, the previous ex-

| | Train | | | Dev | | | Test | | |
|---|---|---|---|---|---|---|---|---|---|
| | Patients | Notes | Words | Patients | Notes | Words | Patients | Notes | Words |
| Ovary | 26 | 1,675 | 1,183,632 | 8 | 562 | 308,814 | 8 | 559 | 257,116 |
| Breast | 33 | 1,002 | 465,644 | 16 | 499 | 225,588 | 35 | 1,333 | 786,896 |
| Melanoma | 10 | 233 | 124,924 | 3 | 211 | 178,308 | 10 | 229 | 156,083 |

Table 1: Gold labeled dataset: number of patients, notes, and words across train/dev/test sets. "Words" denotes the tokens delimited by white spaces.

| | Train | | | Dev | | | Test | | |
|---|---|---|---|---|---|---|---|---|---|
| | EVENT | TIMEX3 | TLINK | EVENT | TIMEX3 | TLINK | EVENT | TIMEX3 | |
| Ovary | 1,168 | 597 | 494 | 790 | 312 | 226 | 664 | 381 | |
| Breast | 1,023 | 576 | 455 | 279 | 146 | 113 | 2,560 | 1,118 | |
| Melanoma | 147 | 78 | 48 | 789 | 261 | 201 | 398 | 193 | |

Table 2: Gold labeled dataset: EVENTs/ TIMEX3s/ TLINKs distribution in the labeled dataset. TIMEX3 and TLINK refer to time expressions and temporal relations respectively.

ample can be converted to:

*<Carboplatin, CONTAINS-1, 2013-09-30>,*

*<Taxol, CONTAINS-1, 2013-09-30>,*

*<Carboplatin, CONTAINS-1, 2013-10-20>,*

*<Taxol, CONTAINS-1, 2013-10-20>.*

With this representation, the construction of chemotherapy treatment timelines can be naturally decomposed into the following stages: chemotherapy event extraction, time expression extraction, temporal relation classification, time expression normalization and patient-level timeline refinement. Time expressions are also referred to as temporal expressions and TIMEX3.

The shared task defined two subtasks. In **Subtask1**, gold chemotherapy event mentions and time expressions are provided (along with the EHR notes). Participants were asked to build the patient-level timelines using gold annotations as input. Thus, the subtask sought to explore the topics of temporal relation extraction and timeline creation if event and time expression input is perfect. In **Subtask2**, which is the realistic real-world setting, only EHR notes are provided. Thus, the subtask aimed at developing an end-to-end system for chemotherapy treatment timeline extraction from patient's EHR notes. Figure 1 is an overview of this task.

### 2.1 Data

The EHR for each patient included all types of available notes regardless of their relevance to the patient's cancer, e.g. radiology reports, pathology

notes, clinical notes, oncology notes, discharge summaries, progress reports, etc. We sampled a subset of patients to create the gold annotations. For the gold annotations, we follow the THYME2 annotation schema (Wright-Bettner et al., 2020; Styler IV et al., 2014) as it is widely used in the clinical temporal relation classification community (Bethard et al., 2015, 2016, 2017; Lin et al., 2019, 2021). Two domain experts created gold annotations of the chemotherapy events, time expressions, and temporal relations. These represent instance-level annotations. These pairwise gold annotations are in the Anafora [†] (Chen and Styler, 2013) xml format. The final gold patient-level timeline was created automatically by merging all instance-level annotations followed by deduplicating and collapsing temporal relations. The gold dataset was split into training, development (dev) and test sets. Table 1 and Table 2 present the distributions of the gold dataset (*the Labeled Dataset*).

Additionally, we provided the *Unlabeled Dataset* which consists of the UPMC EHR notes for 57,530 patients with breast and ovarian cancer, collected between 2004-2020, and 15,946 patients with melanoma, collected between 2010-2020. As implied by its name, this dataset does not have any gold annotations. The *Unlabeled dataset* could potentially be used for continued training of pre-trained language models or for pretraining a language model.

To access both *Labeled* and *Unlabeled* datasets, the PI (Principal Investigator) of each team was

---

[†] https://github.com/weitechen/anafora

required to execute a Data Use Agreement (DUA) with University of Pittsburgh. The process took 3-4 weeks on the average. Upon execution of DUAs, the data were distributed to the teams through Globus [‡] with gated Collections for each split and dataset. Globus provides a secure way of sharing the sensitive patient EHR data.

# 3 Evaluation

We used the standard F1 metric to evaluate system performance, with variations to reflect the real world use case of chemotherapy treatment timelines. In consultation with our oncology domain experts it was determined that the level of granularity most useful for both point of care and translational studies is the month and the year for the chemotherapy treatment; the exact date was not deemed critical.

Therefore, we designed four evaluation strategies with different levels of granularity: strict, relaxed-to-day, relaxed-to-month and relaxed-to-year. Strict evaluation requires all elements in a triplet to match the corresponding ones in the gold annotations to count as a match. In all relaxed evaluations, we consider certain temporal relations interchangeable, and only compare the predicted month (relaxed-to-month) or year (relaxed-to-year) with the gold ones. For instance, under relaxed-to-month evaluation, we consider <TC, BEGINS-ON, 2013-02> correct if the gold timeline is <TC, BEGINS-ON, 2013-02-13>. In this shared task, based on our consultations with our oncology domain experts as described above we use the relaxed-to-month metric as the official score for the leader board and rankings.

Our scoring metrics account for differences in patterns of chemotherapy treatments. Most, but not all patients have chemotherapy. Some melanoma patients, for example, are treated surgical with no chemotherapy. To handle these differences, we used two types of scores based on relaxed-to-month results as motivated above:

- Type A: F1 where all patients are included regardless of whether they have chemotherapy gold timelines.

- Type B: F1 where patients with no chemotherapy timelines are excluded.

Type A score aims to catch false positives for these patients. Type B score measures the effectiveness of the methods on patients with confirmed chemotherapy treatments. The F1 score for each patient was computed and the final F1 score for each type is the average across all patients. The Official score used for the rankings in the Leader Board is the average of Type A and Type B. A link to the evaluation script[§] is posted on the shared task website.

Teams uploaded their systems output into their gated Globus collection and the organizers ran the evaluation script to produce the results posted on the Leader Board on the shared task website. Each team was allowed to upload up to three submissions for each task.

# 4 Baseline Systems

The shared task organizers provide baseline results for Subtask1 and Subtask2.

For both subtasks we used Apache cTAKES[¶] (Savova et al., 2010) for sentence boundary detection, tokenization, and pipelining of software components via the Python bridge to Java (`ctakes-pbj`) module. We use Huggingface Transformers (Wolf et al., 2019) for model training and inference, and CLUlab Timenorm's synchronous context free grammar module (Bethard, 2013) for normalizing time expressions to ISO standard. The system processes all the patients and notes for a given cancer type and split of the dataset. We processed patients by cancer type and dataset split since there are overlapping patient identifiers across different cancer types and splits (although the patients are different).

## 4.1 Subtask1

We used cTAKES' default tokenization and sentence splitting stack, then loaded chemotherapy event mentions and time expressions from the annotated gold data. We normalized as many time expressions as possible using Timenorm. Taking all the relevant pairs of chemotherapy event mentions and normalized time expressions, i.e. within a certain number of tokens from each other, we generated instances for classification by our temporal relation model (described below). Following (Lin et al., 2021), we used tags to distinguish the chemotherapy event mentions from the time expressions, e.g. *The patient received <e> paclitaxel*

---

*</e> on <t> February 2nd, 2011 </t>*. Note, in the generated instance we used the original text of the time expression, not its normalized form from Timenorm (i.e. `2011-02-02`). The normalized form is associated with its source time expression in a data structure within cTAKES and is used later when collecting instances for summarization and scoring.

For the temporal relation classification model we used Microsoft Research's PubMedBERT (Gu et al., 2020), and first fine-tuned on the THYME2 clinical temporal relation dataset (Wright-Bettner et al., 2020), then continue fine-tuned on the shared task training set to produce the type of temporal relation. Finally, when all the pairs have been classified, we generated a text table, with a row for each classified pair. Each row contains the original text of the chemotherapy mention, the normalized form of its paired time expression, their predicted temporal relation, and the identifier of the patient with whom this instance is associated. We then processed this table into a collection of summarized patient-level timelines for each patient.

To derive the patient-level timelines, we refined the pairwise temporal relations by 1) deduplication, and 2) choosing the most specific temporal relation between a chemotherapy treatment and a time expression following a predefined label hierarchy (BEGINS-ON/ENDS-ON > CONTAINS/CONTAINS-1 > BEFORE). In addition, for generic chemotherapy mentions such as "chemo" and "chemotherapy", we added them to the final timelines only if there was not a more specific chemotherapy treatment (e.g. Taxol) having the same temporal relation with the exact same time expression.

### 4.2 Subtask2

Here we also used cTAKES' default sentence detection and tokenization stack. For detecting chemotherapy mentions, we used cTAKES' dictionary lookup module with a customized dictionary of common chemotherapy terms collected from the training split of the shared task gold annotated corpus to identify potential chemotherapy mentions in each note. For detecting time expressions, we used the SVM-based tagger in the cTAKES' temporal module to identify potential time expressions, then normalize as many potential time expressions with Timenorm as possible. As in Subtask1, we generated instances for temporal relation classification from all relevant pairs of chemotherapy mentions

and normalized time expressions, along with a table of the classified instances and relevant associated information for further summarization and evaluation. We used the same model for temporal relation classification as in Subtask1. We provided a docker implementation [‖] of the baseline system for Subtask 2 as a resource on the shared task website.

## 5 Participating Systems

In this section, we briefly describe the approaches of participating systems. Details of each system can be found in the separate papers by each of the team.

The participants explored a variety of methods, including Logistic Regression, TF-IDF, n-grams, transformer models, zero-shot prompting with Large Language Models (LLMs), and instruction tuning. Table 3 summarizes all teams' approaches.

**BioCom** participated in Subtask 1. They utilized SciSpacy for Named Entity Extraction (NER) and Logistic Regression to classify temporal relations. They used unigram Term Frequency-Inverse Document Frequency (TF-IDF) to get features from the input text.

**ClinicalRxMiners** submitted two systems for Subtask 1. In submission 1, ClinicalRxMiners utilized a machine learning (non-deep learning) approach and employed n-grams as features of the input, with a soft voting classifier as the model for making predictions. In submission 2, Clinical-RxMiners utilized a pretrained Language Model (LM) named GLiNER (Zaratiana et al., 2023), which is specialized for NER.

**KCLab** (Tan et al., 2024) utilized a hybrid method, employing cTAKES (Savova et al., 2010) for preprocessing and PubMedBERT (Gu et al., 2020) for post-processing. Their system was built on top of the baseline model provided by the organizers. Additionally, KCLab used the UMLS (Bodenreider, 2004) database. KCLab participated in both Subtask1 and Subtask2.

**LAILab** (Haddadan et al., 2024) utilized two approaches: supervised fine-tuning of language models and a pipeline approach combining rule-based NER with deep learning based relation classification. For Subtask 1, they finetuned

---

[‖] https://github.com/HealthNLPorg/chemoTimelinesBaselineSystem

561

| Teams | Approach | LM or Algorithm | Task |
|---|---|---|---|
| BioCom_submission1 | Machine Learning | Logistic Regression | 1 |
| ClinicalRXMiners_submission1 | Machine Learning | Soft voting classifier | 1 |
| ClinicalRXMiners_submission2 | Deep Learning | GLiNER Base | 1 |
| KCLab_submission1 | Finetuned LM | PubMedBert | 1, 2 |
| LAILab_submission1,2,3 | Finetuned LM | flan-T5-xxl, bart-large | 1, 2 |
| Lexicans-submission1,2,3 | Zero-shot Prompting | Llama2, Mistral, Zephyr, Meditron, and Mixtral | 1 |
| NLPeers_submission1 | Finetuned LM | deberta-v3-base | 1 |
| NLPeers_submission2 | Few-shot Prompting | Mixtral-8X7B-Instruct-v0.1 | 1 |
| NYULangone_submission1 | Zero-shot prompting | Mixtral 8x7B | 2 |
| UTSA-NLP_submission1,2,3 | Instruction tuning LM, continued pretraining LM | OpenChat-3.5-7B | 1, 2 |
| Wonder_submission1,2,3 | Finetuned LM | Bio-LM | 1, 2 |

Table 3: Characteristics of participating systems.

flan-T5-XXL (Chung et al., 2022). For Subtask 2, they used a sequence-to-sequence approach in the first two submissions, and a lookup table for chemotherapy event extraction with a deep learning method for temporal relation classification in the third submission.

**Lexicans** (Sharma et al., 2024) used LLMs with zero-shot prompting to extract relations. They also utilized the THYME ontology to formalize the representation of entities and their relationships. A few LLMs such as Llama2, Mistral, Zephyr, Meditron, and Mixtral (Touvron et al., 2023; Jiang et al., 2023; Tunstall et al., 2023; Chen et al., 2023) were tested under various settings. Additionally, a data normalization step was performed to transform time entities into absolute date-time formats.

**NLPeers** (Bannour et al., 2024) developed two systems, both submitted for Subtask1. For submission 1, NLPeers fine-tuned the `microsoft/deberta-v3-base` model and used it for temporal relation classification. Additionally, the Heideltime library[**] (Strötgen and Gertz, 2010) and an LLM-based prompt with the OpenChat 3.5 model (Wang et al., 2024a) were used to normalize time expressions. For submission 2, the NLPeers team applied few-shot prompting with the `Mixtral-8X7B-Instruct-v0.1` model (Jiang et al., 2023), the prompt was chosen by DSPy (Khattab et al., 2023), a framework for

algorithmically optimizing LM prompts. A Chain-Of-Thought (Wei et al., 2022) approach was integrated during the prompt searching step by DSPy. For time expression normalization, Heideltime was also used in submission 2.

**NYULangone** employed an LLM-based prompt approach with minimal pre- and post-processing. NYULangone participated only in Subtask2, which means the team did not use the gold annotation provided in Subtask1.

**UTSA-NLP** (Zhao and Rios, 2024) presented an instruction-tuning based approach. The UTSA-NLP team reformulated the task into a question-answering (QA) dataset for both the entity extraction step and temporal relation classification step, then instruction-tuned an LLM, `OpenChat-3.5-7B`, on the QA dataset. The team continued pre-training the instruction-tuned model on a portion of the *Unlabeled* dataset in one of their submissions. For the temporal relation classification step, they used an open-sourced LLM to generate reasoning for the answer.

**Wonder** (Wang et al., 2024b) participated in Subtasks 1 and 2. They employed a supervised fine-tuning approach, formulating the task as a multi-class sentence classification task, where the input was the text between the event and time expression. For Subtask 2, MedTagger[††] was used to identify all the potential EVENT-TIMEX3 pairs. Time ex-

---

[**]https://github.com/HeidelTime/heideltime

[††]https://github.com/OHNLP/MedTagger

| Submission | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission1 | 0.94 | 0.86 | 0.90 |
| LAILab_submission2 | 0.94 | 0.86 | 0.90 |
| LAILab_submission3 | 0.94 | 0.86 | 0.90 |
| Baseline_subtask1 | 0.93 | 0.85 | 0.89 |
| Wonder_submission2 | 0.90 | 0.78 | 0.84 |
| Wonder_submission1 | 0.89 | 0.77 | 0.83 |
| Wonder_submission3 | 0.88 | 0.73 | 0.80 |
| NLPeers_submission1 | 0.85 | 0.70 | 0.77 |
| BioCom_submission1 | 0.84 | 0.64 | 0.74 |
| Lexicans_submission1 | 0.81 | 0.61 | 0.71 |
| UTSA-NLP_submission3 | 0.80 | 0.58 | 0.69 |
| UTSA-NLP_submission1 | 0.80 | 0.58 | 0.69 |
| Lexicans_submission2 | 0.79 | 0.57 | 0.68 |
| UTSA-NLP_submission2 | 0.80 | 0.56 | 0.68 |
| NLPeers_submission2 | 0.76 | 0.52 | 0.64 |
| KCLab_submission1 | 0.76 | 0.49 | 0.63 |
| Lexicans_submission3 | 0.75 | 0.47 | 0.61 |
| ClinicalRXMiners_submission1 | 0.51 | 0.28 | 0.40 |
| ClinicalRXMiners_submission2 | 0.56 | 0.21 | 0.38 |

Table 4: Evaluation results of Subtask1 (test set). All scores are macro-F1 of relaxed-to-month setting. We compute two types of scores: F1 with patients with no gold timelines (Type A) and F1 without patients with no gold timelines (Type B). Official score is the average of Type A and Type B, which is used for the rankings in the leader board. See details in section 3.

| Submission | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission2 | 0.76 | 0.63 | 0.70 |
| Baseline_subtask2 | 0.67 | 0.48 | 0.58 |
| LAILab_submission1 | 0.65 | 0.47 | 0.56 |
| KCLab_submission1 | 0.63 | 0.45 | 0.54 |
| Wonder_submission3 | 0.59 | 0.46 | 0.53 |
| Wonder_submission2 | 0.59 | 0.46 | 0.52 |
| Wonder_submission1 | 0.58 | 0.46 | 0.52 |
| LAILab_submission3 | 0.47 | 0.47 | 0.47 |
| NYULangone_submission1 | 0.26 | 0.21 | 0.23 |
| UTSA-NLP_submission1 | 0.22 | 0.22 | 0.22 |

Table 5: Evaluation results, Subtask 2 (test set). All scores are macro-F1 of relaxed-to-month setting. We compute two types of scores: F1 with patients with no gold timelines (Type A) and F1 without patients with no gold timelines (Type B). Official score is the average of Type A and Type B, which is used for the rankings in the leader board. See details in section 3.

pressions were normalized with MedTime (Sohn et al., 2013).

## 6 Results and Discussion

Overall results are presented in Table 4 and 5. Results per type of cancer are presented in Table 6 and 7 in the Appendix.

Most teams employed deep-learning-based methods for this shared task. Two teams used non-deep-learning models: ClinicalRXMiners submission 1 used a machine learning model, BioCom trained a Logistic Regression system. For the event mention extraction step, the Wonder team and the baseline system used off-the-shelf tools for time expression extraction and normalization. LAILab used a lookup table for chemotherapy event identification

in one of their submissions. The other approaches employed in this shared task include end-to-end timeline building (meaning no separate steps for event mention extraction), supervised event mention extraction model, and zero-shot prompting.

**Finetuning LMs:** For both subtasks, the top teams, i.e. LAILab and Wonder, employed fine-tuned pretrained language models as the core technology. LAILab finetuned `Flan-T5-xxl` (Chung et al., 2022) and `Bart-large` (Lewis et al., 2020a), which have 11B and 400M parameters respectively. They achieve best performance on all subtasks (overall and per type of cancer) except for Subtask2, breast cancer. The Wonder team finetuned Bio-LM (Lewis et al., 2020b), yielding top 3 results across all subtasks (excluding the baseline system). The other two teams with good results are NLPeers and KCLab, who finetuned `deberta-v3-base` (He et al., 2023) and Pub-MedBert (Gu et al., 2020) respectively. Overall, the commendable performances of those teams suggest that finetuning LMs remains the optimal approach for optimizing system performance if gold labeled data and computing resources are available.

**Prompting LLMs:** A few teams took the approach of prompting LLMs. The Lexicans team experimented with zero-shot prompting of 5 different LLMs, namely LLAMA2 , Mistral, Zephyr, Meditron, and Mixtral (Touvron et al., 2023; Jiang et al., 2023; Tunstall et al., 2023; Chen et al., 2023). NYULangone applied zero-shot prompting with the Mixtral model. Submission 2 from the NLPeers team prompted the `Mixtral-8X7B-Instruct-v0.1` model in a few-shot fashion.

The gap in performance between prompting LLMs and finetuning smaller-sized LMs indicates that for a challenging task such as patient-level chemotherapy timeline extraction, more sophisticated LLMs or prompting techniques are necessary in order to achieve optimal results. The state-of-the-art results for the 2024 Chemotherapy Treatment Timeline Extraction shared task are established by fine-tuning smaller LMs.

A comparison of the scores between Subtask1 and Subtask2 shows a substantial drop of at least 0.2 F1 Official Score when gold event and time expressions (thus perfect input) are provided. This gap, surprisingly, implies that what is considered the easier task of event and time expression extraction is not a solved problem while the task of

temporal relation extraction holds strong.

# 7 Conclusion

The 2024 Shared Task on Chemotherapy Treatment Timeline Extraction is unique in both (1) focusing on a highly complex task, and (2) providing a large corpus of EHR data to the participants. The community embraced the task with enthusiasm and employed diverse methodologies, thus enabling robust comparison of approaches. Perhaps surprising in our current era of very large LMs, fine-tuned smaller LMs achieved superior performance. This discrepancy between prompting LLMs and finetuning smaller-sized LMs suggests that more sophisticated LLMs or prompting techniques are necessary in order to achieve optimal results for challenging tasks such as patient-level chemotherapy timeline extraction.

# 8 Acknowledgements

# Limitations

There are different types of cancer treatments, such as Immunotherapy, Radiation Therapy, Surgery and Targeted Therapy. In this shared task, we only focus on chemotherapy treatments. We leave the timeline construction of other types of therapy for future research.

# Ethics Statement

All the data used in this shared task are de-identified patient notes. The access the data, the PI of each team was required to execute a Data Use Agreement with University of Pittsburgh. The data were distributed through Globus, which provides a secure way of sharing senstitive data such as patient EHRs. Participants were also required to

submit the final timelines via Globus, to protect the patient privacy.

# References

Nesrine Bannour, Judith Jeyafreeda Andrew, and Marc Vincent. 2024. Team nlpeers at chemotimelines 2024: Evaluation of two timeline extraction methods, can generative llm do it all or is smaller model fine-tuning still relevant? In *Proceedings of the 6th Clinical NLP Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Steven Bethard. 2013. A synchronous context free grammar for time normalization. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 821–826, Seattle, Washington, USA. Association for Computational Linguistics.

Steven Bethard, Leon Derczynski, Guergana Savova, James Pustejovsky, and Marc Verhagen. 2015. SemEval-2015 task 6: Clinical TempEval. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 806–814, Denver, Colorado. Association for Computational Linguistics.

Steven Bethard, Guergana Savova, Wei-Te Chen, Leon Derczynski, James Pustejovsky, and Marc Verhagen. 2016. SemEval-2016 task 12: Clinical TempEval. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1052–1062, San Diego, California. Association for Computational Linguistics.

Steven Bethard, Guergana Savova, Martha Palmer, and James Pustejovsky. 2017. SemEval-2017 task 12: Clinical TempEval. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 565–572, Vancouver, Canada. Association for Computational Linguistics.

Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32 Database issue:D267–70.

Wei-Te Chen and Will Styler. 2013. Anafora: A web-based general purpose annotation tool. In *Proceedings of the 2013 NAACL HLT Demonstration Session*, pages 14–19, Atlanta, Georgia. Association for Computational Linguistics.

Zeming Chen, Alejandro Hern'andez Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Kopf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023. Meditron-70b: Scaling medical pretraining for large language models. *ArXiv*, abs/2311.16079.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models.

Noémie Elhadad, Sameer Pradhan, Sharon Gorman, Suresh Manandhar, Wendy Chapman, and Guergana Savova. 2015. SemEval-2015 task 14: Analysis of clinical text. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 303–310, Denver, Colorado. Association for Computational Linguistics.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.

Shohreh Haddadan, Tuan-Dung Le, Thanh Duong, and Thanh Q Thieu. 2024. Lailab at chemotimelines 2024: Finetuning sequence-to-sequence language models for temporal relation extraction towards cancer patient undergoing chemotherapy treatment. In *Proceedings of the 6th Clinical NLP Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pretraining with gradient-disentangled embedding sharing.

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

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*.

Egoitz Laparra, Xin Su, Yiyun Zhao, Özlem Uzuner, Timothy Miller, and Steven Bethard. 2021. SemEval-2021 task 10: Source-free domain adaptation for semantic processing. In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 348–356, Online. Association for Computational Linguistics.

Egoitz Laparra, Dongfang Xu, Ahmed Elsayed, Steven Bethard, and Martha Palmer. 2018. SemEval 2018 task 6: Parsing time normalizations. In *Proceedings of the 12th International Workshop on Semantic Evaluation*, pages 88–96, New Orleans, Louisiana. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoyanov. 2020b. Pretrained language models for biomedical and clinical tasks: Understanding and extending the state-of-the-art. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 146–157, Online. Association for Computational Linguistics.

Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2019. A BERT-based universal model for both within- and cross-sentence clinical temporal relation extraction. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 65–71, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2021. EntityBERT: Entity-centric masking strategy for model pretraining for the clinical domain. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 191–201, Online. Association for Computational Linguistics.

Guergana K Savova, James J Masanz, Philip V Ogren, Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-Schuler, and Christopher G Chute. 2010. Mayo clinical text analysis and knowledge extraction system (ctakes): architecture, component evaluation and applications. *Journal of the American Medical Informatics Association*, 17(5):507–513.

Vishakha Sharma, Andres Fernandez, Andrei Constantin Ioanovici, David Talby, and Frederik Buijs. 2024. Lexicans at chemotimelines 2024: Chemotimeline chronicles - leveraging large language models (llms) for temporal relations extraction in oncological electronic health records. In *Proceedings of the 6th Clinical NLP Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Sunghwan Sohn, Kavishwar B. Wagholikar, Dingcheng Li, Siddhartha R. Jonnalagadda, Cui Tao, K. E. Ravikumar, and Hongfang Liu. 2013. Comprehensive temporal information detection from clinical text: medical events, time, and tlink identification. *Journal of the American Medical Informatics Association : JAMIA*, 20 5:836–42.

Jannik Strötgen and Michael Gertz. 2010. HeidelTime: High quality rule-based extraction and normalization of temporal expressions. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 321–324, Uppsala, Sweden. Association for Computational Linguistics.

William F. Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, and James Pustejovsky. 2014. Temporal annotation in the clinical domain. *Transactions of the Association for Computational Linguistics*, 2:143–154.

Yukun Tan, Merve Dede, and Ken Chen. 2024. Kclab at chemotimelines 2024: End-to-end system for chemotherapy timeline extraction - subtask2. In *Proceedings of the 6th Clinical NLP Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. Zephyr: Direct distillation of lm alignment.

Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2024a. Openchat: Advancing open-source language models with mixed-quality data. In *The Twelfth International Conference on Learning Representations*.

Liwei Wang, Qiuhao Lu, Rui Li, Sunyang Fu, and Hongfang Liu. 2024b. Wonder at chemotimelines 2024: Medtimeline: An end-to-end nlp system for timeline extraction from clinical narratives. In *Proceedings of the 6th Clinical NLP Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Jeremy L. Warner, Dmitry Dymshyts, Christian G. Reich, Michael J. Gurley, Harry Hochheiser, Zachary H. Moldwin, Rimma Belenkaya, Andrew E. Williams, and Peter C. Yang. 2019. HemOnc: A new standard vocabulary for chemotherapy regimen representation in the OMOP common data model. 96:103239.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.

Kristin Wright-Bettner, Chen Lin, Timothy Miller, Steven Bethard, Dmitriy Dligach, Martha Palmer, James H. Martin, and Guergana Savova. 2020. Defining and learning refined temporal relations in the clinical narrative. In *Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis*, pages 104–114, Online. Association for Computational Linguistics.

Urchade Zaratiana, Nadi Tomeh, Pierre Holat, and Thierry Charnois. 2023. Gliner: Generalist model for named entity recognition using bidirectional transformer.

Xingmeng Zhao and Anthony Rios. 2024. Utsa-nlp at chemotimelines 2024: Evaluating instruction-tuned language models for temporal relation extraction. In *Proceedings of the 6th Clinical NLP Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

## 9   Appendix

| Submission (Breast) | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission1 | 0.97 | 0.94 | 0.96 |
| LAILab_submission3 | 0.97 | 0.94 | 0.95 |
| LAILab_submission2 | 0.97 | 0.94 | 0.95 |
| Baseline_subtask1 | 0.95 | 0.91 | 0.93 |
| Wonder_submission1 | 0.94 | 0.87 | 0.90 |
| Wonder_submission2 | 0.93 | 0.87 | 0.90 |
| Wonder_submission3 | 0.93 | 0.87 | 0.90 |
| BioCom_submission1 | 0.92 | 0.85 | 0.88 |
| KCLab_submission1 | 0.84 | 0.68 | 0.76 |
| NLPeers_submission1 | 0.79 | 0.66 | 0.72 |
| UTSA-NLP_submission1 | 0.79 | 0.60 | 0.70 |
| UTSA-NLP_submission3 | 0.79 | 0.60 | 0.69 |
| UTSA-NLP_submission2 | 0.79 | 0.59 | 0.69 |
| Lexicans_submission1 | 0.78 | 0.58 | 0.68 |
| Lexicans_submission2 | 0.77 | 0.55 | 0.66 |
| Lexicans_submission3 | 0.74 | 0.49 | 0.62 |
| NLPeers_submission2 | 0.63 | 0.34 | 0.49 |
| ClinicalRXMiners_submission1 | 0.49 | 0.39 | 0.44 |
| ClinicalRXMiners_submission2 | 0.49 | 0.18 | 0.33 |

| Submission (Melanoma) | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission1 | 0.93 | 0.81 | 0.87 |
| Baseline_subtask1 | 0.92 | 0.81 | 0.87 |
| LAILab_submission2 | 0.91 | 0.79 | 0.85 |
| NLPeers_submission1 | 0.91 | 0.78 | 0.84 |
| Wonder_submission2 | 0.91 | 0.78 | 0.84 |
| Wonder_submission1 | 0.91 | 0.78 | 0.84 |
| LAILab_submission3 | 0.91 | 0.77 | 0.84 |
| Lexicans_submission1 | 0.90 | 0.76 | 0.83 |
| NLPeers_submission2 | 0.89 | 0.73 | 0.81 |
| Lexicans_submission2 | 0.88 | 0.71 | 0.80 |
| Wonder_submission3 | 0.86 | 0.65 | 0.76 |
| UTSA-NLP_submission1 | 0.82 | 0.55 | 0.68 |
| UTSA-NLP_submission3 | 0.82 | 0.54 | 0.68 |
| UTSA-NLP_submission2 | 0.80 | 0.51 | 0.65 |
| BioCom_submission1 | 0.78 | 0.45 | 0.61 |
| KCLab_submission1 | 0.77 | 0.42 | 0.60 |
| Lexicans_submission3 | 0.77 | 0.42 | 0.59 |
| ClinicalRXMiners_submission2 | 0.70 | 0.24 | 0.47 |
| ClinicalRXMiners_submission1 | 0.67 | 0.17 | 0.42 |

| Submission (Ovarian) | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission3 | 0.93 | 0.86 | 0.89 |
| LAILab_submission2 | 0.93 | 0.85 | 0.89 |
| LAILab_submission1 | 0.92 | 0.84 | 0.88 |
| Baseline_subtask1 | 0.92 | 0.83 | 0.88 |
| Wonder_submission2 | 0.84 | 0.69 | 0.77 |
| Wonder_submission3 | 0.83 | 0.67 | 0.75 |
| NLPeers_submission1 | 0.83 | 0.66 | 0.75 |
| Wonder_submission1 | 0.83 | 0.66 | 0.74 |
| BioCom_submission1 | 0.82 | 0.63 | 0.72 |
| UTSA-NLP_submission2 | 0.80 | 0.59 | 0.70 |
| UTSA-NLP_submission3 | 0.80 | 0.59 | 0.70 |
| UTSA-NLP_submission1 | 0.79 | 0.58 | 0.69 |
| NLPeers_submission2 | 0.75 | 0.50 | 0.63 |
| Lexicans_submission3 | 0.74 | 0.49 | 0.62 |
| Lexicans_submission1 | 0.74 | 0.48 | 0.61 |
| Lexicans_submission2 | 0.73 | 0.46 | 0.59 |
| KCLab_submission1 | 0.68 | 0.37 | 0.53 |
| ClinicalRXMiners_submission2 | 0.48 | 0.21 | 0.34 |
| ClinicalRXMiners_submission1 | 0.39 | 0.27 | 0.33 |

Table 6: Subtask 1, per type of cancer (test set). All scores are macro-F1 of relaxed-to-month setting. We compute two types of scores: F1 with patients with no gold timelines (Type A) and F1 without patients with no gold timelines (Type B). Official score is the average of Type A and Type B, which is used for the rankings in the leader board. See details in section 3

| Submission (Breast) | Type A | Type B | Official Score |
|---|---|---|---|
| KCLab_submission1 | 0.71 | 0.65 | 0.68 |
| Wonder_submission2 | 0.70 | 0.57 | 0.64 |
| Wonder_submission1 | 0.70 | 0.57 | 0.63 |
| Wonder_submission3 | 0.69 | 0.57 | 0.63 |
| LAILab_submission2 | 0.68 | 0.55 | 0.62 |
| Baseline_subtask2 | 0.61 | 0.57 | 0.59 |
| LAILab_submission3 | 0.47 | 0.58 | 0.53 |
| LAILab_submission1 | 0.54 | 0.49 | 0.52 |
| UTSA-NLP_submission1 | 0.32 | 0.18 | 0.25 |
| NYULangone_submission1 | 0.17 | 0.21 | 0.19 |

| Submission (Melanoma) | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission2 | 0.78 | 0.70 | 0.74 |
| LAILab_submission1 | 0.68 | 0.45 | 0.57 |
| KCLab_submission1 | 0.64 | 0.35 | 0.49 |
| Baseline_subtask2 | 0.60 | 0.26 | 0.43 |
| Wonder_submission3 | 0.37 | 0.42 | 0.39 |
| Wonder_submission1 | 0.37 | 0.42 | 0.39 |
| Wonder_submission2 | 0.37 | 0.41 | 0.39 |
| LAILab_submission3 | 0.43 | 0.33 | 0.38 |
| NYULangone_submission1 | 0.40 | 0.25 | 0.32 |
| UTSA-NLP_submission1 | 0.12 | 0.30 | 0.21 |

| Submission (Ovarian) | Type A | Type B | Official Score |
|---|---|---|---|
| LAILab_submission2 | 0.83 | 0.65 | 0.74 |
| Baseline_subtask2 | 0.80 | 0.61 | 0.71 |
| LAILab_submission1 | 0.73 | 0.46 | 0.59 |
| Wonder_submission3 | 0.70 | 0.40 | 0.55 |
| Wonder_submission2 | 0.70 | 0.39 | 0.55 |
| Wonder_submission1 | 0.69 | 0.38 | 0.53 |
| LAILab_submission3 | 0.49 | 0.49 | 0.49 |
| KCLab_submission1 | 0.55 | 0.35 | 0.45 |
| UTSA-NLP_submission1 | 0.21 | 0.17 | 0.19 |
| NYULangone_submission1 | 0.21 | 0.16 | 0.18 |

Table 7: Subtask 2, results for each type of cancer (test set). All scores are macro-F1 of relaxed-to-month setting. We compute two types of scores: F1 with patients with no gold timelines (Type A) and F1 without patients with no gold timelines (Type B). Official score is the average of Type A and Type B, which is used for the rankings in the leader board. See details in section 3

# IryōNLP* at MEDIQA-CORR 2024:
# Tackling the Medical Error Detection & Correction Task
# On the Shoulders of Medical Agents

**Jean-Philippe Corbeil**
Microsoft Health & Life Sciences
jcorbeil@microsoft.com

## Abstract

In natural language processing applied to the clinical domain, utilizing large language models has emerged as a promising avenue for error detection and correction on clinical notes, a knowledge-intensive task for which annotated data is scarce. This paper presents MedReAct'N'MedReFlex, which leverages a suite of four LLM-based medical agents. The MedReAct agent initiates the process by observing, analyzing, and taking action, generating trajectories to guide the search to target a potential error in the clinical notes. Subsequently, the MedEval agent employs five evaluators to assess the targeted error and the proposed correction. In cases where MedReAct's actions prove insufficient, the MedReFlex agent intervenes, engaging in reflective analysis and proposing alternative strategies. Finally, the MedFinalParser agent formats the final output, preserving the original style while ensuring the integrity of the error correction process. One core component of our method is our RAG pipeline based on our ClinicalCorp corpora. Among other well-known sources containing clinical guidelines and information, we preprocess and release the open-source MedWiki dataset for clinical RAG application. Our results demonstrate the central role of our RAG approach with ClinicalCorp leveraged through the MedReAct'N'MedReFlex framework. It achieved the ninth rank on the MEDIQA-CORR 2024 final leaderboard.

## 1 Introduction

In natural language processing applied to the clinical domain, the accurate detection and correction of medical errors are paramount tasks with profound implications for patient care and safety. This paper introduces the multi-agent framework MedReAct'N'MedReFlex, meticulously engineered to tackle medical error detection and correction, as delineated in the MEDIQA-CORR 2024 competition.

Our framework integrates four distinct types of medical agents: MedReAct, MedReFlex, MedEval, and MedFinalParser, each playing a specialized role in the error identification and rectification process. Drawing inspiration from existing frameworks like ReAct (Yao et al., 2023) and Reflexion (Shinn et al., 2023), our framework orchestrates a structured approach to error handling.

Leveraging a Retrieval-Augmented Generation (RAG) framework (Lewis et al., 2020) based on MedRAG (Xiong et al., 2024) and MedCPT (Jin et al., 2023), our approach operates over Clinical-Corp, an extensive corpora curated to encompass crucial clinical guidelines. Additionally, we introduce *MedWiki*, a collection of medical articles from Wikipedia. By integrating these resources, our approach seeks to advance state-of-the-art clinical NLP by offering a comprehensive solution tailored to the intricate nuances of medical error handling. Furthermore, this paper documents the construction and release of *MedWiki*, a substantial repository comprising over 1.3 million article chunks. Additionally, we detail the assembly of the ClinicalCorp, a comprehensive corpus comprising *MedWiki* along with other clinical guideline datasets, such as parts of the MedCorp corpora (Xiong et al., 2024) and parts of the guidelines (Chen et al., 2023).

Our main contributions are:

- We designed a multi-agent framework named *MedReAct'N'MedReFlex* to solve the medical error detection & correction task (MEDIQA-CORR 2024) based on four types of medical agents: *MedReAct*, *MedReFlex*, *MedEval* and *MedFinalParser*. We deployed this framework on ClinicalCorp using a retrieval-augmented generation approach.

---

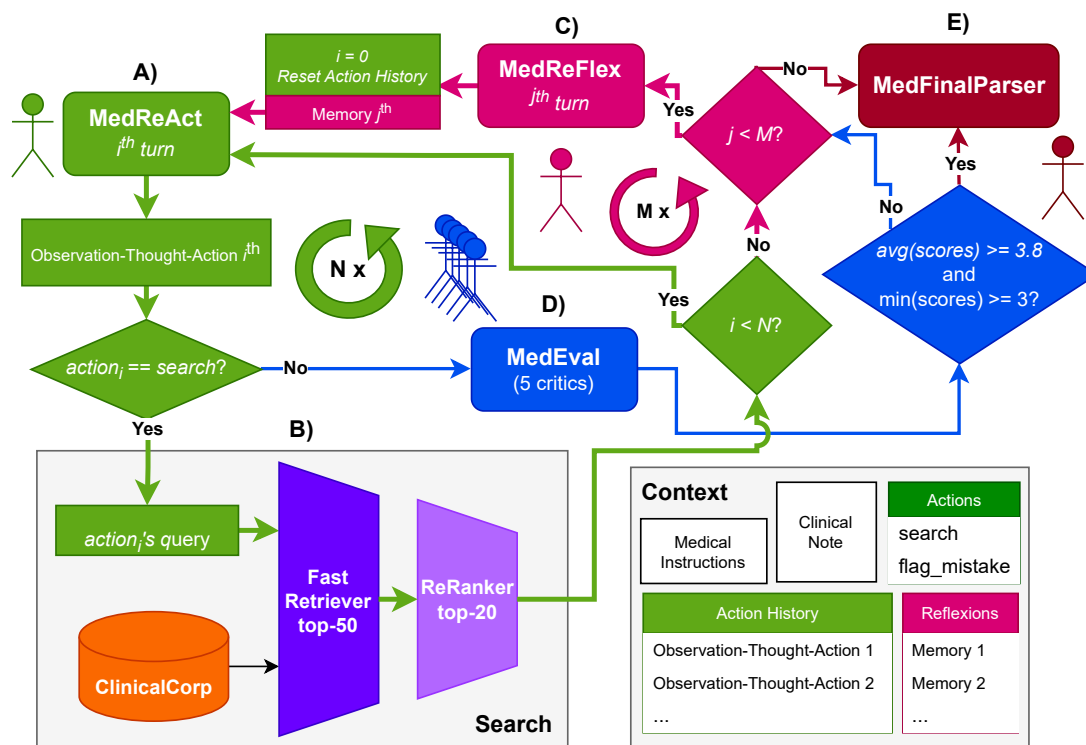*The team name *iryō* comes from the japanese for medical or healthcare.

Figure 1: Schema of *MedReAct'N'MedReFlex* along the context of the clinical error correction task accessible to all medical agents: *MedReAct*, *MedReFlex*, *MedEval* and *MedFinalParser*. A) The *MedReAct* agent first provides an observation, a thought and an action. B) In the case of a *search* action, it triggers a semantic search over ClinicalCorp using *MedReAct*'s query. Then, the MedReAct agent loops up to $N$ times (green inner loop) or until a *final_mistake* action is provided. C) After $N$ unsuccessful searches from *MedReAct*, the *MedReFlex* agent reflects on the current situation and suggests a solution (pink outer loop). Then, *MedReAct* might start again. D) Once *MedReAct* selects the *final_mistake* action, the five *MedEval* agents review the answer and give a score between 1 and 5 (blue line). E) If the average equals or surpasses 3.8 and the minimum above or equal to 3, the *MedFinalParser* agent formats the final answer into a JSON object. If the answer is unsatisfactory, *MedReFlex* is triggered instead. If *MedReFlex* reaches unsuccessfully the $M^{th}$ turns, *MedFinalParser* concludes that there is no error.

- We released the open-source *MedWiki*[1], a version of Wikipedia 2022-12-22 focused solely on medical articles. This RAG-ready dataset contains about 1.3M chunks from more than 150K articles, which represents about 3% of the original corpus.

- We provided the recipe to assemble our large corpora *ClinicalCorp* for RAG applications in the clinical domain, containing more than 2.3M chunks.

- We released a RAG-ready version[2] of the open-source *guidelines* used to pre-train *Meditron* (Chen et al., 2023), containing more than 710K chunks across eight open-source datasets.

- We released our codebase on GitHub[3].

## 2 Related Work

### 2.1 Medical Large Language Models

Since the emergence of ChatGPT by OpenAI in December 2022, the landscape of large language models (LLMs) has witnessed a proliferation of both private and public initiatives, leading to the development of increasingly sophisticated models. OpenAI's journey from the GPT3.5-turbo architecture, as reported by Brown et al. and Ouyang et al., culminated in the release of GPT-4 and its turbo variant (Achiam et al., 2023). Similarly, Google introduced Gemini, available in Nano, Pro, and Ultra configurations (Team et al., 2023), alongside its open-source Gemma model (Team et al., 2024). Anthropic contributed to this landscape with

---

Claude3, offered in three sizes, ranging from Haiku to Opus. Other notable LLMs include Mistral and Mixtral (Jiang et al., 2023, 2024), as well as Llama 2 (Touvron et al., 2023) and Yi (Young et al., 2024). These general-purpose LLMs, such as GPT-4, have demonstrated solid in-context learning capabilities in the medical field Nori et al. (2023).

Researchers have developed various open-source LLMs with diverse capabilities in the medical NLP domain. Examples include ClinicalCamel (Toma et al., 2023), Med42 (Christophe et al., 2023), PMC-Llama (Wu et al., 2023a), BioMedGPT (Zhang et al., 2023), Meditron (Chen et al., 2023), Apollo (Wang et al., 2024), OpenMedLM (Garikipati et al., 2024), and BioMistral (Labrak et al., 2024). Google also contributed Med-PaLM 2, a specialized LLM tailored for medical tasks (Singhal et al., 2023).

In this study, we employed OpenAI's GPT-4, specifically version turbo 0125, due to its proven state-of-the-art performances in various domains, its functional capabilities, and its large context window of 128K tokens. These attributes make it an ideal foundation for our approach. For instance, Nori et al. (2023) demonstrated that utilizing in-context learning with GPT-4 — relying on prompt engineering (i.e. few-shot learning (Brown et al.), chain-of-thought (Wei et al., 2022; Kojima et al., 2022), self-consistency (Wang et al., 2022) and shuffling multiple choice (Ko et al., 2020)) — achieves state-of-the-art performances on medical question-answering tasks, surpassing specialized models like Med-PaLM 2. We relied on a similar approach as our early baseline for medical error detection and correction, discarding the self-consistency and the shuffling techniques since both do not apply to generative tasks. Nonetheless, we have observed low results from which we hypothesized that this approach using only parametric knowledge is lacking reliable knowledge (Mallen et al., 2023; Ovadia et al., 2023; Kandpal et al., 2023), which we addressed by applying agentic methods in a retrieval-augmented generation framework.

## 2.2 Agentic Methods

Researchers have devised several agentic methods to enhance LLMs' responses and reasoning capabilities, such as ReAct (Yao et al., 2023), Reflexion (Shinn et al., 2023), DSPy (Khattab et al., 2023) and self-discovery (Zhou et al., 2024). Additionally, multi-agent paradigms (Wu et al., 2023b)

have found application in the medical domain (Tang et al., 2023). Our approach draws inspiration from the Reflexion framework (Shinn et al., 2023), which we adapted into our *MedReFlex* agent. Specifically, we implemented a *MedReAct* agent — inspired by the ReAct approach (Yao et al., 2023) — to generate trajectories in our environment. However, this agent realizes its sequence of actions in a different order (i.e., observation, thought, and action), enabling streamlined execution.

Given the reliance of the Reflexion framework on feedback mechanisms, we incorporated an LLM-based metric into our *MedEval* medical agents. Evaluation metrics based on prompting strong LLMs, such as GPT-4 (Liu et al., 2023), have demonstrated a high correlation with human judgment. Similar findings have been reported in the medical NLP literature (Xie et al., 2023). Our evaluation protocol involves prompting five GPT-4 reviewers with task-specific criteria: validity, preciseness, confidence, relevance, and completeness. The average and minimum of their scores are both utilized as success criteria, capturing an unbiased final score and the evaluators' confidence, respectively.

## 2.3 Retrieval-Augmented Generation

Before the advent of LLMs, authors have proposed the retrieval-augmented generation (RAG) framework as a mechanism to incorporate non-parametric memory for knowledge-intensive tasks. This framework, as elucidated by Lewis et al. (Lewis et al., 2020), leverages both sparse (Robertson et al., 2009) and dense (Reimers and Gurevych, 2019) retrieval methods. In the medical NLP domain, MedCPT (Jin et al., 2023) serves as a prominent retrieval approach, augmented by a reranking stage based on a cross-encoder model. Notably, Xiong et al. (Xiong et al., 2024) conducted a comprehensive study on RAG applications in the medical domain, culminating in developing the MedRAG framework and the MedCorp corpora. Our approach builds upon these foundations, employing the MedCPT retrieval techniques and two corpora from MedCorp.

A pivotal aspect of RAG is its search engine's collection of indexed documents. The *guidelines* corpora, part of the *GAP-replay* corpora, was curated to train Meditron (Chen et al., 2023). This corpus comprises web pages describing medical guidelines from reputable healthcare websites like the World Health Organization. The *StatPearls* and

*Textbooks* datasets, included in the *MedCorp* corpora used in MedRAG (Xiong et al., 2024), encompass documents from clinical decision support tools and medical textbooks (Jin et al., 2021). While *Wikipedia* and *PubMed* datasets within *MedCorp* offer extensive data (i.e. more than 55M documents), we opted for efficiency by focusing on the smaller *PubMed* subset in the *guidelines* corpora and our *MedWiki* corpus.

## 3 Methodology

### 3.1 MEDIQA-CORR Task

The goal of the medical error detection and correction task (Ben Abacha et al., 2024a) from the clinical note is threefold: detect the presence of an error, locate the sentence containing the error and generate a corrected version of that sentence. The input of the dataset (Ben Abacha et al., 2024b) is a clinical note of several sentences containing a medical description of a patient's condition, test results, diagnosis, treatment and other aspects. There are two parts for the validation and test sets: *MS* from Microsoft and *UW* from the University of Washington. As a primary evaluation metric, the organizers asked to utilize the aggregation score defined by Abacha et al. (2023) over Rouge-1, BertScore and BLEURT, demonstrating a higher correlation with human judgement.

### 3.2 ClinicalCorp Corpora

Our corpus is detailed in Table 1.

**guidelines** We aggregated 13 datasets — which are open-source or closed-source — from the *guidelines* corpora. We adapted and ran the scrappers from the Meditron GitHub repository to gather the closed-source datasets. Then, we chunked the resulting documents using LangChain's recursive-character text splitter (Chase, 2022) with a chunk size of 1,000 characters and an overlap of 200 characters, as used for *StatPearls* (see next section).

**MedCorp** We gathered two of the four datasets contained in *MedCorp* from MedRAG (Xiong et al., 2024): *StatPearls* and *Textbooks*. The former was downloaded, cleaned and chunked using *MedRAG* GitHub repository, while the latter was readily available on the *HuggingFaceHub*[4].

**MedWiki** We filtered the 2022-12-22 Wikipedia dump[5] pre-processed into chunks by *Cohere* for

medical articles only. To select the medical articles, we leveraged an available fine-tuned BerTopic[6] (Grootendorst, 2022), trained on the same Wikipedia dump. We associated its 2,3K topics to the medical domain based on the topics' word representations — e.g. topic *1850* is related to the medical field, and it corresponds to the word representations: shingles, herpesvirus, chickenpox, herpes, smallpox, zoster, immunity, infectivity, inflammation, and viral. We made these predictions by prompting *GPT3.5-turbo 0613* with a temperature of 1.0 followed by a majority vote over five predictions. If at least four were positive, we declared the topic medically relevant. In the manual verification of about 50 diverse medical terms on the resulting collection, we observed a near-perfect coverage of Wikipedia's articles related to diseases, treatments, bacteria, or drugs. Only two topics were missing[7], corresponding to one single example from the manual test. Given that our goal is to reduce the size of this dataset and use it in an RAG application, we added these topics manually. We obtained a corpus of 150K articles and nearly 1.4M chunks.

### 3.3 Semantic Search

We followed the MedCPT approach (Jin et al., 2023) in two stages (see step *B* in Figure 1), which is composed of a fast bi-encoder retrieving stage followed by a cross-encoder reranking stage.

We implemented the first stage on a ChromaDB instance, in which we loaded *ClinicalCorp*. This stage aims to find relevant documents while maintaining a good accuracy/latency trade-off. This vector database embeds documents using a fast bi-encoder model (Reimers and Gurevych, 2019). Then, we provide a query to fetch the closest documents under a given distance, computed with the hierarchical navigable small world approximation (HNSW, by Malkov and Yashunin (2018)). We experimented with three bi-encoders from the *HuggingFaceHub*: *sentence-transformers/all-MiniLM-L6-v2* (default), *NeuML/pubmedbert-base-embeddings-matryoshka* and MedCPT original Query/Article encoders. According to our initial experiments, we discarded *all-MiniLM-L6-v2* because we noticed a critical lack of knowledge about medical terminology hindering its accuracy despite a very low latency. NeuML's model and MedCPT's

---

[4]`hf.co/datasets/MedRAG/textbooks`
[5]`hf.co/datasets/Cohere/wikipedia-22-12`

[6]`hf.co/MaartenGr/BERTopic_Wikipedia`
[7]Index *509* related to biological taxonomy and *806* related to yeasts.

Table 1: Datasets gathered to construct ClinicalCorp.

| Dataset | Source | Status | # Documents | # Chunks |
|---|---|---|---|---|
| Guidelines (Chen et al., 2023) | WikiDoc | open | 33,058 | 360,070 |
| | PubMed (guidelines only) | open | 1,627 | 124,971 |
| | National Institute for Health and Care Excellence | open | 1,656 | 87,904 |
| | Center for Disease Control and Prevention | open | 621 | 70,968 |
| | World Health Organization | open | 223 | 33,917 |
| | Canadian Medical Association | open | 431 | 18,757 |
| | Strategy for Patient-Oriented Research | open | 217 | 11,955 |
| | Cancer Care Ontario | open | 87 | 2,203 |
| | Drugs.com | close | 6,711 | 37,255 |
| | GuidelineCentral | close | 1,285 | 2,451 |
| | American Academy of Family Physicians | close | 60 | 130 |
| | Infectious Diseases Society of America | close | 54 | 7,785 |
| | Canadian Paediatric Society | close | 43 | 1,123 |
| MedCorp (Xiong et al., 2024) | StatPearls | close | 9,379 | 307,187 |
| | Textbooks (Jin et al., 2021) | open | 18 | 125,847 |
| **ClinicalCorp** (Ours) | MedWiki | open | 150,380 | 1,139,464 |
| | **All** | **mix** | **205,850** | **2,331,987** |

are Bert-based models of 768 hidden dimensions and 12 layers, a slow architecture to generate sentence embeddings. However, NeuML fine-tuned a recent model using the Matryoshka Representation Learning technique (Kusupati et al., 2022), allowing to truncate dimensions down to 256 dimensions of the 768 embeddings, which significantly accelerated the computations. Our experiments employ this MRL encoder with truncation at 256 dimensions as a trade-off between accuracy and latency.

We implemented the reranking stage following the cross-encoder approach from MedCPT (Jin et al., 2023). Our early experimentation demonstrated the superiority of this model compared to NeuML's MRL bi-encoder with all 768 dimensions as a reranker.

# 4 MedReAct'N'MedReFlex Framework

Unlike previous multi-agent frameworks (Wu et al., 2023b; Tang et al., 2023), our approach diverges from a free conversation format to adopt a fixed design schema, as illustrated in Figure 1. Within this structured framework, each medical agent intervenes at a specific step, facilitating a systematic

and coordinated approach to address the error detection and correction task. Central to our methodology are four distinct medical agents: MedReAct, MedReFlex, MedEval, and MedFinalParser.

## 4.1 MedReAct Agent

The MedReAct agent (see step *A* in Figure 1), inspired by the ReAct framework (Yao et al., 2023), operates cyclically, beginning with an observation of the current context, followed by a thoughtful analysis, and concluding with an action (*search* or *final_mistake*). This agent generates a trajectory of up to $N$ steps if the action is always a *search* with different queries.

We also experimented with adding two other actions (*get_doc_by_id* and *next_results_from_query*), but MedReAct systematically ignored them.

## 4.2 MedEval Agent

Upon MedReAct's selection of the *final_mistake* action, the MedEval agents (see step *D* in Figure 1), akin to the GPT-Eval approach (Liu et al., 2023), evaluate the proposed solution. Five GPT-4-based evaluators assess the answer based on criteria such

as validity, preciseness, confidence, relevance, and completeness. The ensemble of evaluators ensures comprehensive and unbiased feedback, contributing to robust error detection and correction. We leverage the average final score as well as the minimum review score. We added this condition on the minimum score to capture the confidence of the evaluation. If one reviewer gave a much lower score than the others, we experimentally observed that it was often a signal of lower confidence in the final answer.

### 4.3 MedReFlex Agent

In scenarios where MedReAct's actions fail to yield satisfactory outcomes, the MedReFlex agent (see step $C$ in Figure 1), drawing from the Reflexion framework (Shinn et al., 2023), intervenes. This agent engages in reflective analysis to reassess the situation. By considering contextual cues, past interactions and all five reviews, MedReFlex proposes alternative strategies to address the identified challenges. This iterative process allows for adaptive decision-making and fosters resilience in error detection and correction tasks.

### 4.4 MedFinalParser Agent

Suppose the average score provided by the MedEval agents exceeds or equals 4, and the minimum score surpasses or equals 3. In that case, the MedFinalParser agent (see step $D$ in Figure 1) proceeds to format the final answer into a JSON object. This agent also ensures the conservation of the original style of the clinical note, which the MedReAct agent tends to disrupt by copying the writing style of the search documents. Conversely, if the answer falls short of the predetermined thresholds, MedReFlex is triggered for further refinement. If MedReFlex's interventions prove ineffective after the $M^{th}$ turn, the MedFinalParser agent concludes that no errors exist, ensuring the integrity of the error correction process.

## 5 Results

### 5.1 Results for the Competition

MedReAct'N'MedReFlex achieved the $9^{th}$ rank during the MEDIQA-CORR 2024 official testing period, corresponding to an aggregation score of 0.581. Nonetheless, we thoroughly optimize our method in the following sections. To complete these experiments in a reasonable amount of time,

we randomly sample 50 examples from the MS validation set.

### 5.2 Agentic Method Comparison

In Table 2, we compared the MedReAct agent only against using our proposed method MedReAct'N'MedReFlex. Our approach achieves more than a few absolute percent across metrics. We also experimented with a baseline inspired from Nori et al. (2023) (i.e. "MedPrompt") with in-context learning prompting alone, but the results were drastically lower.

| Metric | MedReAct | MedReAct'N'MedReFlex |
|--------|----------|----------------------|
| ROUGE-1 | 0.504 | 0.568 |
| BERTScore | 0.580 | 0.642 |
| BLEURT | 0.531 | 0.588 |
| Aggregate | 0.539 | 0.599 |

Table 2: Comparison between MedReAct agent only with up to 10 turns against MedReAct'N'MedReFlex with 4 turns for MedReAct and 5 turns for MedReFlex, leveraging the optimal search configuration (retrieval top-k at 50 and reranking top-k at 20).

### 5.3 Semantic Search Optimization

After the end of the MEDIQA-CORR 2024 shared task, we carried out a thorough analysis of our semantic search engine. The main parameters to tune are retrieval top-k, reranking top-k and the source included in ClinicalCorp.

#### 5.3.1 Retrieval Top-K

In Figure 2, we illustrate the performances across many retrieval top-k values employing a fixed reranking top-k of 20. For the official ranking of the MEDIQA-CORR 2024, we set this value to 300. However, we observe here that this setting is sub-optimal. A retrieval top-k of 50 improves the final performances by a few absolute percent. We interpret this observation as indicative of a misalignment between our task and the fine-tuning of the MedCPT reranker. The more documents we provide to the reranking model (e.g. 200 or 300), the more low-relevance documents are put in the top 20 by the reranker output.

Nonetheless, a reranking without surplus documents — i.e. retrieval top-k of 20 with a reranking top-k of 20 — remains sub-optimal, mainly in contrast to using 50 documents. In Figure 3, we provide the associated average latency for one react step in seconds. We notice that the latency

seems to scale with the order of magnitude of the retrieval top-k, with a value of 20 and 50 having 17 seconds on average, while 100, 200 and 300 are around 20 seconds. We expected that the reranking of 300 examples against 100, for instance, would lead to noticeable latency, but it is negligible in contrast to the retrieval from ChromaDB over our 2.3M chunks.
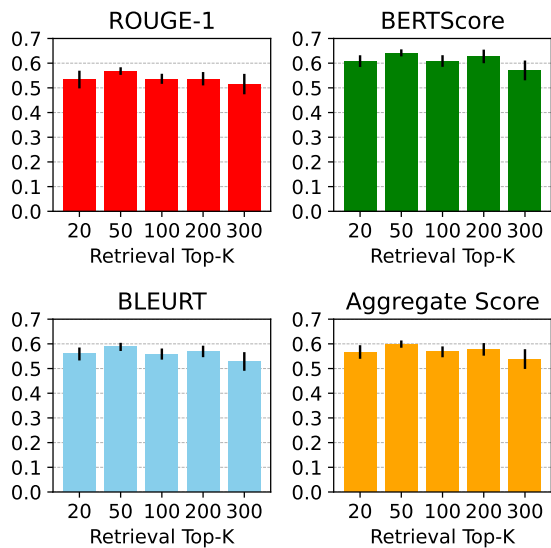


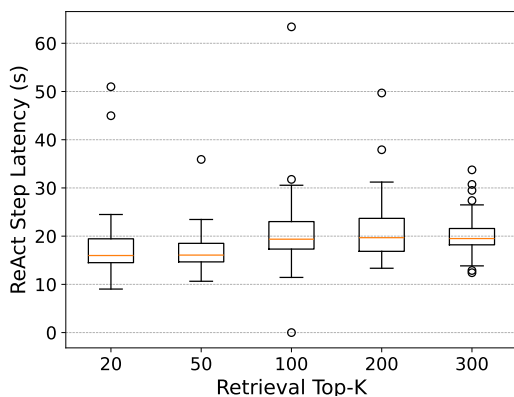Figure 2: Performances across many retrieval top-k values with a reranking top-k set at 20 over 3 runs.



Figure 3: ReAct step average latency per retrieval top-k with a reranking top-k set at 20.

We also show the average amount of MedReFlex turns and the average of the sum of all MedReAct turns in Figure 4. Overall, the trends are similar, with 4.8 total ReAct turns on average, but there is a slight increase in the average and variance for top-k values of 200 and 300. Therefore, these settings are underperforming and slower regarding latency



Figure 4: Average turns of MedReAct and MedReFlex according to various retrieval top-k with a reranking top-k set at 20.

and the number of turns needed to reach an answer.

Overall, the retrieval top-k of 50 leads to higher performances across all metrics and reduced latency and number of turns required by our algorithm.

### 5.3.2 Reranker Top-K

In Figure 5, we fix the retrieval top-k at 300 and compute the performances across three reranking top-k values: 5, 10 and 20. Since the context window of the LLM limits us, we constraint the maximum of $K$ to 20, given that these $K$ documents are injected in the prompt up to $N$ times for each MedReAct step. According to Figure 5, we observe that the more documents we provide in the prompt, the more we increase the performances — the aggregate score gains close to 10% absolute when augmenting from 5 to 20 documents.
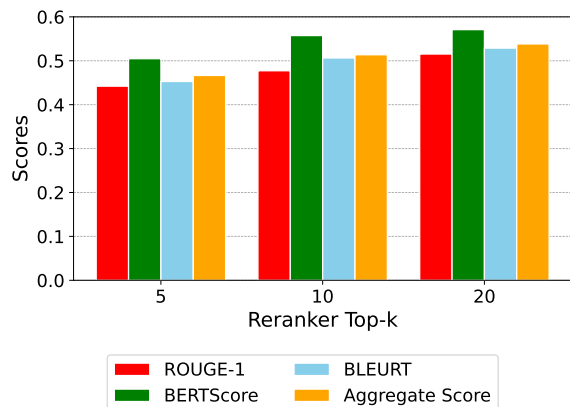


Figure 5: Reranker top-K with a retrieval top-k set at 300.

576

### 5.3.3 Sources in ClinicalCorp

We measure the impact of each source in Clinical-Corp in Figure 6. First, we observe that *MedWiki* is the lowest-performing source of documents with an aggregation score of nearly 0.47. *guidelines* and *Textbooks* provide a similar accuracy at about 0.51 in aggregate score. Finally, *StatPearls* leads to the highest score close to the full ClinicalCorp. Given our small validation set of 50 examples, we consider it a better practice to keep all ClinicalCorp for our task since more edge cases might appear at test time.



Figure 6: Performances per source in ClinicalCorp with the retrieval top-k set at 50 and the reranking top-k set at 20.

We show in Figure 7 the distributions of sources from ClinicalCorp in general in comparison to the distributions of sources' chunks used by one run of MedReAct'N'MedReFlex. We observe, in general, a much larger utilization of the *StatPearls'* chunks in contrast to the *MedWiki*'s chunks, while we remark similar distributions for the other two datasets. These results align with the previous analysis demonstrating a higher performance from using only *StatPearls*.

### 5.4 MedEval Aggregation Thresholds

In Figure 8, we show the impact of applying different thresholds to the average and minimum review scores on the performance. For the minimum score criterion, we choose the integer values of 2.0, 3.0 and 4.0. We select values for the average score criterion: 3.0, 3.2, 3.5, 3.8, 4.0 and 4.2. We do not compute the performances for combinations where the minimum threshold is higher than the average threshold for mathematical consistency. We observe an optimal setting for a minimum evaluation score of 3.0 with a range of average evaluation scores in [3.5, 3.8].
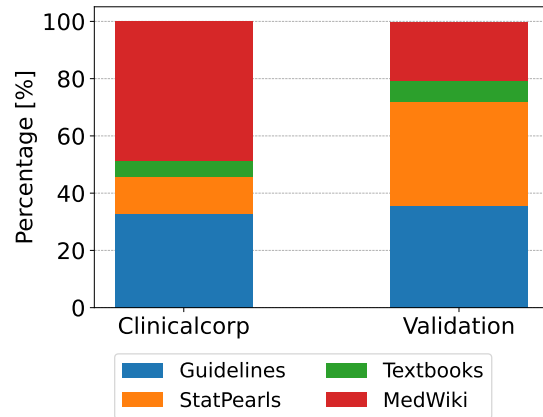


Figure 7: Distribution of sources' chunks in Clinical-Corp against appearances of these chunks' sources in one run of MedReAct'N'MedReFlex.
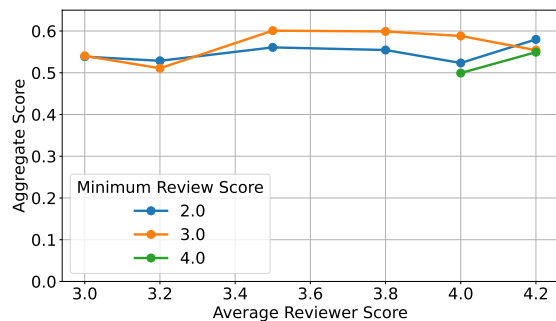


Figure 8: Aggregate scores across different MedEval's average and minimum thresholds with the retrieval top-k set at 50 and the re-ranker top-k set at 20. We omitted average thresholds that are strictly lower for consistency for a given minimum threshold.

## 6 Conclusion

In this paper, we introduced MedReAct'N'MedReFlex, a multi-agent framework developed for the MEDIQA-CORR 2024 competition aimed at medical error detection and correction in clinical notes. The framework incorporates four specialized medical agents: MedReAct, MedReFlex, MedEval, and MedFinal-Parser, leveraging the RAG framework and our ClinicalCorp. We detail the construction of our ClinicalCorp, including diverse clinical datasets such as *guidelines*, *Textbooks*, and *StatPearls*. Additionally, we released MedWiki, a corpus comprising Wikipedia medical articles. Our framework achieved the ninth rank in the competition with an aggregation score of 0.581. Through optimization experiments, we identified sub-optimal settings at the time, demonstrating substantial performance

improvements with a retrieval top-k of 50, a reranking top-k of 20, an average review threshold of 3.8, and a minimum review threshold of 3. As future work, we envision refining the chunking strategy on the ClinicalCorp, applying further prompt engineering of the medical agents, and conducting a deeper analysis of the interactions between the MedReAct'N'MedReFlex's agents.

## Acknowledgments

## References

Asma Ben Abacha, Wen-wai Yim, Griffin Adams, Neal Snider, and Meliha Yetisgen-Yildiz. 2023. Overview of the mediqa-chat 2023 shared tasks on the summarization & generation of doctor-patient conversations. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 503–513.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.

Harrison Chase. 2022. Langchain. https://github.com/langchain-ai/langchain. Version 1.2.0, Released on October 17, 2022.

Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. 2023. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*.

Clément Christophe, Avani Gupta, Nasir Hayat, Praveen Kanithi, Ahmed Al-Mahrooqi, Prateek Munjal, Marco Pimentel, Tathagata Raha, Ronnie Rajan, and Shadab Khan. 2023. Med42 - a clinical large language model.

Anurag Garikipati, Jenish Maharjan, Navan Preet Singh, Leo Cyrus, Mayank Sharma, Madalina Ciobanu, Gina Barnes, Qingqing Mao, and Ritankar Das. 2024. Openmedlm: Prompt engineering can out-perform fine-tuning in medical question-answering with open-source large language models. In *AAAI 2024 Spring Symposium on Clinical Foundation Models*.

Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, Lana Yeganova, W John Wilbur, and Zhiyong Lu. 2023. Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval. *Bioinformatics*, 39(11):btad651.

Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pages 15696–15707. PMLR.

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, Heather Miller, et al. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*.

Miyoung Ko, Jinhyuk Lee, Hyunjae Kim, Gangwoo Kim, and Jaewoo Kang. 2020. Look at the first sentence: Position bias in question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1109–1121, Online. Association for Computational Linguistics.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-taka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.

Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham M Kakade, Prateek Jain, et al. 2022. Matryoshka representation learning. In *Advances in Neural Information Processing Systems*.

Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. *arXiv preprint arXiv:2402.10373*.

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.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: Nlg evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522.

Yu A Malkov and Dmitry A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):824–836.

Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822.

Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *Medicine*, 84(88.3):77–3.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback.

Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. 2023. Fine-tuning or retrieval? comparing knowledge injection in llms. *arXiv preprint arXiv:2312.05934*.

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.

Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.

Xiangru Tang, Anni Zou, Zhuosheng Zhang, 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*.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Augustin Toma, Patrick R Lawler, Jimmy Ba, Rahul G Krishnan, Barry B Rubin, and Bo Wang. 2023. Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding. *arXiv preprint arXiv:2305.12031*.

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*.

Xidong Wang, Nuo Chen, Junyin Chen, Yan Hu, Yidong Wang, Xiangbo Wu, Anningzhe Gao, Xiang Wan, Haizhou Li, and Benyou Wang. 2024. Apollo: Lightweight multilingual medical llms towards democratizing medical ai to 6b people. *arXiv preprint arXiv:2403.03640*.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain

of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.

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.

Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023a. Pmc-llama: Further fine-tuning llama on medical papers. *arXiv preprint arXiv:2304.14454*.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023b. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.

Yiqing Xie, Sheng Zhang, Hao Cheng, Zelalem Gero, Cliff Wong, Tristan Naumann, and Hoifung Poon. 2023. Enhancing medical text evaluation with gpt-4. *arXiv preprint arXiv:2311.09581*.

Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. *arXiv preprint arXiv:2402.13178*.

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)*.

Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*.

Kai Zhang, Jun Yu, Zhiling Yan, Yixin Liu, Eashan Adhikarla, Sunyang Fu, Xun Chen, Chen Chen, Yuyin Zhou, Xiang Li, et al. 2023. Biomedgpt: A unified and generalist biomedical generative pre-trained transformer for vision, language, and multimodal tasks. *arXiv preprint arXiv:2305.17100*.

Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V Le, Ed H Chi, Denny Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. 2024. Self-discover: Large language models self-compose reasoning structures. *arXiv preprint arXiv:2402.03620*.

# Overview of the MEDIQA-M3G 2024 Shared Task on Multilingual Multimodal Medical Answer Generation

**Wen-wai Yim, Asma Ben Abacha**
Microsoft Health AI
{yimwenwai,abenabacha}@microsoft.com

**Yujuan Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen**
University of Washington
{velvinfu,zhaoyis,fxia,melihay}@uw.edu

**Martin Krallinger**
Barcelona Supercomputing Center
martin.krallinger@bsc.es

## Abstract

Remote patient care provides opportunities for expanding medical access, saving healthcare costs, and offering on-demand convenient services. In the MEDIQA-M3G 2024 Shared Task, researchers explored solutions for the specific task of dermatological consumer health visual question answering, where user generated queries and images are used as input and a free-text answer response is generated as output. In this novel challenge, eight teams with a total of 48 submissions were evaluated across three language test sets. In this work, we provide a summary of the dataset, as well as results and approaches. We hope that the insights learned here will inspire future research directions that can lead to technology that deburdens clinical workload and improves care.

## 1 Introduction

Driven by long patient wait times, high medical costs and physician burnout, remote patient care delivery (e.g. e-visits, e-mails) provides a cost effective solution that lowers facility expenditures, allows flexible schedules for both clinicians and patients, and expands health care access(Bishop et al., 2024). The trend, already in motion due to the maturation of telecommunication technologies and the proliferation of health portals, was massively accelerated by the onset of the global COVID 19 epidemic in 2019. Five years later, today, while remote technologies allow for the conveniences of patient care delivered conveniently from one's own home, this poses new challenges for providers who need to meet the new demand, where patients can request services at any time of the day, creating a perception of "never-ending work"(Sinsky et al., 2024).

Automatic response generation may alleviate doctor burden by providing suggestions when answering patient queries, speeding up response throughput. In this work, we present the MEDIQA

2024 Multilingual & Multimodal Medical Answer Generation (M3G) Shared Task, which is focused on the problem of multimodal answer generation in the space of dermatology, evaluated in multiple languages (English, Chinese, and Spanish). Specifically, a health related question along with one or more images is posed; the expected task is to generate an appropriate answer response.

Previous editions of the MEDIQA shared tasks have featured radiology-related visual question answering(Lau et al., 2018; Ben Abacha et al., 2019) and text-only consumer health answer generation(Ben Abacha et al., 2017). Other prior work in medical VQA includes images in the space of pathology and GI-tract (He et al., 2021; Hicks et al., 2023). Meanwhile, previous work in dermatological image classification focused on image-only input and multi-class classification(Daneshjou et al., 2021; Groh et al., 2021). This is the first shared task to incorporate visual question answering for user generated health queries and images.

## 2 Task

### 2.1 Description

In this task, participants were given textual inputs which may include clinical history and a query, along with one or more associated images. The task objective consisted of generating a relevant textual answer response. An example instance is shown in Table 1.

The training set contained multiple possible gold standard responses. Each response included information related to its author validation level (e.g. real-id verified, medical doctor verified) and a ranking based on their platform contribution from 0-8 levels, the higher the better. English and Spanish translations were automatically generated from original Chinese (in simplified Chinese characters) using GPT4.

In the validation and test sets, each text response

Table 1: Example from the DermaVQA IIYI data subset. The original posts are in Chinese, which are translated into English and Spanish by GPT4 (if they are in the training set) or medical translators (otherwise).

| Query | Responses |
|---|---|
|  帮忙诊断一下:三个月前出现如下图，自己用达克能宁喷雾两个月无明显效果，之后去乡村诊所，医生指导用鸡眼膏，之后出现变红变多，请帮忙诊断下 <br><br> Please help with the diagnosis: Three months ago, the condition shown in the picture below appeared. The patient used Daknening spray for two months without any noticeable effect. Afterwards, they went to a rural clinic, where the doctor advised them to use corn ointment. Subsequently, the condition turned red and worsened. Please help with the diagnosis. <br><br> Por favor, ayude con el diagnóstico : Hace tres meses, apareció la condición mostrada en la imagen de abajo. El paciente utilizó el spray Daknening durante dos meses sin ningún efecto notable. Posteriormente, acudió a una clínica rural, donde el médico le aconsejó que utilizara pomada de maíz. Posteriormente, la condición se volvió roja y empeoró. Por favor, ayude con el diagnóstico. | RESPONSE1: <br> 是鸡眼。 <br> It's a corn. <br> Es un callo. <br><br> RESPONSE2: <br> 考虑：跖疣 <br> Consideration: Plantar wart <br> Consideración: ¿Verruga plantar <br><br> RESPONSE3: <br> 是跖疣，不是鸡眼，激光治疗。 <br> It's a plantar wart, not a corn. Laser treatment is recommended. <br> Es una verruga plantar, no un callo. Se recomienda el tratamiento con láser. |

was also given a human rating for completeness and whether an answer is one that was the most frequent. The rating guide is as follows: 0.0 for no, 0.5 for partial and, 1.0 for yes. English and Spanish versions were human translated by medical translators.

## 2.2 Dataset

The dataset here was constructed by using content from a Chinese online medical platform 爱爱医[1] for posts related to dermatological problems. In the platform, users may post a question with images; doctors on the platform may respond. Thus, in our dataset, in each instance, the input is the original query and images provided by the original poster. The answer is the set of answers provided by medical experts who responded to the query.

Encounters were filtered out if it met at least one of the following exclusion criteria: (a) images that included identifying features (e.g. full faces), (b) no medical answers were given, (c) queries were not seeking information (e.g. "look at my tatoo"), and (d) images contained annotations (e.g. drawn arrows). Train/validation/test sets included 842/56/100 instances, respectively. Table 2 shows summary statistics of the data. A query can in-

volve multiple anatomic locations and medical topics (calculated by counting terms identified using QUICKUMLS(Soldaini and Goharian, 2016) on the English). The test set required at least two responses.

The data here used a subset of the DermaVQA dataset, for which the full description can be found in (Yim et al., 2024b).

## 2.3 Evaluation

We evaluate the system responses by comparing with the multiple gold standard responses per query. We used relevant multi-reference metrics/variants including:

***deltaBLEU.*** A variant of SacreBLEU developed for response generation, a case in which many diverse gold standard responses are possible (Galley et al., 2015). The metric incorporates human-annotated quality rating and assigns higher weights to n-grams from responses rated to be of higher quality. The authors have shown this method produces higher correlation with human rankings compared to previous BLEU metrics. In our system, we assign response weights according to four criteria: (a) if user expertise level is 4 or above (out of 9), (b) if user is formally validated as a medical doctor by the platform, (c) if the response answer is the most

---

[1]iiyi.com

Table 2: DermaVQA IIYI Subset Data Characteristics. (encs=encounters, encs-x img=number of encounters with x images, encs-x resp=number of encounters with x responses)

|  | TRAIN | VALID | TEST | TOTAL |
|---|---|---|---|---|
| N | 842 | 56 | 100 | 998 |
| IMAGES |  |  |  |  |
| total count | 2473 | 157 | 314 | 2944 |
| mean count | 2.9 | 2.8 | 3.1 | 2.95 |
| encs-1 img | 196 | 11 | 18 | 225 |
| encs-2 img | 233 | 22 | 30 | 285 |
| encs-3 img | 171 | 11 | 18 | 200 |
| encs->=4 img | 242 | 12 | 34 | 288 |
| RESPONSES |  |  |  |  |
| total count | 5871 | 417 | 926 | 7214 |
| mean count | 7.0 | 7.4 | 9.3 | 7.2 |
| encs-1 resp | 66 | 0 | 0 | 66 |
| encs-2 resp | 80 | 6 | 5 | 91 |
| encs-3 resp | 100 | 4 | 6 | 110 |
| encs->=4 resp | 596 | 46 | 89 | 731 |
| LENGTH (words/char) |  |  |  |  |
| per query(en/es/zh) | 80.4/81.8/89.0 | 75.0/71.9/79.0 | 76.0/74.3/81.0 | 79.6/80.5/87.6 |
| per response(en/es/zh) | 11.9/12.7/16.4 | 14.9/15.2/19.6 | 10.8/10.7/14.0 | 11.9/12.6/16.3 |
| MEDICAL TOPICS |  |  |  |  |
| Diagnosis | 610 | 196 | 137 | 695 |
| Tests | 39 | 10 | 13 | 46 |
| Treatments | 494 | 123 | 104 | 567 |
| LOCATIONS |  |  |  |  |
| Arm region | 162 | 6 | 19 | 187 |
| Back region | 85 | 10 | 9 | 104 |
| Chest/Abdomen region | 107 | 4 | 13 | 124 |
| Foot region | 129 | 8 | 15 | 152 |
| Hand region | 221 | 19 | 31 | 271 |
| Head region | 178 | 12 | 13 | 203 |
| Leg region | 198 | 12 | 21 | 231 |
| UNSPECIFIED | 161 | 9 | 25 | 195 |

frequent answer, and (d) if the response answers the query completely. The former two were manually assigned to the validation and test sets by two NLP scientists. The test set was double-reviewed. Out of a 0.0-1.0 scale, if (d) is not met, the score is discounted to 0.9; for the other 3 criteria, 0.1 is discounted for every missing element to reach the final weight.

**BERTScore.** BERTScore[2](Zhang et al., 2019) averages the maximum word embedding similarity scores between two texts based on BERT embeddings. This metric has been shown to work well on a variety of tasks, including image captioning and machine translation. The maximum was taken over multiple candidate scores.

**MEDCON.** For the English and Chinese datasets, we also used a medical information extraction based metric, here known as MEDCON. The metric uses QuickUMLS[3] to identify medical concepts in conjunction with an in-house llama-based assertion classifier. Concepts identified by QuickUMLS are normalized according to a curated concept map. Precision, recall, and F1 were calculated based on combined concept and assertion statuses. The maximum was taken over multiple candidate scores.

Evaluation code can be found in our github repo[4].

---

[2]github.com/Tiiiger/bert_score

[3]github.com/Georgetown-IR-Lab/QuickUMLS
[4]https://github.com/wyim/MEDIQA-M3G-2024

Table 3: Participant Information

| TEAM | AFFLIATION | LANGUAGES |
|---|---|---|
| Intel-PA | United Kingdom | en,zh |
| LTRC-IIITH (Thomas et al., 2024) | International Institute of Information Technology - Hyderabad, India | en |
| MediFact (Saeed, 2024) | Pakistan | en,es,zh |
| MIPAL (Kim et al., 2024) | Seoul National University / Republic of Korea | en |
| NEUI (García and Lithgow-Serrano, 2024) | The Dalle Molle Institute for Artificial Intelligence (IDSIA), Switzerland | en |
| IKIM (Bauer et al., 2024) | Institute for AI in Medicine (IKIM), University Hospital Essen, Germany | en,es,zh |
| UMass-BioNLP (Vashisht et al., 2024) | University of Massachusetts Amherst, USA | en |
| WangLab (Toma et al., 2024) | University of Toronto, Canada | en |

## 3 Baselines

For reference, we provide two baselines based on vision-language models and an additional translation step: **Baseline1:** GPT4V and GPT4(gpt, 2024), **Baseline2:** Gemini-Pro-Vision and Gemini-Pro(gem, 2024) are used.

The prompt for generating content was "请以专业人士皮肤科医生来回答，答案只限17字，不用加客气的说话。{post_title}: {post_content}" (Translation: Please answer as a professional dermatologist, answer limited to 17 characters, do not include pleasantries). The English and Spanish baselines were translations of the Chinese output. The prompt for the machine translation was "Instructions: Translate the following medical text faithfully from Chinese into {TARGET_LANGUAGE}."

## 4 Official Results

### 4.1 Participating teams

The shared task included 52 of registered participants. The final number of teams that submitted runs was 8 teams, with a total of 48 submissions. Participating teams came from various regions including Europe (3), North America (2), South Asia (2), and East Asia (1). The number of teams and submissions were 8 and 36 for English, 3 and 12 for Chinese, and 3 and 6 for Spanish. We limited the number of runs to 10. Details of the participating teams are shown in Table 3. deltaBLEU was used for official ranking.

### 4.2 Approaches and Results

Tables 4, 5, 6 detail the results for the English, Chinese, and Spanish test sets respectively. The BLEU scores ranged between 0.231-12.855, 2.171-7.053, and 0.446-1.355 for English, Chinese, and Spanish test sets. It is notable that the

magnitude of scores for both Chinese and Spanish test sets did not vary widely, the top three scores for English was significantly higher than other systems with the difference between the third best system and fourth at 7 BLEU points. BERTScore had higher ranges for English (0.800-0.886), and lower ranges for Chinese (0.685-0.764) and Spanish (0.764-0.818). In general the MEDCON scores were low, with the highest number at 0.287.

**Fine-tuned Vision-language Models Systems:** Three teams–Team MIPAL, IKIM, and LTRC-IIITH–relied fine-tuning visual-language models. The models included MedVInT(Zhang et al., 2023) and LLaVA(Liu et al., 2023), LLaVA-Med(Li et al., 2023), ViLT(Kim et al., 2021), respectively. The score variation, ranging from 0.457 to 3.827 BLEU suggests the combination of model, prompts, and fine-tuning strategy lead to large differences in results.

**Pre-trained Vision-language Systems:** As multiple submissions were allowed, the previous teams also submitted non-fined-tuned model outputs as shown in the FINE_TUNED columns of Tables 4, 5, 6. For non open models, in one submission, Team WangLab experimented with Claude3 Opus(ant, 2024), using two calls - one for candiate generation another for a final response, with competitive results. Likewise, the UMass-BioNLP used pre-trained models without fine-tuning in a multi-step fashion. The team first employed GPT-4/GPT-4-Vision(Wu et al., 2023) to generate initial hypotheses; secondly they generated image descriptors from the disease candidates of the previous step. Afterwards, they selected possible diagnosis by comparing image descriptors similarities of the disease candidates and that of the image descriptors from

Table 4: Results (English) - Top 3 Results per Team

| RANK | TEAM | FINE_TUNED | MODELS | deltaBLEU | BERTScore | MEDCON |
|------|------|------------|--------|-----------|-----------|--------|
| 1 | WangLab | FALSE | clip | 12.855 | 0.882 | 0.222 |
| 2 | WangLab | FALSE | Claude. based prompt engineering | 12.159 | 0.886 | 0.287 |
| 3 | WangLab | FALSE | fine-tuned clip | 11.979 | 0.862 | 0.125 |
| 4 | MIPAL | TRUE | PMC-VQA(PMC-CLIP, PMC-LLaMA) | 3.827 | 0.872 | 0.139 |
| 5 | MIPAL | TRUE | PMC-VQA(PMC-CLIP, PMC-LLaMA) | 3.263 | 0.872 | 0.139 |
| 6 | MIPAL | TRUE | PMC-VQA(PMC-CLIP, PMC-LLaMA) | 3.263 | 0.872 | 0.139 |
| 7 | IKIM | TRUE | llava-med + mixtral-instruct | 2.662 | 0.858 | 0.123 |
| 8 | IKIM | TRUE | llava-med, mixtral | 2.662 | 0.858 | 0.123 |
| 9 | NEUI | FALSE | Phi1 | 2.133 | 0.850 | 0.131 |
| 10 | Intel-PA | FALSE | BLIP2 | 1.758 | 0.852 | 0.155 |
| 11 | Intel-PA | FALSE | Intel-PA-run8 | 1.505 | 0.849 | 0.180 |
| 12 | UMass-BioNLP | FALSE | GPT4 | 0.923 | 0.852 | 0.159 |
| 13 | UMass-BioNLP | FALSE | GPT4 | 0.823 | 0.851 | 0.131 |
| 14 | MediFact | TRUE | VGG16-CNN-SVM | 0.717 | 0.842 | 0.148 |
| 15 | Intel-PA | FALSE | BLIP2 | 0.711 | 0.837 | 0.086 |
| 16 | UMass-BioNLP | FALSE | GPT4 | 0.670 | 0.821 | 0.158 |
| 17 | NEUI | FALSE | Phi1 | 0.595 | 0.851 | 0.205 |
| 18 | MediFact | TRUE | VGG16-CNN-SVM | 0.588 | 0.845 | 0.163 |
| 19 | MediFact | FALSE | BART, SVM, TF-IDF | 0.588 | 0.838 | 0.054 |
| 20 | IKIM | FALSE | llava med on chinese data + translation | 0.554 | 0.860 | 0.057 |
| 21 | LTRC-IIITH | FALSE | Vision-and-Language Transformer (ViLT) model - dandelin/vilt-b32-mlm | 0.457 | 0.829 | 0.016 |
| 22 | neui | TRUE | Phi1 | 0.231 | 0.810 | 0.065 |
| - | baseline1 | FALSE | GPT4 | 0.813 | 0.867 | 0.083 |
| - | baseline2 | FALSE | GEMINI | 1.094 | 0.800 | 0.157 |

Table 5: Results (Chinese) - All Results

| RANK | TEAM | FINE_TUNED | MODELS | deltaBLEU | BERTScore | MEDCON |
|------|------|------------|--------|-----------|-----------|--------|
| 1 | IKIM | TRUE | llava-med, mixtral-instruct | 7.053 | 0.764 | 0.067 |
| 2 | IKIM | FALSE | llava-med, mixtral | 7.053 | 0.764 | 0.074 |
| 3 | IKIM | FALSE | llava-med, Biomistral | 7.053 | 0.764 | 0.060 |
| 4 | Intel-PA | FALSE | – | 6.976 | 0.756 | 0.031 |
| 5 | Intel-PA | FALSE | BLIP2 | 6.976 | 0.756 | 0.029 |
| 6 | Intel-PA | FALSE | – | 5.166 | 0.757 | 0.017 |
| 7 | Intel-PA | FALSE | BLIP2 | 5.032 | 0.741 | 0.027 |
| 8 | MediFact | TRUE | VGG16-CNN-SVM | 4.503 | 0.763 | 0.106 |
| 9 | MediFact | TRUE | VGG16-CNN-SVM | 4.503 | 0.763 | 0.105 |
| 10 | Intel-PA | FALSE | BLIP2 | 4.073 | 0.731 | 0.036 |
| 11 | Intel-PA | FALSE | BLIP2 | 2.426 | 0.712 | 0.015 |
| 12 | MediFact | FALSE | BART, SVM, TF-IDF | 2.171 | 0.707 | 0.075 |
| - | baseline1 | FALSE | GPT4 | 7.025 | 0.735 | 0.016 |
| - | baseline2 | FALSE | GEMINI | 9.311 | 0.685 | 0.107 |

Table 6: Results (Spanish) - All Results

| RANK | TEAM | FINE_TUNED | MODELS | deltaBLEU | BERTScore |
|------|------|-----------|--------|-----------|-----------|
| 1 | IKIM | TRUE | llava-med, mixtral | 1.355 | 0.818 |
| 2 | NEUI | FALSE | Phi1 | 0.974 | 0.814 |
| 3 | NEUI | FALSE | Phi1 | 0.974 | 0.814 |
| 4 | MediFact | TRUE | VGG16-CNN-SVM | 0.918 | 0.806 |
| 5 | MediFact | TRUE | VGG16-CNN-SVM | 0.823 | 0.809 |
| 6 | MediFact | FALSE | BART, SVM, TF-IDF | 0.446 | 0.802 |
| - | baseline1 | FALSE | GPT4 | 0.979 | 0.822 |
| - | baseline2 | FALSE | GEMINI | 1.355 | 0.764 |

the encounter images outputted by GPT-4-Vision.

**Multi-step Mixed Model Systems:** The teams, Teams Intel-PA, NEUI, MediFact, and WangLab, experimented with a series of multiples steps using both fine-tuned and pre-trained models in a pipeline. Team Intel-PA uses a BLIP2 model, taking the output layer and combining word embeddings. These combined vectors were then fed to a large language model for text generation. Team NEUI used a fine-tuned visual language model, Moondream (https://moondream.ai/), to generate candidates. Then candidates were given as input into a BioMistral-7B-DARE(Labrak et al., 2024) to produce the final output. Team MediFact experimented with various image embedding methods, e.g. CLIP and VGG16, with a prediction task to classify a training answer response label using an SVM. The previous output combined with the query information was then fed into a reading comprehension model, Medical-QA-deberta-MRQA-COVID-QA(mrq, 2024), to generate an intermediate output. The final response is chosen by leveraging CLIP and finding the highest similarity of the image and QA output to a trained response. Google translator was used to generate the Chinese and Spanish versions. Team WangLab experimented embedding images using a fine-tuned CLIP model. The highest similarity to the test set was retrieved; the label selected from multiple gold responses in the test set was determined using GPT4. Finally, the retrieved labels were post-processed to an expected sentence format.

**Multilingual Generation Approaches**

Three patterns emerged for handling of multiple languages: (a) separate fine-tuning for each language, (b) prompt-adjustment as in Team NEUI, e.g. instructing output to be in Spanish, (c) a separate machine translation step as in Team IKIM, MediFact.

While Team IKIM fine-tuned on the Chinese dataset, then translated to English and Spanish separately using a Mixtral-8x7B-instruct model(Jiang et al., 2024); Team NEUI focused on English, translating to Spanish. The performance gap between IKIM and NEUI in English was at 0.529 BLEU, and 0.37 BLEU in Spanish. Though they used different systems, the relative scoring gap was preserved, suggesting that the two methods (b) and (c) are comparable.

The comparative effect of fine-tuning on automatically translated text prior to training versus using the original language and translating after generation requires further study.

### 4.3 Discussion and Related Work

The baseline systems using out-of-the box GPT-4-Vision and Gemini-Pro-Vision showed highly competitive performance for its original Chinese language at 7.025 and 9.311 BLEU (Table 5). However, this performance drops considerably when the same text is translated to English and Spanish; then evaluated on those test sets. Part of this drop may be due to automatic translation error, however this difference can also be partly attributed to the n-gram treatment of Chinese characters compared to latin words; which allows more partial credit. BERTScores were more stable across other languages, however are the comparatively higher compared to other metrics. MEDCON, a relatively simple, but strict metrics showed lower scores across datasets, suggesting much room for further improvement.

Although scores here are modest compared to previous Visual Question Answering (VQA) tasks. On further examination, this difference is due to the

nature of expected answers. Prior VQA datasets have question types with 1 or 2 fixed expected categorical responses. In fact, except for one work, all previous VQA tasks report accuracy as a metric. For the three cases of prior work that also report BLEU, average answer length was around 2 words. BLEU-3 scores for PathVQA were at most 17.4, even with at least half the corpus including a yes/no question type. BLEU ranges for the VQA-RAD, with more open-ended questions, achieved scores ranging from a modest 0.0058 to 0.1047 BLEU. This is consistent with recent studies which have shown that when queries are converted from a closed question-answering setting, e.g. multiple choice, to an open question-answering setting, this leads to significant degradations in performance, as much as 20%(Yim et al., 2024a).

A comparison with prior dermatological image classification tasks with user generated images also lend a helpful landscape. In Glock et al(Glock et al., 2021), with two classification categories an accuracy 95% was achieved; however for a dataset like SD-128, 128 categories, accuracy was at 52%. In a direct comparison, the authors of the Fitzpatrik 17k dataset study found a 20% accuracy when using 114 skin conditions which rises to 62% when simplifying to three categories (non-neoplastic, benign, and malignant)(Groh et al., 2021). As our gold standard responses were not generated using a fixed vocabulary, all the possible types and subtypes of diagnosis, treatments, and recommendations contributed to the difficulty of the task.

## 5 Conclusion and Future Work

Open-ended consumer health visual question answering remains a challenging problem. This shared tasks highlights several areas for future work.

One aspect is related to the generation of a dermatology common problem gold standard. Here we used a dataset with multiple references, some with varying opinions. For the dermatological specialty, a true gold standard with pathological lab confirmation is difficult to obtain in real life. This reflects the realities of current healthcare technology and costs – biological sampling and assays are only reserved for the most severe cases. Thus, datasets with biopsy observations are highly biased towards problems suspected to be malignant; whereas the plethora of other common-place maladies will remain unconfirmed. Textbook images and diagnosis

labels, on the other hand, will not include user-generated queries. This is a non-trivial hurdle if an unequivocal dermatological VQA gold standard beyond medical doctor opinion is to be achieved. Furthermore, the dataset here limits responses to queries to a single turn - however multiple turns are necessary for clarification purposes in real clinical settings.

Another future direction is the development of mature evaluation methods when multiple references of varying quality is available. In past TREC competitions, one evaluation strategy included the employment of expert humans who would annotate each participant system based on answer quality(Ben Abacha et al., 2017). Ratings include categories: (a) Correct and Complete Answer, (b) Correct but Incomplete, (c) Incorrect but Related, and (d) Incorrect. In this task, we sought to incorporate this automatically in terms of weighing response answers for BLEU. However, although this sidesteps a need for a human expert to rate each system output, this method still relies on some human annotation of the gold standard instances. As well, the final scoring depends heavily on the quality and variety of existing answers; this leaves room for metric exploitation given the data biases. For example, on observation of the test set, although responses may include a variety of responses including recommended diagnosis, treatments, and test suggestions; since most responses at least give a diagnosis, it is advantageous to optimize for a short disease response instead of try to add more details and possibly incur penalties with an incorrect suggestion. Furthermore, mentioned medical concepts may have hierarchical relations with those the gold standard for current metrics do not take into account for well. For example, atopic dermatitis is equivalent to eczema and is a subtype of dermititis – however, eczema is not the same as contact dermatitis. Depending on the available combinations of gold responses, the same system output may receive different scores.

In this shared task, a variety of solutions were explored to provide solutions for the dermatological VQA. We hope that the benchmarks provided here, the insights from different systems, and the identified methological problems will inspire future research directions.

## Limitations

The paper does not cover all types of possible methods and models for the generation of dermatological consumer health queries. The challenge datasets are limited in terms of size and coverage of diseases, treatments, and question types. The scope of the dataset only covers single turn responses. Further experiments and evaluations are needed to validate the best performing methods on other datasets and scenarios.

## Acknowledgments

## References

2024. Anthropic. https://www.anthropic.com/news/claude-3-family. Accessed: 2024-04-24.

2024. Gemini models. https://ai.google.dev/gemini-api/docs/models/gemini. Accessed: 2024-04-24.

2024. Gpt-4 turbo and gpt-4. https://platform.openai.com/docs/models/gpt-4-turbo-and-gpt-4. Accessed: 2024-04-24.

2024. Medical-qa-deberta-mrqa-covid-qa. https://huggingface.co/longluu/Medical-QA-deberta-MRQA-COVID-QA. Accessed: 2024-04-24.

Marie Bauer, Amin Dada, Constantin Marc Seibold, and Jens Kleesiek. 2024. Ikim at mediqa-m3g 2024: Multilingual visual question-answering for dermatology through vlm fine-tuning and llm translations. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Eugene Agichtein, Yuval Pinter, and Dina Demner-Fushman. 2017. Overview of the medical question answering task at trec 2017 liveqa. In *TREC 2017*.

Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Dina Demner-Fushman, and Henning Müller. 2019. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. In *Proceedings of CLEF (Conference and Labs of the Evaluation Forum) 2019 Working Notes*. 9-12 September 2019.

Tara F. Bishop, Matthew J. Press, Jayme L. Mendelsohn, and Lawrence P. Casalino. 2024. Electronic communication improves access, but barriers to its widespread adoption remain. *Health affairs (Project Hope)*, 32(8):10.1377/hlthaff.2012.1151.

Roxana Daneshjou, Kailas Vodrahalli, Weixin Liang, Roberto A. Novoa, Melissa Jenkins, Veronica Rotemberg, Justin M. Ko, Susan M. Swetter, Elizabeth E. Bailey, Olivier Gevaert, Pritam Mukherjee, Michelle Phung, Kiana Yekrang, Bradley Fong, Rachna Sahasrabudhe, James Zou, and Albert S. Chiou. 2021. Disparities in dermatology ai performance on a diverse, curated clinical image set. *Science Advances*, 8.

Michel Galley, Chris Brockett, Alessandro Sordoni, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 445–450, Beijing, China. Association for Computational Linguistics.

Ricardo Omar Chàvez García and Oscar William Lithgow-Serrano. 2024. Neui at mediqa-m3g 2024: Medical vqa through consensus. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Kimberly Glock, Charlie Napier, Todd Gary, Vibhuti Gupta, Joseph Gigante, William Schaffner, and Qingguo Wang. 2021. Measles rash identification using transfer learning and deep convolutional neural networks. In *2021 IEEE International Conference on Big Data (Big Data)*, pages 3905–3910.

Matthew Groh, Caleb Harris, Luis Soenksen, Felix Lau, Rachel Han, Aerin Kim, Arash Koochek, and Omar Badri. 2021. Evaluating deep neural networks trained on clinical images in dermatology with the fitzpatrick 17k dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1820–1828.

Xuehai He, Zhuo Cai, Wenlan Wei, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. 2021. Towards visual question answering on pathology images. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 708–718, Online. Association for Computational Linguistics.

Steven Hicks, Andrea M. Storås, Pål Halvorsen, Thomas de Lange, M. Riegler, and Vajira Lasantha Thambawita. 2023. Overview of imageclefmedical 2023 - medical visual question answering for gastrointestinal tract. In *Conference and Labs of the Evaluation Forum*.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. *Preprint*, arXiv:2401.04088.

Hyeonjin Kim, MIN KYU KIM, Jae Won Jang, KiYoon Yoo, and Nojun Kwak. 2024. Team mipal at mediqa-m3g 2024: Large vqa models for dermatological diagnosis. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. *Preprint*, arXiv:2102.03334.

Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. *Preprint*, arXiv:2402.10373.

Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. 2018. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. Llavamed: Training a large language-and-vision assistant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc.

Nadia Saeed. 2024. Medifact at mediqa-m3g 2024: Medical question answering in dermatology with

multimodal learning. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Christine A. Sinsky, Tait D. Shanafelt, and Jonathan A. Ripp. 2024. The electronic health record inbox: Recommendations for relief. 37(15):4002–4003.

Luca Soldaini and Nazli Goharian. 2016. Quickumls: a fast, unsupervised approach for medical concept extraction. In *MedIR workshop, sigir*, pages 1–4.

Jerrin John Thomas, Sushvin Marimuthu, and Parameswari Krishnamurthy. 2024. Ltrc-iiith at mediqa-m3g 2024: Efficient medical visual question answering with fine-tuned lightweight models. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Augustin Toma, Ronald Xie, Steven Palayew, Gary D. Bader, and BO WANG. 2024. Wanglab at mediqa-m3g 2024: Multimodal medical answer generation using large language models. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Parth Vashisht, Abhilasha Lodha, Mukta Maddipatla, Zonghai Yao, Avijit Mitra, Zhichao Yang, Junda Wang, Sunjae Kwon, and Hong Yu. 2024. Umass-bionlp at mediqa-m3g 2024: Dermprompt - a systematic exploration of prompt engineering with gpt-4v for dermatological diagnosis. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Yang Wu, Shilong Wang, Hao Yang, Tian Zheng, Hongbo Zhang, Yanyan Zhao, and Bing Qin. 2023. An early evaluation of gpt-4v(ision). *Preprint*, arXiv:2310.16534.

Wen-wai Yim, Yujuan Fu, Asma Ben Abacha, and Meliha Yetisgen. 2024a. To err is human, how about medical large language models? comparing pre-trained language models for medical assessment errors and reliability. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation*.

Wen-wai Yim, Yujuan Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *ArXiv*, abs/1904.09675.

Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-vqa: Visual instruction tuning for medical visual question answering. *Preprint*, arXiv:2305.10415.

# EM_Mixers at MEDIQA-CORR 2024: Knowledge-Enhanced Few-Shot In-Context Learning for Medical Error Detection and Correction

**Swati Rajwal[1], Eugene Agichtein[1], Abeed Sarker[2]**

[1]Department of Computer Science & Informatics, Emory University
[2]Department of Biomedical Informatics, Emory University
**Correspondence:** srajwal@emory.edu

## Abstract

This paper describes our submission to MEDIQA-CORR 2024 shared task for automatic identification and correction of medical errors in a given clinical text. We report results from two approaches: the first uses a few-shot in-context learning (ICL) with a Large Language Model (LLM) and the second approach extends the idea by using a knowledge-enhanced few-shot ICL approach. We used Azure OpenAI GPT-4 API as the LLM and Wikipedia as the external knowledge source. We report evaluation metrics (accuracy, ROUGE, BERTScore, BLEURT) across both approaches for validation and test datasets. Of the two approaches implemented,[1] our experimental results show that the knowledge-enhanced few-shot ICL approach with GPT-4 performed better with error flag (subtask A) and error sentence detection (subtask B) with accuracies of 68% and 64%, respectively on the test dataset. These results positioned us fourth in subtask A and second in subtask B, respectively in the shared task.

## 1 Introduction

An estimated 795,000 Americans either become permanently disabled or die each year across various healthcare settings due to misdiagnoses of serious diseases, as reported by Newman-Toker et al. (2024). The key process failures, especially in the emergency department, are errors in diagnostic assessment, test ordering, and test interpretation (Newman-Toker et al., 2023). Therefore there is a growing interest to assist clinicians in automatic medical error identification, if any, in a clinical note. The MEDIQA-CORR 2024 shared task (Ben Abacha et al., 2024a), hosted by the 6th Clinical Natural Language Processing Workshop at NAACL 2024, was proposed to encourage research

in medical error identification and correction in clinical texts. From a human perspective, these errors require medical expertise and knowledge to be both identified and corrected. Here we describe our submission to the three sub-tasks: error detection, error sentence identification, and error correction. We explore two approaches; the first uses LLM for error detection and correction while the second extends the approach by integrating an additional layer of information retrieval. We selected GPT-4 since it has shown good performance on a variety of medical tasks, according to various recent studies (Nori et al., 2023; Waisberg et al., 2023; Gertz et al., 2024). Out of the two approaches discussed here and implemented, we observed that the second approach performed better as measured by the evaluation metrics (section 4). The results for error flag (sub-task A) and error sentence detection (sub-task B) by our proposed system (approach 2) ranked fourth and second, respectively, in the shared task.
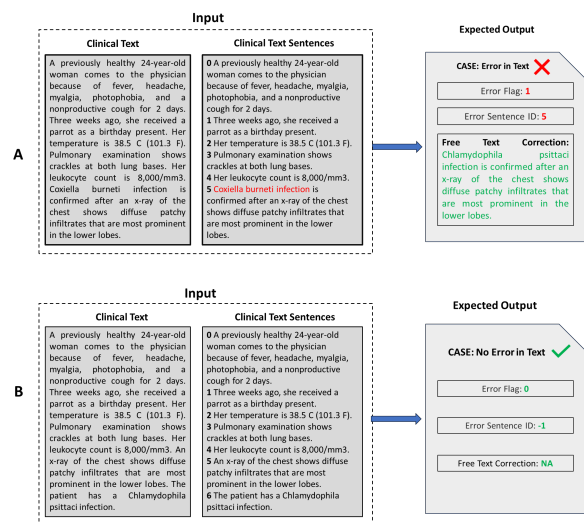


Figure 1: Example of clinical texts and clinical text sentences from the training set (Ben Abacha et al., 2024b) that have (A) a medical error and (B) no medical error.
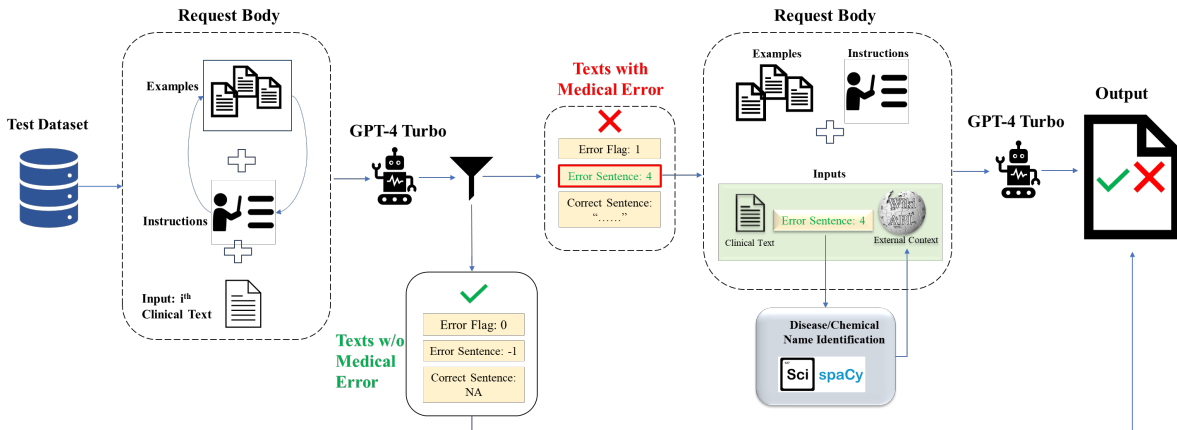
---

Figure 2: Outline of the proposed approach, illustrating the LLM and information retrieval components.

## 2 Shared Task and Dataset

MEDIQA-CORR 2024 proposed the following three sub-tasks. Each sub-task builds upon the previous one, creating a sequential process for detecting, identifying, and correcting errors in medical texts.

1. **Sub-Task A (Medical Error Identification/Binary Classification)**: Given a patient's clinical text, the task is to detect whether the text includes a medical error.

2. **Sub-Task B (Erroneous Sentence/Span Identification)**: If an error is identified in the given clinical text, the next task is to identify the text span associated with the error if a medical error exists.

3. **Sub-Task C (Correction of Erroneous Sentence)**: If the given clinical text has a medical error, this task requires rectifying or correcting the erroneous text span and providing a free text correction.

### 2.1 Dataset

The dataset (Ben Abacha et al., 2024b) was provided by two institutions: Microsoft (MS) and the University of Washington (UW). Specifically, the training dataset (MS) consists of 2,189 examples. The validation dataset contains 734 examples (574 from MS and 160 from UW, respectively) and the test set contains 925 samples. Each sample contains "Text ID" (unique), "Text" (clinical note), and "Sentences" (clinical note divided into sentences with IDs). Additionally, the training and validation dataset contains ground truth values under the columns: "Error Flag" (0 for no error, 1 otherwise),

"Error Sentence ID", "Error Sentence", "Corrected Sentence", and "Corrected Text". The mean length of a clinical text in the training dataset is 781 words (Fig. 3). The clinical text contains critical information such as symptoms, clinical examination findings, patient history, and other details. Figure 1 shows the two possible cases in the dataset—either there is a medical error in the given clinical text or there is none.
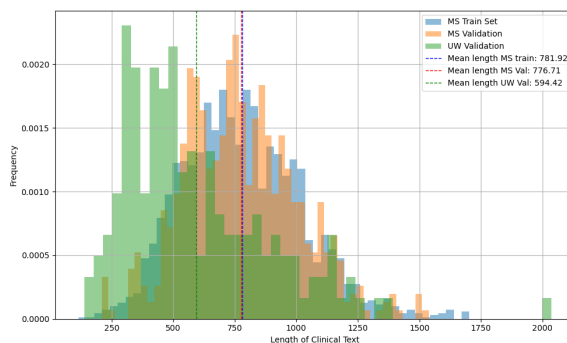


Figure 3: Clinical text lengths across datasets.

## 3 Proposed Approach

Figure 2 shows the entire framework and the following is the description of the two approaches to this shared task. We used GPT-4 as the LLM (Achiam et al., 2023) and designed a prompt to call the Microsoft Azure OpenAI GPT-4 Turbo (gpt-4-1106-preview) API[2]. This model has a context window of 128,000 tokens and returns a maximum of 4,096 output tokens. We set the temperature parameter to 0 and top_p to 0.95, respectively. For additional information and access

---

[2]https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo (last accessed: 04/24/2024)

to the code used in our study, please refer to the GitHub repository we have made publicly available.

## 3.1 Approach 1: Few-Shot In-Context Learning

For each clinical text, a request body for the GPT-4 model API is constructed as a set of instructions that outline the task of analyzing clinical text to identify and correct diagnostic errors. We provided 7 examples in the prompt to guide the LLM model in performing the analysis, followed by the actual clinical text and sentences to be evaluated. Fig. 4 shows the final prompt template which was curated over multiple manual iterations. Also, the examples in the prompts were taken from the training dataset only and remained constant across all the subsequent calls to the LLM API. The results for each clinical text were returned as a JSON object containing the error flag, erroneous sentence ID, and the corrected sentence (if any).
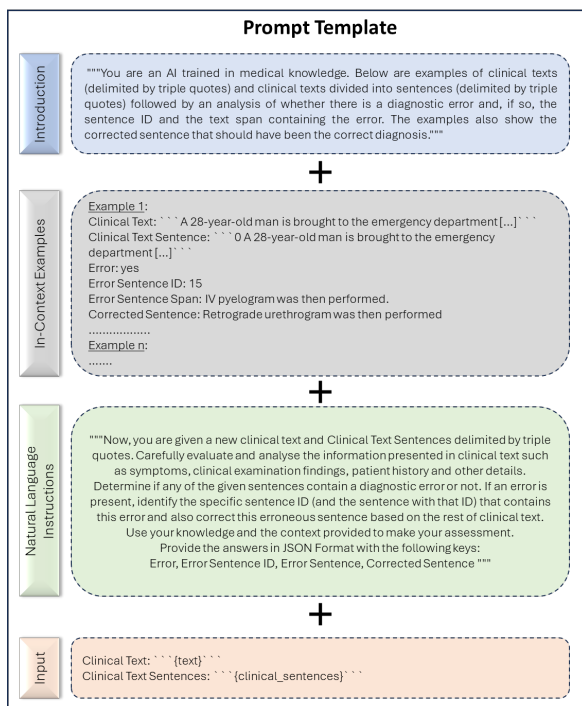


Figure 4: Prompt Template for the ICL-based approach.

## 3.2 Approach 2: Knowledge-Enhanced Few-Shot In-Context Learning

The first approach as described previously resulted in many positive predictions, especially false positives (i.e., predicted an error when there is none). Therefore, we decided to extend approach 1 (Fig. 2) by re-evaluating the instances that were previously identified by the GPT-4 model as positive (indicating the presence of an error). For such instances, an erroneous sentence was also predicted that forms a basis of re-evaluation in our approach 2.

To enrich the context for the GPT-4 model during its re-evaluation process, we integrated an additional layer of information retrieval. Specifically, we identified disease or chemical keywords within the sentence flagged as erroneous by using specialized models 'en_core_sci_scibert' and 'en_ner_bc5cdr_md' from Scispacy (Neumann et al., 2019). Then, we fetched related content from Wikipedia (English Wikimedia Wiki Endpoint[3]) and provided this external knowledge to GPT-4 alongside the original clinical text. The intention behind this strategy was to supply the model with a broader context to enable it to make more informed decisions regarding the presence of medical errors. Also, note that if there is no Wikipedia page for a particular keyword, then vanilla prompting is used (i.e., no context).

## 3.3 Rationale behind Scispacy models

ScispaCy is a Python package designed for processing biomedical, scientific, or clinical text using spaCy models. We utilized all eight available models to analyze the training dataset, specifically focusing on detecting keyword in the sentence marked as erroneous errors flagged by GPT-4. Our goal was to identify disease and chemical names in sentences where GPT-4 predicted errors. In our analysis, two models, 'en_core_sci_scibert' and 'en_ner_bc5cdr_md', worked well for keyword identification. Sometimes, en_core_sci_scibert' missed certain keywords that 'en_ner_bc5cdr_md' could detect, and vice versa. Consequently, we decided to use both models to ensure comprehensive keyword detection. As an example, take a look at Figure A.1, which shows that most of the keywords of concern are detected by one or the other model.

## 3.4 Final Submission

Our final submission for the shared task included combined analysis through an ensemble method: For each instance if the error flag from Approach 1 is set to 0, the process moves to the next instance. If both approaches agree on the presence of an error (error flag = 1), the final result (dataframe in

---

[3]https://en.wikipedia.org/w/api.php (last accessed: 04/24/2024)

Table 1: Comparison of Approach 1 (few-shot in-context learning) and Approach 2 (knowledge-enhanced few-shot in-context learning) on validation and test datasets.

| Metric | Validation Dataset | | Test Dataset | |
| --- | --- | --- | --- | --- |
| | Approach 1 | Approach 2 | Approach 1 | Approach 2 |
| **Accuracy** | | | | |
| Error Flags Accuracy | 0.622 | 0.648 | 0.626 | **0.680**[a] |
| Error Sentence Detection Accuracy | 0.598 | 0.638 | 0.562 | **0.640**[b] |
| **ROUGE Scores** | | | | |
| R1F_subset_check | 0.488 | 0.550 | 0.540 | 0.571 |
| R2F_subset_check | 0.375 | 0.439 | 0.444 | 0.478 |
| RLF_subset_check | 0.481 | 0.543 | 0.534 | 0.565 |
| R1FC | 0.369 | 0.516 | 0.429 | 0.542 |
| R2FC | 0.313 | 0.484 | 0.388 | 0.512 |
| RLFC | 0.365 | 0.514 | 0.426 | 0.540 |
| **BERTScore** | | | | |
| BERTSCORE_subset_check | 0.566 | 0.620 | 0.574 | 0.595 |
| BERTC | 0.407 | 0.537 | 0.444 | 0.550 |
| **BLEURT** | | | | |
| BLEURT_subset_check | 0.569 | 0.607 | 0.580 | 0.596 |
| BLEURTC | 0.409 | 0.533 | 0.446 | 0.550 |
| **Average Composite Score** | | | | |
| aggregate_subset_check | 0.541 | 0.592 | 0.565 | 0.587 |
| AggregateC | 0.395 | 0.529 | 0.440 | 0.548 |

[a] Fourth and [b] Second best accuracy in the shared task results among 17 participating teams.

Python) is updated with the error flag with the sentence ID from Approach 1, and the corrected sentence as identified. If Approach 1 flags an error but Approach 2 does not, the instance is left unchanged, moving on to the next. This methodical combination of inferences from both approaches forms our final solution for error identification and correction mechanism essentially giving more weightage to the knowledge-enhanced approach.

## 3.5 Evaluation

The official evaluation script[4] provided by the organizers was used for model evaluation. The test set results were released after system submission on codalab. The proposed systems predictions are evaluated for binary accuracy of error detection and a multi-dimensional evaluation of text correction quality against the provided ground truth notes with the following metrics: ROUGE (Lin, 2004), BERTScore (Zhang* et al., 2020), and BLEURT (Sellam et al., 2020).

## 4 Results

Table 1 shows the results on the validation and test dataset for multiple evaluation metrics. Refer to Appendix A.1 for the detailed definition of each metric variable name.

**Accuracy Metrics**: Experimental results show that Approach 2 improved error flag accuracy by about 2.6% on the validation dataset and 5.4% on the test dataset. Similarly, for Error Sentence Detection Accuracy, Approach 2 shows an improvement of approximately 4% and 7.8% on the validation and test datasets, respectively. This suggests that providing external context around the disease/chemical name is useful (to a certain extent) for GPT-4 in making sound decisions.

**ROUGE Scores**: Approach 2 demonstrates a higher score compared to Approach 1, with improvements of approximately 6.2% and 3.2% on the validation and test datasets, respectively. Similar performance improvements were observed for BERTScore, BLEURT and Average Composite scores.

## 5 Discussion

Across multiple evaluation metrics and datasets, Approach 2 consistently outperforms Approach 1. This indicates that the addition of external knowledge is potentially leveraging more effective strategies for both error detection and error correction.

### 5.1 Error Analysis

We studied the misclassified examples in the dataset. It appears that the model found it challenging to recognize rare or complex conditions (e.g., Picornavirus, being less commonly referenced in lay texts). Although external information from Wikipedia is used to provide context, GPT-4's interpretation of this supplementary data is still limited by its ability to integrate and analyze it effectively within the clinical scenario presented. This process might have been complicated due to Wikipedia content being too general to aid in accurate analysis.

### 5.2 Limitations & Future Directions

Automatic evaluation metrics such as ROUGE, BERTScore, and BLEURT may not accurately reflect human judgment. Therefore, in real-life settings, it is necessary to conduct an expert human evaluation to validate the results. Furthermore, our current approach uses Wikipedia as the external source of information which, while a rich source of information, might not be very specialized for medical knowledge. In the future, we plan to utilize other sources of medical knowledge, such as PubMed. During the second approach, we rely solely on the sentence that has been predicted by GPT-4 to be erroneous. This might be wrong since there were cases when GPT-4 correctly identified that there was an error in the clinical text but incorrectly identified the erroneous sentence span which is the basis of our knowledge-retrieval component.

## 6 Conclusion

In this paper, we present our submission to the MEDIQA-Corr shared task for Medical Error Detection and Correction. We evaluated two approaches: one with in-context learning (ICL) and the other an extension with knowledge-enhanced few-shot ICL. Based on the evaluation metric results, we conclude that knowledge-enhanced few-shot in-context learning is a promising path toward medical error detection and correction. For future work, we plan to experiment the proposed pipeline with other sources of medical information for comparative analysis.

## Acknowledgement

## Data & Code Availability

https://github.com/
swati-rajwal/EM_Mixers_
MEDIQA-CORR-NAACL-ClinicalNLP-2024

## Ethics Statement

Design and development of an automated system for medical error detection and correction can raise many ethical issues. For instance, the system design should address issues of data bias and fairness to avoid unfair medical error detection for certain patient groups. Also, transparency about the system's capabilities and limitations is key, allowing users to understand and trust our AI's decisions. We also emphasize the importance of sourcing credible information, particularly when integrating external content like Wikipedia, to maintain the accuracy and relevance of our corrections.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Roman Johannes Gertz, Thomas Dratsch, Alexander Christian Bunck, Simon Lennartz, Andra-Iza Iuga, Martin Gunnar Hellmich, Thorsten Persigehl, Lenhard Pennig, Carsten Herbert Gietzen, Philipp Fervers, David Maintz, Robert Hahnfeldt, Jonathan Kottlors, and Linda Moy. 2024. Potential of gpt-4 for detecting errors in radiology reports: Implications for reporting accuracy. *Radiology*, 311(1):e232714. PMID: 38625012.

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.

Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. 2019. ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 319–327, Florence, Italy. Association for Computational Linguistics.

David E Newman-Toker, Najlla Nassery, Adam C Schaffer, Chihwen Winnie Yu-Moe, Gwendolyn D Clemens, Zheyu Wang, Yuxin Zhu, Ali S Saber Tehrani, Mehdi Fanai, Ahmed Hassoon, et al. 2024. Burden of serious harms from diagnostic error in the usa. *BMJ Quality & Safety*, 33(2):109–120.

David E Newman-Toker, Susan M Peterson, Shervin Badihian, Ahmed Hassoon, Najlla Nassery, Donna Parizadeh, Lisa M Wilson, Yuanxi Jia, Rodney Omron, Saraniya Tharmarajah, et al. 2023. Diagnostic errors in the emergency department: a systematic review.

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*.

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.

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.

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

## A  Appendix

### A.1  Evaluation Metrics

'aggregate_subset_check' is the mean score of all individual metric scores combined for each subset of data.

'R1F_subset_check' is the $F_1$-score of the ROUGE-1 metric and assesses how many of the same words are used in both texts, adjusted for both precision and recall.

'R2F_subset_check' is the $F_1$-score of the ROUGE-2 metric, focusing on the overlap of bigrams between the generated and reference texts.

Figure A.1: ScispaCy models for entity detection.



'RLF_subset_check' is the score for the ROUGE-L metric and measures the longest common subsequence between the generated and reference texts.

'R1FC', 'R2FC', and 'RLFC' are composite scores for the ROUGE-1, ROUGE-2, and ROUGE-L metrics, respectively, adjusted for the total number of texts, including those correctly identified as no error ("NA" cases). These scores balance between correctly generated corrections and correctly identified non-correction scenarios.

'BERTSCORE_subset_check' reflects the mean BERTScore $F_1$ metric and uses BERT's contextual embeddings to compare the generated text against references. 'BERTC' is the composite score for BERTScore, taking into account the entire dataset and adjusting for "NA" cases similar to the ROUGE composite scores.

'BLEURT_subset_check' represents the mean BLEURT score for the subsets of data. BLEURT is a learned metric that compares generated text to reference texts, fine-tuned on human judgments.

'BLEURTC' is the composite score for BLEURT, adjusted for the total dataset including "NA" scenarios.

'AggregateC' is the average composite score of all individual metrics (ROUGE-1 $F_1$, BERTSCORE, BLEURT), providing a single, consolidated measure of the NLG system's performance across the entire evaluation framework.

# Overview of the MEDIQA-CORR 2024 Shared Task
# on Medical Error Detection and Correction

**Asma Ben Abacha, Wen-wai Yim**

Microsoft Health AI

{abenabacha,yimwenwai}@microsoft.com

**Yujuan Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen**

University of Washington

{velvinfu,zhaoyis,fxia,melihay}@uw.edu

## Abstract

Automatic detection and correction of medical errors enables a more rigorous validation of medical documentation as well as clinical notes generated by large language models. Such solutions can ensure the accuracy and medical coherence of clinical texts and enhance patient care and health outcomes. The MEDIQA-CORR 2024 shared task focused on detecting and correcting different types of medical errors in clinical texts. Seventeen teams participated in the shared task and experimented with a broad range of approaches and models. In this paper, we describe the MEDIQA-CORR task, datasets, and the participants' results and methods.

## 1 Introduction

A recent survey study from three US health care organizations showed that 1 in 5 patients who read a clinical note reported finding a mistake and 40% perceived the mistake as serious. Among the very serious errors reported by patients, the most common category of mistakes was related to current or past diagnoses. Other very serious patient-reported mistakes included inaccurate description of medical history, medications or allergies, physical examination, test results, notes on the wrong patient, and sidedness (left vs right) (Bell et al., 2020).

Giardina et al. (2018) focused on diagnostic errors and analyzed patient- and family-reported error narratives to explore factors that contribute to diagnostic errors. Problems related to patient-physician interactions emerged as major contributors.

The probability of such errors is expected to increase in medical documents and clinical notes generated by Large Language models (LLMs) to assist healthcare professionals in their daily documentation tasks.

On a general, coarse-grained, level, LLMs have shown the ability to imitate clinical reasoning while forming mostly accurate diagnoses (Savage et al.,

2024). However, one of the main challenges in integrating LLMs in medical workloads is their potential to generate misleading or incorrect information (Tang et al., 2023). Rigorous validation processes are essential to mitigate these risks and make LLMs safe(r) to use for medical content generation (Karabacak and Margetis, 2023).

One important aspect of this validation is medical common-sense checking to validate the coherence and soundness of the generated medical reasoning. However, most previous studies on common sense detection have focused on the general domain (Wang et al., 2020; Onoe et al., 2021).

In this task, we tackle the problem of identifying and correcting (common sense) medical errors in clinical notes. From a human perspective, identifying and correcting these errors requires medical expertise, specialized knowledge, and sometimes practical experience. To the best of our knowledge, this task is the first to address the automatic validation and correction of clinical notes.

## 2 Task Description

The MEDIQA-CORR 2024 shared task[1] addresses the problem of identifying and correcting (common sense) medical errors in clinical notes. From a human perspective, identifying and correcting these errors require medical expertise and knowledge.

In the task data, each clinical text is either correct or contains one error. The task consists of three subtasks:

**A:** Predicting the error flag (1: the text contains an error, 0: the text has no errors)

**B:** Extracting the sentence that contains the error for flagged texts (-1: the text contains no error; Sentence ID: if the text contains an error)

---

[1] https://sites.google.com/view/mediqa2024/mediqa-corr

| DIAGNOSIS | CAUSAL ORGANISM | MANAGEMENT | TREATMENT | PHARMACOTHERAPY |
|---|---|---|---|---|

**ERROR**

A 17-year-old boy is brought to the physician by his mother because of increasingly withdrawn behavior for the last two years. His mother reports that in the last 2-3 years of high school, her son has spent most of his time in his room playing video games. He does not have any friends and has never had a girlfriend. He usually refuses to attend family dinner and avoids contact with his siblings. The patient states that he prefers being on his own. When asked how much playing video games means to him, he replies that "it's okay." When his mother starts crying during the visit, he appears indifferent. Physical and neurologic examinations show no other abnormalities. **Suspected of autism spectrum disorder.** On mental status examination, his thought process is organized and logical. His affect is flattened.

A 64-year-old man is brought to the emergency department because of fever, chills, shortness of breath, chest pain, and a productive cough with bloody sputum for the past several days. He has metastatic pancreatic cancer and is currently undergoing polychemotherapy. His temperature is 38.3 C (101 F). Pulmonary examination shows scattered inspiratory crackles in all lung fields. A CT scan of the chest shows multiple nodules, cavities, and patchy areas of consolidation. **Histoplasma capsulatum was determined as the causal pathogen.** A photomicrograph of a specimen obtained on pulmonary biopsy is shown.

A 42-year-old woman comes to the physician because of a low-grade fever and generalized fatigue for a week. During this period, she has passed decreased amounts of urine. Two months ago, she underwent a renal allograft transplant because of reflux nephropathy. There is no family history of serious illness (...) Oral fluconazole is administered. **Patient was recommended intravenous immunoglobulin therapy as a next step in management.**

A 47-year-old woman comes to the physician because of easy bruising and fatigue. She appears pale. Her temperature is 38 C (100.4 F). Examination shows a palm-sized hematoma on her left leg. Abdominal examination shows an enlarged liver and spleen. **Based on the following findings, patient was treated with platelet transfusion.** Her hemoglobin concentration was 9.5 g/dL, leukocyte count was 12,300/mm3, platelet count was 55,000/mm3, and fibrinogen concentration was 120 mg/dL. Cytogenetic analysis of leukocytes showed a reciprocal translocation of chromosomes 15 and 17.

A 67-year-old man with type 2 diabetes mellitus and benign prostatic hyperplasia comes to the physician because of a 2-day history of sneezing and clear nasal discharge. He has had similar symptoms occasionally in the past. His current medications include metformin and tamsulosin. Examination of the nasal cavity shows red, swollen turbinates. **The patient is given diphenhydramine.**

**Correction**

Suspected of schizoid personality disorder.

Aspergillus fumigatus was determined as the causal pathogen.

Patient was recommended methylprednisolone therapy as a next step in management.

Based on the following findings, patient was treated with all-trans retinoic acid.

The patient is given desloratadine.

Figure 1: Examples from the MEDIQA-CORR MS training set.

**C:** Generating a corrected sentence for flagged texts

## 3 Data Creation

We created a new dataset of 3,848 clinical texts with injected errors such as diagnosis, causal organism, management, treatment, and pharmacotherapy (Ben Abacha et al., 2024). The dataset includes two types of texts: clinical texts from publicly available data (MS collection) and de-identified clinical notes from the University of Washington Medical Center (UW collection). The UW dataset was built using new de-identified notes and requires signing a data usage agreement. The MS dataset was built by transforming the MedQA medical question-answering dataset (Jin et al., 2020) with manual error injections and text modifications that leveraged the clinical notes and the multiple-choice questions.

The MS training set contains 2,189 clinical texts. Figure 1 presents examples from the MS training data. The MS validation set contains 574 clinical texts and the UW validation set contains 160 clini-

cal texts. The final test set consists of 597 clinical texts from the MS collection and 328 clinical texts from the UW dataset.

## 4 Evaluation

### 4.1 Evaluation Metrics

We rely on Accuracy for Error Flag Prediction (subtask A) and Error Sentence Detection (subtask B).

For the evaluation of Sentence Correction (subtask C), we selected three automatic metrics that highly correlate with human judgments on clinical texts based on recent studies (Ben Abacha et al., 2023a,b). These metrics are: ROUGE-1 (Lin, 2004), BLEURT (Sellam et al., 2020), and BERTScore (Zhang et al., 2020).

Similar to MEDIQA-Chat (Ben Abacha et al., 2023) and MEDIQA-SUM 2023 (Yim et al., 2023), we used the aggregate (average) score from ROUGE-1, BLEURT-20, and BERTScore (microsoft/deberta-xlarge-mnli) as the main score to rank the participating systems.

We also computed a Composite score as follows

597

| | Team | Affiliation | Subtasks | Paper | Code |
|---|---|---|---|---|---|
| 1 | WangLab | University of Toronto, Canada | 1, 2, 3 | (Toma et al., 2024) | 1 |
| 2 | PromptMind | Google, USA | 1, 2, 3 | (Gundabathula and Kolar, 2024) | 2 |
| 3 | HSE NLP | Higher School of Economics University, Russia | 1, 2, 3 | (Valiev and Tutubalina, 2024) | 3 |
| 4 | KU-DMIS | Korea University | 1, 2, 3 | (Hwang et al., 2024) | 4 |
| 5 | Maven | Pune Institute of Computer Technology, India | 1, 2, 3 | (Jadhav et al., 2024) | 5 |
| 6 | Edinburgh Clinical NLP | University of Edinburgh, Scotland | 1, 2, 3 | (Gema et al., 2024) | 6 |
| 7 | knowlab_AIMed | University College London & The University of Hong Kong | 1, 2, 3 | (Wu et al., 2024) | 7 |
| 8 | EM_Mixers | Emory University, USA | 1, 2, 3 | (Rajwal et al., 2024) | 8 |
| 9 | IryoNLP | Microsoft, Canada | 1, 2, 3 | (Corbeil, 2024) | 9 |
| 10 | IKIM | Institute for AI in Medicine, Germany | 1, 2, 3 | - | 10 |
| 11 | CLD-MEC | Princess Sumaya University for Technology, Jordan | 1, 2, 3 | (Alzghoul et al., 2024) | 11 |
| 12 | romarcg | IDSIA, Switzerland | 1, 2, 3 | - | 12 |
| 13 | mekki | Um6p College Of Computing, Morocco | 1, 2, 3 | - | 13 |
| 14 | MediFact | National University of Computer and Emerging Sciences, Pakistan | 1, 2, 3 | (Saeed, 2024) | 14 |
| 15 | harivm | University of California, Los Angeles (UCLA), USA | 1, 2, 3 | - | 15 |
| 16 | VerbaNexAI | Pontificia Universidad Javeriana, Colombia | 1, 2 | (Pajaro et al., 2024) | 16 |
| 17 | nlp-lab-iu | Indiana University Bloomington, USA | 1, 2 | - | 17 |

1 https://github.com/bowang-lab/mediqacorr24
2 https://github.com/satyakesav/medical-error-detection-and-correction
3 https://github.com/Rebell-Leader/mediqa-corr
4 https://github.com/HwangHyeoni/MEDIQA-CORR-2024
5 https://github.com/abhayshanbhag2003/MEDIQA-NAACL
6 https://github.com/aryopg/mediqa
7 https://github.com/wuzl01/Knowlab_MEDIQA-CORR-2024
8 https://github.com/swati-rajwal/EM_Mixers_MEDIQA-CORR-NAACL-ClinicalNLP-2024
9 https://github.com/jpcorb20/mediqa-corr-llm
10 https://github.com/dadaamin/MEDIQA-CORR-2024
11 https://github.com/Renadzghoul/CLD-MEC
12 https://github.com/OWLmx/mediqa2024_medicorr
13 https://github.com/4mekki4/MEDIQA-CORR-2024
14 ttps://github.com/NadiaSaeed/MediFact-MEDIQA-CORR-2024
15 https://github.com/Hari-vm-01
16 https://github.com/DavidVilem/Caoba
17 https://github.com/dhananjay-srivastava/MEDIQA-CORR

Table 1: MEDIQA-CORR 2024: Participating teams, subtasks, papers, and codes.

for each text: (i) 1 point if both the system correction and the reference correction are "NA": i.e., both the reference and system agree that the text has no errors, (ii) 0 points if only one of the system or the reference is "NA" (i.e., disagreement on error presence), and (iii) Aggregate-Score if both the system and reference agree that the sentence has an error.

Our evaluation scripts are available online[2].

## 4.2 Code Verification

For additional validation, we required the submission of the code in addition to the models' outputs/runs. The participants shared their private

codes with the organizers on GitHub following provided guidelines.

## 4.3 Baseline System

We built a GPT-4-based baseline system, with deterministic outputs (temperature=0), using the following prompt for the three subtasks:

- *The following is a medical narrative about a patient. You are a skilled medical doctor reviewing the clinical text. The text is either correct or contains one error. The text has a sentence per line. Each line starts with the sentence ID, followed by a pipe character then the sentence to check. Check every sentence of the text. If the text is correct return the following output: CORRECT. If the text has*

| Team | Error Flag Accuracy | Error Sentence Detection Accuracy |
|---|---|---|
| WangLab * | 0.8649 | 0.8357 |
| PromptMind | 0.6216 | 0.6086 |
| HSE NLP | 0.5222 | 0.5200 |
| KU-DMIS | 0.6346 | 0.6151 |
| Maven * | 0.5600 | 0.5200 |
| Edinburgh Clinical NLP | 0.6692 | 0.6108 |
| knowlab_AIMed | 0.6941 | 0.6195 |
| EM_Mixers | 0.6800 | 0.6400 |
| IryoNLP | 0.6714 | 0.6097 |
| IKIM | 0.6778 | 0.5903 |
| CLD-MEC | 0.5665 | 0.4908 |
| romarcg | 0.5016 | 0.3784 |
| mekki | 0.5395 | 0.3632 |
| MediFact | 0.7373 | 0.6000 |
| harivm | 0.5027 | 0.1924 |
| nlp-lab-iu | 0.5124 | 0.0497 |
| VerbaNexAI | 0.5103 | 0.4865 |
| Baseline (GPT-4) | 0.6562 | 0.5503 |

Table 2: Official Results of Error Flag Prediction (Subtask A) and Error Sentence Detection (Subtask B). * Potential use of MS test data.

*a medical error, return the sentence id of the sentence containing the error, followed by a space, and a corrected version of the sentence.*

# 5 Official Results

## 5.1 Participating Teams

The MEDIQA-CORR 2024 shared task attracted 112 registered teams from academy and industry. Among them, seventeen teams submitted their codes and runs following the challenge rules. Table 1 presents the teams that participated in the three subtasks. We limited the number of submitted runs to 20 runs per team.

## 5.2 Results & Approaches

The main results of the challenge are presented in Table 2 and Table 3.

The WangLab team (Toma et al., 2024) achieved the best Accuracy of 0.8649 in Error Flag Prediction (subtask A) and 0.8357 in Error Sentence Detection (subtask B). They also achieved the best Aggregate-Score of 0.7891 and Aggregate-Composite of 0.7746 in Sentence Correction (subtask C). The WangLab team used two different methods for the MS and UW datasets. They leveraged the MedQA medical question-answering dataset (Jin et al., 2020) to isolate questions resembling those in the MS data. This likely led to test data leakage as the MedQA dataset was used to build the MS subset.

They employed DSPy (Khattab et al., 2023), a framework for automating the optimization of LLM programs, to refine a series of modules aimed at detecting and correcting errors. They also implemented a distinct set of DSPy modules to develop LLM-based programs for error identification and correction in the UW dataset.

The PromptMind team (Gundabathula and Kolar, 2024) achieved the second best aggregate score of 0.7866 in error sentence correction with 0.6216 error flag accuracy and 0.6086 error sentence detection accuracy using a prompt-based in-context learning strategy. They combined the results of GPT-4 and Claude-3 Opus models to generate the error flag, error sentence ID, and corrected sentence.

The third best aggregate score was obtained by the HSE NLP team (Valiev and Tutubalina, 2024) with an in-prompt ensemble approach with named entity recognition and knowledge graph for medical error checking. Their approach consists of three key components: entity extraction, prompt engineering, and ensemble. First, they automatically extract biomedical entities such as therapies, diagnoses, and biological species. Next, they explore few-shot learning techniques and incorporate graph information from the MeSH database for the identified entities. Finally, they investigate two methods for ensembling: (i) combining the predictions of three previous LLMs using an AND strat-

| Team | ROUGE1 | BERTSCORE | BLEURT | AggregateComposite | AggregateScore | Rank |
|---|---|---|---|---|---|---|
| WangLab * | 0.7755 | 0.8087 | 0.7831 | 0.7746 | 0.7891 | 1 |
| PromptMind | 0.8070 | 0.8058 | 0.7470 | 0.5739 | 0.7866 | 2 |
| HSE NLP | 0.7795 | 0.8059 | 0.7564 | 0.5117 | 0.7806 | 3 |
| KU-DMIS | 0.7288 | 0.7672 | 0.7047 | 0.5709 | 0.7336 | 4 |
| Maven * | 0.7031 | 0.7437 | 0.7522 | 0.5239 | 0.7330 | 5 |
| Edinburgh Clini-calNLP | 0.6780 | 0.7435 | 0.7111 | 0.5629 | 0.7109 | 6 |
| knowlab_AIM | 0.6435 | 0.6767 | 0.6542 | 0.5731 | 0.6581 | 7 |
| EM_Mixers | 0.5713 | 0.5952 | 0.5959 | 0.5475 | 0.5875 | 8 |
| IryoNLP | 0.5607 | 0.5916 | 0.5905 | 0.5283 | 0.5810 | 9 |
| IKIM | 0.5233 | 0.5644 | 0.5882 | 0.5500 | 0.5587 | 10 |
| CLD-MEC | 0.4273 | 0.4837 | 0.5318 | 0.3448 | 0.4809 | 11 |
| romarcg | 0.4323 | 0.4574 | 0.4608 | 0.3227 | 0.4501 | 12 |
| mekki | 0.4180 | 0.4592 | 0.4679 | 0.3997 | 0.4483 | 13 |
| MediFact | 0.4540 | 0.4441 | 0.4386 | 0.5353 | 0.4456 | 14 |
| harivm | 0.1431 | 0.1345 | 0.2563 | 0.1766 | 0.1780 | 15 |
| Baseline (GPT-4) | 0.5559 | 0.5801 | 0.5900 | 0.4726 | 0.5754 | - |

Table 3: Official Results of Error Sentence Correction (Subtask C). The teams are ranked according to AggregateScore. * Potential use of MS test data.

egy within a prompt, and (ii) integrating the previous predictions into the prompt as separate 'expert' solutions, accompanied by trust scores representing their performance. The latter system ranked second in BERTScore (0.8059) and third in aggregated score (0.7806), with an error flag accuracy of 0.5222 and an error sentence detection accuracy of 0.5200.

The KU-DMIS team (Hwang et al., 2024) generated a Chain-of-Thought reasoning dataset using GPT-4 and MEDIQA-CORR dataset. Subsequently, they fine-tuned Meerkat-7B with this generated dataset to enhance its error detection and correction capabilities. The fine-tuned model achieved an aggregate score of 0.7336 in error sentence correction, with a 0.6346 error flag accuracy and 0.6151 error sentence detection accuracy.

The Maven team (Jadhav et al., 2024) conducted Named Entity Recognition (NER) using GEMINI to identify words representing diseases or vaccines in the text. After masking these identified words, the team implemented the Retrieval-Augmented Generation (RAG) model on external datasets . If the RAG score fell below a certain threshold, they passed the input to the model, which was created by quantizing Palmyra-20b (Team, 2023) using 4-bit quantization and then fine-tuned it using the QLoRA technique on MedQA data (possible test data leakage). If the word provided by Palmyra or

RAG model matched the word detected by NER, no error was detected. Otherwise, if a different word was obtained, it was replaced with the masked word identified by NER. Finally the error sentence is mapped with the sentence Id to get the output in desired format.

The Edinburgh Clinical NLP team (Gema et al., 2024) evaluated multiple prompting strategies such as In-context Learning (ICL) and Chain-of-Thought (CoT) to improve LLMs' performance. To aid the error correction LLM, they experimented with integrating a relatively smaller language model (i.e. BioLinkBERT) as an error-span predictor. They integrated the predicted error span in two ways; presenting it as a hint for the LLM to correct the error or presenting it as multiple-choice questions for the LLM to select the most likely one.

The knowlab_AIMed team (Wu et al., 2024) used two methods: (i) Dynamic In-Context Learning with RAG, CoT, and manual analysis. In this method, they performed manual analysis on a subset of the dataset. They used the RAG model to implement dynamic ICL, incorporating CoT prompts. They also used ICL-augmented examples from the training dataset. In the second method, the team utilized the training dataset to prompt an LLM to deduce reasons about the correctness or incorrectness of the clinical notes. By leveraging the LLM's capabilities, the constructed reasons provided ad-

ditional information and insights into the errors present in the notes. These reasons, along with the ICL examples, were used to train the model for error detection, span identification, and error correction tasks.

The IKIM team (Amdada et al., 2024) trained a linear classifier on embeddings from the model pritamdeka/S-PubMedBert-MS-MARCO to predict whether a sentence potentially contains an error in the MS dataset. They also clustered these sentence embeddings. For each cluster, they leveraged GPT-4 to generate a chain of thought that describes the medical reasoning for a sample from the training dataset. For test predictions, they gave GPT-4 the sentence predicted by the linear classifier along with a chain of thought from the cluster to which the sentence belongs, and prompted it to predict whether the sentence was wrong and to provide a correction if needed. They directly prompted GPT-4 with few-shot examples and a chain of thought prompt for UW samples, without clustering or sentence selection.

The MediFact team (Saeed, 2024) employed weakly-supervised SVM and extractive QA for observed errors, alongside pre-trained QA models for unseen errors in clinical text correction. The team achieved the second best score in error flag detection with an accuracy of 0.7373, and an aggregate score of 0.4456 in error sentence correction.

## 6 Conclusion

The MEDIQA-CORR shared task was tackled by a wide variety of approaches from the participating teams. Ranging from algorithmic reasoning approaches leveraging the LLMs as intermediate extraction tools (e.g., for NER) to approaches that are fully controlled by LLMs and prompting techniques. The best performing methods were dataset-dependent, i.e., different methods or parameters were used for each dataset. Generalized, dataset-agnostic, approaches fared reasonably well in comparison. A key challenge was in detecting correctly which text and which sentence contained errors, with only two teams reaching an accuracy above 70% in text flagging, and only one team reaching an accuracy greater than 65% in detecting the sentence containing the error. The detection accuracy impacted the quality of the corrected texts (e.g., providing corrections when the sentence contained no errors) but the correction results were less contrasted in general with six teams reaching an aggregate score greater than 70%.

Moving forward, optimizing the dataset-agnostic approaches is likely to be a key focus as it has the most impact on production-grade models/systems for clinical note generation/validation. The data provided by MEDIQA-CORR can be leveraged for that as they showed to be sufficiently challenging to be used as a benchmark for generalized approaches.

## 7 Limitations

The paper does not cover all types of possible methods and models for the detection and correction of medical errors. The MS and UW datasets are also limited in terms of size and types of medical errors. Further experiments and evaluations are needed to validate the best performing methods on other datasets and scenarios.

## Acknowledgements

## References

Renad M. Alzghoul, Abdulrahman Tabaza, Aya Abdelhaq, and Ahmad Altamimi. 2024. Cld-mec at mediqa- corr 2024 task: Gpt-4 multi-stage clinical chain of thought prompting for medical errors detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Sigall K. Bell, Tom Delbanco, Joann G. Elmore, Patricia S. Fitzgerald, Alan Fossa, Kendall Harcourt, Suzanne G. Leveille, Thomas H. Payne, Rebecca A. Stametz, Jan Walker, and Catherine M. DesRoches. 2020. Frequency and types of patient-reported errors in electronic health record ambulatory care notes. *JAMA Netw Open*, 3(6).

Asma Ben Abacha, Wen wai Yim, Velvin Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Asma Ben Abacha, Wen-wai Yim, Griffin Adams, Neal Snider, and Meliha Yetisgen. 2023. Overview of the mediqa-chat 2023 shared tasks on the summarization & generation of doctor-patient conversations. In *Proceedings of the 5th Clinical Natural Language Processing Workshop, ClinicalNLP@ACL 2023, Toronto, Canada, July 14, 2023*, pages 503–513. Association for Computational Linguistics.

Asma Ben Abacha, Wen-wai Yim, Yadan Fan, and Thomas Lin. 2023a. An empirical study of clinical note generation from doctor-patient encounters. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2291–2302, Dubrovnik, Croatia. Association for Computational Linguistics.

Asma Ben Abacha, Wen-wai Yim, George Michalopoulos, and Thomas Lin. 2023b. An investigation of evaluation metrics for automated medical note generation. In *ACL (Findings) 2023*, Toronto, Canada. Association for Computational Linguistics.

Jean-Philippe Corbeil. 2024. Iryonlp at mediqa-corr 2024: Tackling the medical error detection & correction task on the shoulders of medical agents. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Aryo Pradipta Gema, Chaeeun Lee, Pasquale Minervini, Luke Daines, T. Ian Simpson, and Beatrice Alex. 2024. Edinburgh clinical nlp at mediqa-corr 2024: Guiding large language models with hints. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Traber Davis Giardina, Shailaja Menon Helen Haskell, Julia Hallisy, Frederick S. Southwick, Urmimala Sarkar, Kathryn E. Royse, and Hardeep Singh. 2018. Learning from patients' experiences related to diagnostic errors is essential for progress in patient safety. *Health Affairs*, 37(11).

Satya Kesav Gundabathula and Sriram R Kolar. 2024. Promptmind team at mediqa-corr 2024: Improving clinical text correction with error categorization and llm ensembles. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Hyeon Hwang, Taewhoo Lee, Hyunjae Kim, and Jaewoo Kang. 2024. Ku-dmis at mediqa-corr 2024: Exploring the reasoning capabilities of small language models in medical error correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Suramya Jadhav, Abhay Shanbhag, Sumedh Joshi, Atharva Date, and Sheetal S. Sonawane. 2024. Maven at mediqa-corr 2024: Leveraging rag and medical llm for error detection and correction in medical notes. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? A large-scale open domain question answering dataset from medical exams. *CoRR*, abs/2009.13081.

Mert Karabacak and Konstantinos Margetis. 2023. Embracing large language models for medical applications: Opportunities and challenges. *Cureus*, 15(5).

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*.

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.

Yasumasa Onoe, Michael J. Q. Zhang, Eunsol Choi, and Greg Durrett. 2021. CREAK: A dataset for commonsense reasoning over entity knowledge. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.

Juan Pajaro, Edwin Puertas, David Villate, Laura Estrada, and Laura Tinjaca. 2024. Verbanexai at mediqa-corr: Efficacy of gru with biowordvec and clinicalbert in error correction in clinical notes. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Swati Rajwal, Eugene Agichtein, and Abeed Sarker. 2024. Em_mixers at mediqa-corr 2024: Knowledge-enhanced few-shot in-context learning for medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Nadia Saeed. 2024. Medifact at mediqa-corr 2024: Why ai needs a human touch. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath Rangan, and Jonathan H. Chen. 2024. Diagnostic reasoning prompts reveal the potential for large language model interpretability in medicine. *npj Digit. Medicine*, 7(1).

Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. BLEURT: learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics,*

*ACL 2020, Online, July 5-10, 2020*, pages 7881–7892. Association for Computational Linguistics.

Liyan Tang, Zhaoyi Sun, Betina Ross S. Idnay, Jordan G. Nestor, Ali Soroush, Pierre A. Elias, Ziyang Xu, Ying Ding, Greg Durrett, Justin F. Rousseau, Chunhua Weng, and Yifan Peng. 2023. Evaluating large language models on medical evidence summarization. *npj Digit. Medicine*, 6.

Writer Engineering Team. 2023. Palmyra-large parameter autoregressive language model. https://dev.writer.com.

Augustin Toma, Ronald Xie, Steven Palayew, Gary D. Bader, Patrick Lawler, and BO WANG. 2024. Wanglab at mediqa-corr 2024: Optimized llm-based programs for medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Airat Valiev and Elena Tutubalina. 2024. Hse nlp team at mediqa-corr 2024 task: In-prompt ensemble approach with named entity recognition and knowledge graph for medical error checking. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Cunxiang Wang, Shuailong Liang, Yili Jin, Yilong Wang, Xiaodan Zhu, and Yue Zhang. 2020. SemEval-2020 task 4: Commonsense validation and explanation. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 307–321, Barcelona (online). International Committee for Computational Linguistics.

Zhaolong Wu, Abul Hasan, Jinge Wu, Yunsoo Kim, Jason Pui-Yin Cheung, Teng Zhang, and Honghan Wu. 2024. Knowlab_aimed at mediqa-corr 2024: Chain-of-though (cot) prompting strategies for medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wen-wai Yim, Asma Ben Abacha, Griffin Adams, Neal Snider, and Meliha Yetisgen. 2023. Overview of the mediqa-sum task at imageclef 2023: Summarization and classification of doctor-patient conversations. In *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2023), Thessaloniki, Greece, September 18th to 21st, 2023*, volume 3497 of *CEUR Workshop Proceedings*, pages 1347–1360. CEUR-WS.org.

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.

603

# UTSA-NLP at ChemoTimelines 2024: Evaluating Instruction-Tuned Language Models for Temporal Relation Extraction

**Xingmeng Zhao** and **Anthony Rios**
Department of Information Systems and Cyber Security
The University of Texas at San Antonio
{Xingmeng.Zhao, Anthony.Rios}@utsa.edu

## Abstract

This paper presents our approach for the 2024 ChemoTimelines shared task. Specifically, we explored using Large Language Models (LLMs) for temporal relation extraction. We evaluate multiple model variations based on how the training data is used. For instance, we transform the task into a question-answering problem and use QA pairs to extract chemo-related events and their temporal relations. Next, we add all the documents to each question-answer pair as examples in our training dataset. Finally, we explore adding unlabeled data for continued pretraining. Each addition is done iteratively. Our results show that adding the document helps, but unlabeled data does not yield performance improvements, possibly because we used only 1% of the available data. Moreover, we find that instruction-tuned models still substantially underperform more traditional systems (e.g., EntityBERT).

## 1 Introduction

Extracting chemotherapy treatment timelines from clinical notes is crucial in Clinical Natural Language Processing (ClinicalNLP) for enhancing patient care and advancing cancer research (Cui et al., 2023). Researchers can construct detailed treatment timelines within Electronic Health Records (EHR) across various medical domains by identifying and extracting events related to chemotherapy treatments and their temporal information from medical documents. This work aims to develop an end-to-end system utilizing Large Language Models (LLMs) in a Question-Answer format for chemotherapy timeline extraction. Such a system will aid healthcare professionals in comprehending patient histories, thereby improving clinical text-mining efforts and assisting physicians in making more informed care decisions. Additionally, it will contribute to research in personalized cancer treatment development.

The main approach in clinical entity and relation extraction tasks heavily relies on pre-trained, domain-specific models like BioBERT (Lee et al., 2020), ClinicalBERT (Alsentzer et al., 2019), Pub-MedBERT (Gu et al., 2021), and EntityBERT (Lin et al., 2021). These models are trained on a broad range of biomedical corpora, like PubMed articles and clinical notes, to understand the complex language of the clinical domain, which is often succinct and laden with phrases, jargon, non-standard expressions, abbreviations, assumptions, and specialized knowledge. These models are then adapted or fine-tuned for specific tasks such as named entity recognition (NER), relation extraction (RE), and event extraction (EE), which often employ strategies like multi-task learning (MTL) and an all-in-one scheme to enhance performance across multiple tasks by leveraging shared knowledge and representations (Luo et al., 2023; Yadav et al., 2020). However, there are still challenges, such as a drop in performance when these models are used for out-of-domain tasks or very different sub-domains in terms of context and terminology, revealing their limitations in adaptability (Košprdić et al., 2023).

Recently, Large Language Models (LLMs) have shown remarkable potential in Natural Language Processing (NLP) tasks, including text generation, reasoning, text classification, summarization, and question answering, through their ability for zero-shot or few-shot learning (Xu et al., 2023). This capability allows them to adapt to new tasks quickly with minimal fine-tuning. This adaptability has resulted in their outstanding application performance, including NER and RE within the general domain. Models like CoT-ER (Ma et al., 2023), GPT-RE (Wan et al., 2023), and PromptNER (Ashok and Lipton, 2023) show that through few-shot learning or zero-shot learning, these generative LLMs can achieve performance levels competitive with the state-of-the-art methods in entity or relation extraction (Li et al., 2023; Brown et al., 2020; Wei et al.,

2022; Liu et al., 2023). This achievement is primarily due to their capability to use task-specific and concept-level knowledge stored during pre-training, which is then effectively leveraged through prompting to generate relevant evidence for the tasks.

However, challenges arise when adapting LLMs from the general to the medical domain, primarily due to their lack of domain-specific knowledge and the difficulty in incorporating new, relevant factual information over time (Jiang et al., 2024; Brokman and Kavuluru, 2024; Li and Zhang, 2023). While LLMs have shown potential in biomedical natural language processing tasks through innovative in-context learning strategies, their application in specific tasks like NER and RE remains problematic. This is partly because current few-shot learning methods, which trained on large amounts of source data and fine-tuning on exemplars from the target domain, do not perform well in the medical context (Gutiérrez et al., 2022; Keloth et al., 2024; Ma et al., 2023). The discrepancy arises from significant differences in entity and relationship definitions between general and medical texts (Das et al., 2022). To address these challenges, researchers have explored various approaches, including the development of domain-specific generative LLMs like BioGPT (Luo et al., 2022), BioMedLM (Bolton et al., 2024), and BioBART (Yuan et al., 2022). These models are trained from scratch using medical corpora such as PubMed or are continually fine-tuned on medical data. Basically, fine-tuning is required for adequate performance on biomedical NLP tasks. These efforts represent steps toward bridging the gap in domain adaptation for LLMs. However, updating these models for the rapidly changing medical field is still non-trivial due to the risk of catastrophic forgetting during fine-tuning (Ren et al., 2024), highlighting the need for better training methods tailored to medical knowledge.

To address this, we explored instruction-tuning methods for large language models, focusing on an open-source language model. Traditional fine-tuning methods risk forgetting previous knowledge, so we adopted a novel training strategy, gradually extending training to include associated documents and unlabeled datasets. Initially, we instruction-tuned on Question-Answer (QA) pairs before integrating complete EHR documents. Then, we trained on QA pairs and documents simultaneously. Finally, we continue pre-training on the large un-labeled corpus. Jiang et al. (2024) demonstrates that this integration strategy ensures the retention of acquired knowledge. In the inference stage, our system directly generates output relations from input questions for subtask 1. For subtask 2, we first extract event entities and time expressions before predicting relationships between identified entities using different input questions. Our approach provides an end-to-end relation extraction system for extracting Chemotherapy Treatment Timelines. This system formulates the task as a text generation task, using clinical notes as input to generate relational triplets end-to-end, without requiring additional intermediate annotations, as seen in the REBEL method (Cabot and Navigli, 2021).

In summary, this paper makes the following contributions:

1. We introduce a novel approach that combines instruction-based fine-tuning with continuous knowledge acquisition to adapt pre-trained general LLMs to the medical domain, specifically targeting the extraction of chemotherapy treatment timelines.

2. We evaluate the performance of a smaller 7B model, OpenChat-3.5-7B (Wang et al., 2023b), on extracting chemotherapy treatment timelines for breast cancer, ovarian cancer, and melanoma datasets provided by the University of Pittsburgh/UMPC. Additionally, we conduct a detailed analysis of each training component to establish a robust framework for biomedical end-to-end relation extraction, with the potential to apply the same approach to other biomedical NLP tasks.

3. We conduct an error analysis to identify the strengths and weaknesses of our proposed approach, offering insights into areas for potential improvement.

## 2 Methods

In this section, we describe our instruction-tuned LLMs strategy. Figure 1 shows a high-level overview of our approach. We convert the information extraction task to the question-answer instruction format. Our strategy has three main components: 1) Instruction-tuning LLMs on task-specific QA pairs (i.e., Named Entity Recognition (NER) and Relation Classification (RE)); 2) Joint
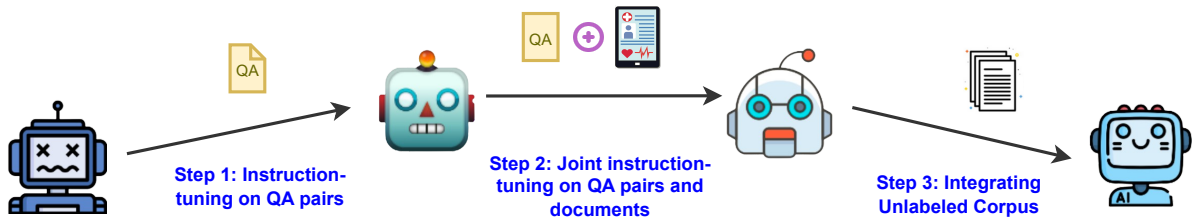
Figure 1: Overview of instruction-tuned LLMs Framework. First, we perform instruction-tuning on LLMs using task-specific QA pairs (e.g., NER and RE). Second, we conduct further instruction-tuning on QA pairs and associated documents to enhance its ability to progressively absorb knowledge from simpler to more complex data. Finally, we continue pre-training the model on an unlabeled corpus to refine its capabilities in the clinical domain further.

instruction-tuning on QA pairs and associated documents to enhance its ability to absorb knowledge progressively from simpler to more complex data; and 3) Continuing pre-training on unlabeled corpus, our intuition is first trained on QA pairs to understand knowledge access patterns, then progresses to training on a combination of QA and document data to align knowledge access through questions and knowledge encoding from documents, this will help absorb information from unlabeled data. We describe each component in the following subsections and how the three components are integrated into a unified training strategy.

## 2.1 Step 1: Instruction-tuning on QA pairs

We fine-tuned the pretrained open-source LLMs (i.e., OpenChat-3.5-7B) for two clinical tasks: classifying TLINK temporal relations and recognizing named entities, including DocTimeRel, EVENTS, and TIMEX3. Subsequent sections detail the labeled datasets used for these instruction-tuning tasks, and Figure 2 illustrates the format used for task-specific question-answer pairs.

**Relationship Classification QA Design.** For the RE QA pairs, Let $C$ represent the input context, and let $e_{\text{event}} \in C$ and $e_{\text{timex3}} \in C$ denote a chemotherapy event entity and a time expression entity, respectively. For a set of predefined relation classes $R$, the goal of relation extraction is to determine the relationship $y \in R$ between the entity pair $(e_{\text{event}}, e_{\text{timex3}})$ within $C$. If no predefined relation exists between them, the model predicts $y = \text{NULL}$. Building on the prior work (Ma et al., 2023), we use a three-step reasoning framework combining concept-level entity knowledge and explicit evidence to design question-answer instructions. This approach aims to maximize the use of knowledge embedded in LLMs to support step-by-step reasoning. For the RE task, questions are for-

mulated with instructions, definitions of potential relations, and the context. Answers are designed as a structured three-step reasoning process. First, we integrate concept-level knowledge about the event entity. Second, we apply a similar approach to the time expression entity. Third, to identify the most suitable relation label for the pair of entities within the context, we explicitly highlight relevant text spans as evidence and subsequently construct a coherent expression that combines the two entities and the relation label. An example using the relation label "CONTAINS" is shown in the last "Answer" in Figure 2. First, Avastin is described as a chemotherapy drug. Next, the TIMEX3 entity (March 2009) is described. Finally, some reasoning is described about a potential relation, and the relation is specified.

How does the model learn this reasoning framework? Inspired by Wan et al. (2023), we implement the OpenChat-3.5-7B model to generate logical reasoning in question-answer pairs. We employ few-shot learning to prompt the LLMs to generate a three-step reasoning process based on the question and corresponding given golden label. For example, we append the query "What are the three-step reasoning processes that lead to the relation between [entity1] and [entity2] being [relation] in the sentence [context]?" to the end of the question and corresponding a gold label. This prompt is then passed to the LLMs to generate the three-step reasoning. Specifically, we generate the reasons using an untrained OpenChat-3.5-7B without fine-tuning for all examples in our training dataset. These reasons are then used during our instruction tuning phase.

**Named Entity Recognition QA Design.** The NER QA instruction design is inspired by Prompt-NER (Ashok and Lipton, 2023), which shows the advantages of enhancing language models' under-

**NER System:** Given the context below, identify a list of possible entities and for each item explain why it is considered as an entity or not. The response should be structured as follows: 'entity name | entity type | True/False | Explanation', where you explain the rationale behind the classification. Output NULL and mark it as False if there is no entity identified.

**Define:** the DOCTIME entity refers to the time expression representing the document creation time, usually found at the start of the document.
**Question:** "{DOCTIME}" Given the context, the DOCTIME entity is:
**Answer:** 20090824 | DOCTIME | True | As it is listed as the "Principal Date" at the start of the document, indicating it as the date the document was created or formalized.

**Define:** The EVENT entity refers to chemotherapy mention in the clinical notes, including general terms like 'chemotherapy' and 'chemo', as well as specific chemotherapy treatments such as 'cytoxan', which involve the use of powerful drugs to target and destroy cancer cells, often administered in cycles to shrink tumors, prevent cancer spread, and potentially achieve remission or alleviate symptoms. Diseases (e.g., "melanoma"), diagnostic scans (e.g., "FDG PET scan," "CT scan") or medications not used in chemotherapy (e.g., "Vicodin" for pain relief, "Zocor" for cholesterol management) are not EVENT entities.
**Question:** "Avastin for 6 cycles through March 2009." Given the context, all relevant EVENT entities are:
**Answer:** Avastin | EVENT | True | As it is a specific type of chemotherapy treatment for breast cancer, the mention of Avastin highlights a particular therapeutic approach within the patient's care.

**Define:** The TIMEX3 entity is a time expression that identifies specific dates, times, and periods like "tomorrow" or "last week" in texts specifically related to chemotherapy mentions in the document.
**Question:** "Avastin for 6 cycles through March 2009." Given the context, all relevant TIMEX3 entities are:
**Answer:** March 2009 | TIMEX3 | True | As it tells us when the patient had chemotherapy, showing it was happened in March 2009.

**RE System:** Please solve the TLINK Relation Extraction task, which is a temporal link only between an EVENT and a TIMEX3. Given the context below, consider what is the most precise relation between two entities belongs to the following N possible relations. The relation to choose must be in these N possible relations: CONTAINS, CONTAINS-1, BEGINS-ON, ENDS-ON. Please output NULL if the EVENT and TIMEX3 entities do not appear in the same sentence or if none of these relations apply.

**Question:** "Avastin for 6 cycles through March 2009." Given the context, what are the relations between the EVENT entity '{source_entity}' and the TIMEX3 entity '{target_entity}'?
**Answer:** 1. EVENT entity 'Avastin' is a chemotherapy drug mentioned in the context as part of a treatment regimen that includes carboplatin and paclitaxel, indicating its use in the patient's cancer therapy.
2. TIMEX3 entity 'March 2009' is a time expression indicating the endpoint of the period during which the chemotherapy treatment regimen, including 'Avastin', was administered for 6 cycles.
3. According to the context, the phrase "Avastin for 6 cycles through March 2009" suggests that the chemotherapy treatment phase encompassing 'Avastin' spans until 'April 2013', denoting the EVENT 'Avastin' contains the TIMEX3 'March 2009' within its duration.
Therefore, Avastin | EVENT | March 2009 | TIMEX3 | CONTAINS

**Doc Question:** The relevant document is:
**Answer:** {document}

Figure 2: An instruction example for clinical document and task-specific QA pairs. Both subtask 1 and subtask 2 use the same training dataset and process. However, subtask 1 focuses on identifying temporal relations by generating specific relation pairs through tailored questions during inference. In contrast, subtask 2 first identifies chemo-related entities with distinct instructions before determining their temporal relationships. Tokens used for computing losses are highlighted in green.

standing of textual logical structures. This understanding is used to improve NER tasks by employing Chain-of-Thought Prompting, guiding the model through a step-by-step reasoning process that leads to entity identification. This technique boosts entity recognition accuracy and offers a versatile framework adaptable to various entity types by adjusting definitions and explanations within the prompting template (Ashok and Lipton, 2023; Wang et al., 2023a). Therefore, in our NER QA instruction design, each question includes instructions and definitions of entities, with answers detailing the chosen entities in the format of "entity name | entity type | True/False | Explanation," where the Explanation includes the rationale behind the NER

type classification. Inspired by Ashok and Lipton (2023), this method employs Chain-of-Thought Prompting to refine our model's understanding of textual logic, enhancing NER tasks by guiding step-by-step reasoning. We've crafted a structured output template for the LLMs that identifies and classifies entities. This structure has the potential to enhance accuracy through outcome supervision using reinforcement learning (Gao et al., 2024). Additionally, the True/False component marks noun phrases that are relevant entities we want to extract (True) or irrelevant (False). In our experiments, we learn to generate relevant entities because we are fine-tuning, hence we only use True. However, we kept the option for False in future work by adding

incorrect entities.

This format displays the model's decision-making process, making it adaptable across different NER types by simply modifying definitions. Similarly, we use the non-finetuned OpenChat-3.5-7B model, employing few-shot learning with manually created demonstrations to generate explanations for all examples in the training data. In general, our NER QA instruction includes three distinct categories of entities: EVENTS, which refer specifically to mentions of chemotherapy treatments; DocTimeRel, which represents the temporal relationship between an event and the time the document was created; and Temporal Expressions (TIMEX3), which are precise references to times linked to chemotherapy treatments. These entities are illustrated in Figure 2, which shows "Avastin" and "March 2009" as example extractions.

## 2.2 Step 2: Joint instruction-tuning on QA pairs and documents

In this training phase, the instruction combines QA pairs with their relevant documents. Intuitively, QA pairs are typically simple, unlike documents, which are usually more complex and dense, containing numerous factual details not available in a single (or few) sentence. Therefore, Jiang et al. (2024) suggests that it is effective to deliberately expose LLMs to QA data before continued pre-training on documents so that the process of encoding knowledge from complex documents considers how this knowledge is accessed through questions. During this phase, LLMs improve at digesting detailed content from documents, building on the QA pairs they've already learned. The training starts with QA pairs to grasp basic knowledge access patterns and then adds documents to enhance question-based knowledge access and document understanding. The instruction is created based on each document; we position all the NER QA pairs, followed by the RE QA pairs. Finally, the document itself is formatted as a QA pair, with the question identifying the document and the answer being the document's content, as illustrated in Figure 2. Jiang et al. (2024) found that placing the documents after the QA pairs leads to better performance than placing them before. We also experimented with positioning the document before and after the QA pairs and tested on the melanoma development set. The results showed that placing the document after the QA pairs yielded better per-

formance. Therefore, we put the document after the QA pairs in our following experiments.

## 2.3 Step 3: Integrating Unlabeled Corpus

In this training phase, we aim to improve how the fine-tuned OpenChat-3.5-7B model handles clinical documents, which are often complex and full of medical terminology. Instead of using instruction-tuning alone, we continued "pre-training" the model on unlabeled documents (i.e., training on unlabeled data after instruction-tuning).[1] This potentially helps the model learn a specialized vocabulary for the clinical domain, capturing important terms such as diseases, symptoms, medications, and medical procedures in their original context (Lin et al., 2021). This approach is crucial for enhancing the model's performance on tasks specific to the clinical field. Based on Jiang et al. (2024), there's a concern that directly continuing pre-training on a vast, unlabeled clinical corpus might lead to the model forgetting previously acquired knowledge. However, by initially training on QA pairs to grasp knowledge access patterns and then moving on to a blend of QA and document data, we can strengthen the model's ability to assimilate document knowledge. This method helps mitigate the issue of catastrophic forgetting by aligning how the model accesses knowledge through questions with how it encodes knowledge from documents (Ouyang et al., 2022; Jiang et al., 2024). Technically, we employed Byte-Pair Encoding (BPE) (Gage, 1994) to break down the text into small context windows, considering the OpenChat-3.5-7B model's 8192 token maximum context limit, setting our windows to 7800 tokens for efficiency. We prepared the training data by joining these pieces with an end-of-sequence (eos) token and then splitting the extended text into sections. This structured training method is designed to make the model more effective at analyzing and interpreting medical documents.

## 3 Experiments

In this section, we provide a brief overview of the dataset, discuss the evaluation metrics, discuss our results on the validation dataset, and briefly mention the final model performance in the competition on the test set.

---

[1]Because of lack of time and limited GPU resources, we were not able to use the entire unlabeled dataset and only learned on less than 1% of the unlabeled data.

## 3.1 Dataset

In this shared task, we use both unlabeled and labeled EHRs, including radiology reports, pathology notes, clinical notes, oncology notes, discharge summaries, and progress reports, from the University of Pittsburgh/UPMC to construct the end-to-end system for Extracting Chemotherapy Treatment Timelines. For the unlabeled data, this included EHR notes from approximately 62,000 patients with breast and ovarian cancer and 16,000 patients with melanoma. For the labeled data, we have gold annotations for 310 patients' histories, focusing on EVENTs, TIMEX3s entities, and temporal relations (TLINKs) between an EVENT and a TIMEX3. The training set includes EHRs for ovarian (26 patients), breast (33 patients), and melanoma (10 patients), while the development set comprises records from ovarian (8 patients), breast (16 patients), and melanoma (3 patients). Additionally, for ethical reasons and to protect patient privacy, the data has been de-identified (Jiarui Yao, 2024).

An EVENT refers to any relevant chemotherapy treatment on the clinical timeline, each with a temporal relation to the document creation time (DocTimeRel), categorized as BEFORE, BEFORE-OVERLAP, OVERLAP, or AFTER. Temporal expressions (TIMEX3) denote discrete references to time, normalizations to a unified format (e.g., "YYYY-MM-DD") using TimeNorm (Laparra et al., 2018; Xu et al., 2019). Additionally, temporal relations (TLINKs) link an EVENT and TIMEX3, including categories such as CONTAINS, CONTAINS-1, BEFORE, BEGINS-ON, and ENDS-ON, where CONTAINS-1 is the inverse of CONTAINS, meaning the Target CONTAINS the Source (Styler IV et al., 2014).

For training, we created positive NER QA pairs from all gold standard examples, even though there were no relations between EVENT and TIMEX3. For RE QA pairs, we randomly selected three pairs of chemo events and time expressions with no temporal relation, where the answer would be NULL.

## 3.2 Hyperparameters

In our experiments, we trained models on 2 Nvidia A6000 GPUs using DeepSpeed Zero stage 2 (Rasley et al., 2020), HuggingFace Accelerate (Gugger et al., 2022), and FlashAttention2 (Dao, 2023) for a maximum of 10 epochs and using

---

[1]https://github.com/clulab/timenorm

Melanoma dev set to select best epoch for all three stage training. Following Jiang et al. (2024), we employed the AdamW optimizer (Loshchilov and Hutter, 2018) with specific parameters ($\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay = 0.1) and set a maximum context length of 1024.

For instruction tuning on question-answer pairs, we used a batch size of 128 and learning rates of 3e-5 for direct pairs and 5e-6 when documents were associated while continuing pre-training on unlabeled datasets at a batch size of 36 and a learning rate of 3e-5. We use spaCy's "en_ner_bc5cdr_md" model for sentence boundary detection and text segmentation. Moreover, we adopted Low-Rank Adaptation (LoRA) fine-tuning (Hu et al., 2021) with a rank of 256, LoRA alpha of 512, and LoRA dropout of 0.05, targeting modules ["q_proj", "o_proj", "k_proj", "v_proj", "gate_proj", "up_proj", "down_proj", "fc_in", "fc_out", "wte"], to optimize specific target modules within pre-trained language models (LLMs), effectively reducing the number of parameters needed for training without altering the original model weights. This approach was facilitated by using the "trl" library from HuggingFace (von Werra et al., 2020), enhancing our model's performance and efficiency.

When training on QA pairs, we compute the average negative log-likelihood loss by focusing only on the tokens within the answer. This approach is inspired by Lin et al. (2024), which suggests that not all tokens are equally important in language model training. We can enhance the model's efficiency and performance by selectively focusing on tokens that align with the desired distribution. For QA + Doc training, we treat the phrase "The relevant document is" as a question and apply next-token prediction loss to the document's tokens, treating them as an expanded answer. This is because the document provides a rich context that informs the model's understanding, enabling it to learn from contextually relevant tokens, as shown in Figure 2.

In the inference stage, we experimented with different settings with temperatures from 0.1 to 0.9, top p values from 0.1 to 0.6, and top k options of 10, 20, and 30. After experimenting, we found that the best settings were a temperature of 0.2, a top p value of 0.5, and a top k of 20.

### 3.3 Evaluation Metrics

The final output of our system employs the following approach to summarize event triples into patient-level timelines: We begin by using gold-standard DOCTIME annotations for subtask 1. In subtask 2, we predict DOCTIME by analyzing the first sentence of each document, discarding any document that lacks a DOCTIME prediction. Next, we normalize all temporal expressions to a standard format using the TimeNorm package (Laparra et al., 2018; Xu et al., 2019), with DOCTIME as the anchor time. We then de-duplicate timeline entries where chemotherapy events, time expressions, and their relations are identical. Using the timeline summarization system described by Jiarui Yao (2024), we prioritize specific temporal labels from a predefined hierarchy (e.g., BEGINS-ON/ENDS-ON → CONTAINS) for chemotherapy events and only include generic terms like "chemotherapy" in the timeline if there is no mention of a specific drug like "cytoxan" on the same day with the same label.

Performance in this shared task is measured by comparing generated patient-level timelines against gold-standard timelines. Specifically, we evaluate the accuracy of identified tuples containing chemotherapy events, their temporal relations, and time expressions ("chemo EVENT", "temporal_relation", "TIMEX3") compared to the correct timelines. The F1 score is calculated for each patient and then averaged across all patients to yield the macro F1 score. This evaluation employs a relaxed criterion, acknowledging certain temporal relations, specifically "contains-1" with "begins-on" and "contains-1" with "ends-on", as equivalent (Jiarui Yao, 2024).

### 3.4 Results

In the inference stage, for subtask 1, we directly fed questions to the model to generate output relations. For subtask 2, the model processes each sentence first to extract the chemo event entity. Inspired by Cui et al. (2023), we adopt a sentence window approach to extract associated time expressions. If the target treatment entity is within the target sentence, the model selects $k$ sentence before and after the target sentence to gather contextual information. Due to constraints in time and computing resources, we initially set the window size to zero. If an event entity is detected, we extract the time expression by reprocessing the sentence through the model. Furthermore, to enhance accuracy for subtask 2,

we implemented rule-based postprocessing. This approach uses regular expressions to identify and remove inaccurate named entity recognition (NER) predictions for EVENTS and TIMEX3, specifically targeting the pattern associated with chemo entities.

Table 1 shows the official results on the dev set for subtask 1. Our best performance is achieved when instruction tuning with QA and associated documents, leading to a slight accuracy improvement across all disease types, with an overall average score of .68. This indicates the benefit of integrating document context into our training regimen. However, we observed a slight decrease in performance for all three disease types when we continued pretraining on the unlabeled dataset. This decline may be attributed to the limited usage of training data, as we only utilized 1% of the unlabeled data. This did not fully explore the potential of continuous training capabilities, possibly explaining the observed performance dip. Further exploration and more extensive use of the unlabeled data might be necessary to fully optimize the model's performance.

Table 2 shows the official results on the dev set for subtask 2. The model shows variable performance across cancer types, struggling notably with ovarian cancer (.17) and achieving a total average precision of .47. This suggests that subtask 2's entity extraction and relation task is more challenging, especially in complex cancer data.

Table 3 shows the official results on the test set for subtask1. Our method ranks in the mid-tier compared to other teams, with a total average precision of .69. This indicates our approach's competitiveness but also highlights a gap to top-performing models and the baseline.

Table 4 shows the official results on the test set for subtask 2. We face significant challenges, with a total average precision of .22, considerably lower than the baseline. This underscores the complexity of subtask 2 and the need for method improvement.

Overall, our method employs generative LLMs, which, despite their innovative approach, encounter difficulties when competing against traditional state-of-the-art (SOTA) BERT methods in specific tasks like NER and RE. The broad capabilities of generative models aimed at creating new content may not directly translate to the high specificity required for these tasks in the medical domain. This discrepancy is evident in our performance on dev

| | Breast | | | Melanoma | | | Ovarian | | | Total Average |
|---|---|---|---|---|---|---|---|---|---|---|
| | Type A | Type B | Average | Type A | Type B | Average | Type A | Type B | Average | |
| train QA | .81 | .50 | .66 | .80 | .70 | .75 | .57 | .57 | .57 | .66 |
| + train QA + DOC | .82 | .51 | .67 | .83 | .74 | .78 | .58 | .58 | .58 | .68 |
| + train on unlabeled corpus | .77 | .39 | .58 | .80 | .70 | .75 | .56 | .56 | .56 | .63 |

Table 1: Official results on the dev set for subtask 1.

| | Breast | | | Melanoma | | | Ovarian | | | Total Average |
|---|---|---|---|---|---|---|---|---|---|---|
| | Type A | Type B | Average | Type A | Type B | Average | Type A | Type B | Average | |
| train QA + Doc | .78 | .41 | .59 | .71 | .57 | .64 | .17 | .17 | .17 | .47 |

Table 2: Official results on the dev set for subtask 2.

| Team | Breast | Melanoma | Ovarian | Total Average |
|---|---|---|---|---|
| LAILab 1 | .96 | .87 | .88 | .90 |
| Wonder 2 | .90 | .84 | .77 | .84 |
| NLPeers 1 | .72 | .81 | .75 | .77 |
| BioCom 1 | .88 | .61 | .72 | .74 |
| Lexicans 1 | .68 | .83 | .61 | .71 |
| UTSA-NLP 1 (Ours) | .70 | .68 | .69 | .69 |
| EmoryClincalRXMiners 1 | .44 | .47 | .34 | .40 |
| Baseline | .93 | .87 | .88 | .89 |

Table 3: Official results on the test set for subtask 1.

| Model | Breast | Melanoma | Ovarian | Total Average |
|---|---|---|---|---|
| LAILab 2 | .62 | .74 | .74 | .70 |
| KCLab 1 | .68 | .49 | .45 | .54 |
| Wonder 3 | .63 | .39 | .55 | .53 |
| NYULangone | .19 | .32 | .18 | .23 |
| UTSA-NLP (Ours) | .25 | .21 | .18 | .22 |
| Baseline | .59 | .43 | .71 | .58 |

Table 4: Official results on the test set for subtask 2.

and test sets, especially for subtask 2, where our approach trails behind the baseline model built based on EntityBERT (Lin et al., 2021). This outcome suggests that leveraging the strengths of generative models for such specific tasks requires a strategic reevaluation of our model's application or methodology.

### 3.5 Error Analysis

Our error analysis shows that the model is prone to generating false positive relation triples. This issue appears to be rooted in insufficient NULL relation examples during training, leading to the model's poor performance in recognizing the absence of a relationship between EVENT and TIMEX3 entities.

> **"Gemcitabine used in August 2010 and cisplatin used from March 2012."**

For instance, in the above sentence:[2] "Gemcitabine used in August 2010 and cisplatin used from March 2012." In this case, two chemotherapy treatment events are linked with specific time expressions. Our approach to relation extraction involves testing every possible combination of EVENT and TIMEX3 entities, such as Gemcitabine with August 2010, Gemcitabine with March 2012, cisplatin with August 2010, and cisplatin with March 2012. Notably, the combinations of Gemcitabine with March 2012 and cisplatin with August 2010 do not have a temporal relation. Nevertheless, our model erroneously predicts a relation for these pairs. This flaw is primarily due to the difficulty in generating high-quality negative examples for creating QA pairs, which is essential for accurately predicting a NULL relationship.

In subtask 2, we also need to identify EVENT entities accurately. However, generative language models (LLMs) struggle with this, often misidentifying unrelated entities as EVENTS. These errors include categorizing diseases (like "melanoma" or "Parkinson"), diagnostic scans ("FDG PET scan," "CT scan"), diagnostic codes ("PD13-007285PD"), people ("Person2"), and non-chemotherapy medications ("Vicodin," "Zocor") as EVENT entities, despite instructions to exclude them. To address these inaccuracies, we use regular expressions to filter and refine our EVENT entity identification, based on a list of valid chemotherapy events extracted from the training and development sets. This use of regular expressions as a post-processing step ensures the exclusion of these inaccurately named entities.

---

[2]All examples have been modified and do not directly match the training data to ensure data privacy.

> "Patient underwent diagnostic CT scans in June 2012 ."

For example, when analyzing the sentence "Patient underwent diagnostic CT scans in June 2012," our model incorrectly classifies "diagnostic CT scans" as a chemotherapy EVENT. Although the model explains that "diagnostic CT scans | EVENT | True | As it is crucial for diagnosing the disease and planning chemotherapy," meaning CT scans are important for diagnosis, not chemotherapy events, the model still wrongly labels them as EVENT entities. This leads to many false positives in identifying entities.

## 4  Related Work

**Continual Knowledge Acquisition.**  In continual knowledge acquisition, several studies have investigated the ability of language models (LMs) to retain and update knowledge over time. Hu et al. (2023) and Ovadia et al. (2023) explore the effectiveness of different pre-training approaches using smaller LMs like BART (Lewis et al., 2020) and EntityBERT (Lin et al., 2021). Zhu and Li (2023); Jiang et al. (2024); Keloth et al. (2024) delve into fine-tuning LMs on QA pairs related to individuals, with a focus on mixed training settings combining biographies and QA pairs. These studies are a foundation for exploring strategies to incorporate QA data before continued pre-training. Additionally, researchers have sought to adapt LMs to specialized domains, such as medicine, with Li and Zhang (2023); Hu et al. (2024); Zhang et al. (2023) proposing various strategies. However, a common challenge in continual knowledge acquisition is the potential for inaccuracies or difficulties in clinical NLP tasks. Models like BioGPT (Luo et al., 2022), BioMedLM (Bolton et al., 2024), and BioBART (Yuan et al., 2022) address these concerns by continuing training specifically within the medical domain.

**Instruction Fine-tuning.**  Recently, instruction tuning, also known as supervised fine-tuning, has gained prominence for its ability to draw out knowledge from Large Language Models (LLMs) using high-quality annotated data or data from proprietary models (Wei et al., 2021; Zhou et al., 2024; Brokman and Kavuluru, 2024; Zhou et al., 2023). This process enhances LLMs' capacity to address user inquiries and improves their factual accuracy, a focal point of our research. Additionally, the zero-shot and few-shot in-context learning capabilities of LLMs, which operate with minimal or no training data, present a significant advantage for efficient learning. These approaches, further discussed by Wei et al. (2021) and highlighted in the works of Wang et al. (2024) and Sanh et al. (2021), underscore the potential of instruction tuning in refining LLMs' factuality and responsiveness.

## 5  Limitation

Due to the constrained timeline and limited resources of the shared task, our exploration was restricted to basic setups. We did not create negative examples for NER QA pairs and only used a limited set of negative examples for RE QA pairs by randomly selecting three unrelated pairs of chemotherapy events and time expressions. Additionally, our limited use of just 1% of the unlabeled dataset resulted in decreased performance across all three disease types, suggesting that we didn't fully exploit the continuous training capabilities.

Furthermore, our experiments only considered entities within the same sentence, overlooking cases where entities span multiple sentences in the ChemoTimelines dataset. This oversight could significantly impact model performance evaluation. NER and RE tasks are sensitive to prompt design, and our initial single prompt strategy may not have been optimal. More comprehensive training and experiments, including ablation tests, will be necessary to evaluate and enhance our system's performance and efficiency thoroughly.

## 6  Conclusion and Future Work

This paper presents our end-to-end system for extracting Chemotherapy Treatment Timelines from the Clinical NLP ChemoTimelines share the task. We explored various instruction tuning strategies for open-source generative LLMs, providing a starting point for developing NER and RE models in the medical domain. Our future work will explore the implementation of outcome supervision and process-based reward mechanisms in reinforcement learning training to address the issue of false positive predictions (Gao et al., 2024).

# References

Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. Publicly available clinical bert embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78.

Dhananjay Ashok and Zachary C Lipton. 2023. Promptner: Prompting for named entity recognition. *arXiv preprint arXiv:2305.15444*.

Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, et al. 2024. Biomedlm: A 2.7 b parameter language model trained on biomedical text. *arXiv preprint arXiv:2403.18421*.

Aviv Brokman and Ramakanth Kavuluru. 2024. How important is domain specificity in language models and instruction finetuning for biomedical relation extraction? *arXiv e-prints*, pages arXiv–2402.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. Rebel: Relation extraction by end-to-end language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370–2381.

Yang Cui, Lifeng Han, and Goran Nenadic. 2023. Medtem2. 0: Prompt-based temporal classification of treatment events from discharge summaries. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*, pages 160–183.

Tri Dao. 2023. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations*.

Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca J Passonneau, and Rui Zhang. 2022. Container: Few-shot named entity recognition via contrastive learning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6338–6353.

Philip Gage. 1994. A new algorithm for data compression. *The C Users Journal*, 12(2):23–38.

Jun Gao, Huan Zhao, Wei Wang, Changlong Yu, and Ruifeng Xu. 2024. Eventrl: Enhancing event extraction with outcome supervision for large language models. *arXiv preprint arXiv:2402.11430*.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23.

Sylvain Gugger, L Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, M Sun, and B Bossan. 2022. Accelerate: Training and inference at scale made simple, efficient and adaptable.

Bernal Jiménez Gutiérrez, Nikolas McNeal, Clayton Washington, You Chen, Lang Li, Huan Sun, and Yu Su. 2022. Thinking about gpt-3 in-context learning for biomedical ie? think again. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4497–4512.

Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Nathan Hu, Eric Mitchell, Christopher D Manning, and Chelsea Finn. 2023. Meta-learning online adaptation of language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4418–4432.

Yan Hu, Qingyu Chen, Jingcheng Du, Xueqing Peng, Vipina Kuttichi Keloth, Xu Zuo, Yujia Zhou, Zehan Li, Xiaoqian Jiang, Zhiyong Lu, et al. 2024. Improving large language models for clinical named entity recognition via prompt engineering. *Journal of the American Medical Informatics Association*, page ocad259.

Zhengbao Jiang, Zhiqing Sun, Weijia Shi, Pedro Rodriguez, Chunting Zhou, Graham Neubig, Xi Victoria Lin, Wen-tau Yih, and Srinivasan Iyer. 2024. Instruction-tuned language models are better knowledge learners. *arXiv preprint arXiv:2402.12847*.

WonJin Yoon Eli Goldner Guergana Savova Jiarui Yao, Harry Hochheiser. 2024. Overview of the 2024 shared task on chemotherapy treatment timeline extraction.

Vipina K Keloth, Yan Hu, Qianqian Xie, Xueqing Peng, Yan Wang, Andrew Zheng, Melih Selek, Kalpana Raja, Chih Hsuan Wei, Qiao Jin, et al. 2024. Advancing entity recognition in biomedicine via instruction tuning of large language models. *Bioinformatics*, page btae163.

Miloš Košprdić, Nikola Prodanović, Adela Ljajić, Bojana Bašaragin, and Nikola Milošević. 2023. A transformer-based method for zero and few-shot biomedical named entity recognition. *arXiv preprint arXiv:2305.04928*.

Egoitz Laparra, Dongfang Xu, and Steven Bethard. 2018. From characters to time intervals: New paradigms for evaluation and neural parsing of time normalizations. *Transactions of the Association for Computational Linguistics*, 6:343–356.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.

Guozheng Li, Peng Wang, and Wenjun Ke. 2023. Revisiting large language models as zero-shot relation extractors. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6877–6892.

Mingchen Li and Rui Zhang. 2023. How far is language model from 100% few-shot named entity recognition in medical domain. *arXiv preprint arXiv:2307.00186*.

Chen Lin, Timothy Miller, Dmitriy Dligach, Steven Bethard, and Guergana Savova. 2021. Entitybert: Entity-centric masking strategy for model pretraining for the clinical domain. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 191–201.

Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao, Nan Duan, et al. 2024. Rho-1: Not all tokens are what you need. *arXiv preprint arXiv:2404.07965*.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35.

Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Ling Luo, Chih-Hsuan Wei, Po-Ting Lai, Robert Leaman, Qingyu Chen, and Zhiyong Lu. 2023. Aioner: all-in-one scheme-based biomedical named entity recognition using deep learning. *Bioinformatics*, 39(5):btad310.

Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 23(6):bbac409.

Xilai Ma, Jing Li, and Min Zhang. 2023. Chain of thought with explicit evidence reasoning for few-shot relation extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2334–2352.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. 2023. Fine-tuning or retrieval? comparing knowledge injection in llms. *arXiv preprint arXiv:2312.05934*.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.

Weijieying Ren, Xinlong Li, Lei Wang, Tianxiang Zhao, and Wei Qin. 2024. Analyzing and reducing catastrophic forgetting in parameter efficient tuning. *arXiv preprint arXiv:2402.18865*.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. 2021. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.

William F Styler IV, Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, et al. 2014. Temporal annotation in the clinical domain. *Transactions of the Association for Computational Linguistics*, 2:143–154.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl.

Zhen Wan, Fei Cheng, Zhuoyuan Mao, Qianying Liu, Haiyue Song, Jiwei Li, and Sadao Kurohashi. 2023. Gpt-re: In-context learning for relation extraction using large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3534–3547.

Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023a. Towards understanding chain-of-thought prompting: An empirical study of what matters. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2717–2739.

Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2023b. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*.

Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. 2024. How far can camels go? exploring the state of instruction tuning on open resources. *Advances in Neural Information Processing Systems*, 36.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Jason Wei, 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.

Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, and Enhong Chen. 2023. Large language models for generative information extraction: A survey. *arXiv preprint arXiv:2312.17617*.

Dongfang Xu, Egoitz Laparra, and Steven Bethard. 2019. Pre-trained contextualized character embeddings lead to major improvements in time normalization: A detailed analysis. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (* SEM 2019)*, pages 68–74.

Hemant Yadav, Sreyan Ghosh, Yi Yu, and Rajiv Ratn Shah. 2020. End-to-end named entity recognition from english speech. *arXiv preprint arXiv:2005.11184*.

Hongyi Yuan, Zheng Yuan, Ruyi Gan, Jiaxing Zhang, Yutao Xie, and Sheng Yu. 2022. Biobart: Pretraining and evaluation of a biomedical generative language model. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 97–109.

Kai Zhang, Bernal Jiménez Gutiérrez, and Yu Su. 2023. Aligning instruction tasks unlocks large language models as zero-shot relation extractors. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 794–812.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

Wenxuan Zhou, Sheng Zhang, Yu Gu, Muhao Chen, and Hoifung Poon. 2023. Universalner: Targeted distillation from large language models for open named entity recognition. In *The Twelfth International Conference on Learning Representations*.

Zeyuan Allen Zhu and Yuanzhi Li. 2023. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*.

This is a section in the appendix.

# WangLab at MEDIQA-CORR 2024: Optimized LLM-based Programs for Medical Error Detection and Correction

**Augustin Toma[1,2] Ronald Xie[1,2] Steven Palayew[1,2] Patrick R. Lawler[1,3] Bo Wang[1,2,4,5]**

[1]University of Toronto    [2]Vector Institute

[3]McGill University    [4]Peter Munk Cardiac Centre, University Health Network

[5]AI Hub, University Health Network

{augustin.toma, ronald.xie, steven.palayew}@mail.utoronto.ca
bowang@vectorinstitute.ai

## Abstract

Medical errors in clinical text pose significant risks to patient safety. The MEDIQA-CORR 2024 shared task focuses on detecting and correcting these errors across three subtasks: identifying the presence of an error, extracting the erroneous sentence, and generating a corrected sentence. In this paper, we present our approach that achieved top performance in all three subtasks. For the MS dataset, which contains subtle errors, we developed a retrieval-based system leveraging external medical question-answering datasets. For the UW dataset, reflecting more realistic clinical notes, we created a pipeline of modules to detect, localize, and correct errors. Both approaches utilized the DSPy framework for optimizing prompts and few-shot examples in large language model (LLM) based programs. Our results demonstrate the effectiveness of LLM based programs for medical error correction. However, our approach has limitations in addressing the full diversity of potential errors in medical documentation. We discuss the implications of our work and highlight future research directions to advance the robustness and applicability of medical error detection and correction systems.

## 1 Introduction

Medical errors pose a significant threat to patient safety and can have severe consequences, including increased morbidity, mortality, and healthcare costs. Detecting and correcting these errors in clinical text is crucial for ensuring accurate medical documentation and facilitating effective communication among healthcare professionals. One of the fastest-growing use cases for artificial intelligence (AI) in healthcare is clinical note generation, often from transcriptions of physician-patient dialogues. However, assessing the quality and accuracy of these notes is challenging, and automated detection and correction of errors could have a significant impact on patient care. The reliability of large language models (LLMs) in critical applications, such as healthcare, is a major concern due to the potential for hallucinations (generating false or nonsensical information) and inconsistencies. Robust solutions to the question of error detection and correction are essential for addressing these concerns and enabling the safe and effective use of LLMs in medical contexts.

The MEDIQA-CORR 2024 (Ben Abacha et al., 2024a) shared task focuses on identifying and correcting medical errors in clinical notes. Each text is either correct or contains a single error. The task involves three subtasks: (1) detecting the presence of an error, (2) extracting the erroneous sentence, and (3) generating a corrected sentence for flagged texts.

In this paper, we present our approach, which achieved the top performance across all three subtasks in the MEDIQA-CORR 2024 competition. We develop a series of LLM-based programs using DSPy, a framework for optimizing prompts and few-shot examples. We provide a detailed description of our methodology and results, followed by a discussion of the implications of our work and future directions in the field of medical error detection and correction.

## 2 Related Work

The use of large language models (LLMs) in medicine has attracted considerable attention in recent years. The release of LLMs such as GPT-4 has led to intensive research in the medical community (Nori et al., 2023), particularly in clinical note generation. The MEDIQA-Chat 2023 (Ben Abacha et al., 2023) competition showcased the performance of automated note generation solutions (Giorgi et al., 2023), and further work has demonstrated that LLMs can sometimes outperform humans on clinical text summarization tasks

(Van Veen et al., 2024).

However, there has been limited research focusing on granular audits of these clinical notes with respect to accuracy and error correction. The MEDIQA-CORR 2024 shared task addresses this gap by providing a platform for researchers to develop and evaluate novel approaches to error detection and correction in clinical text, ultimately contributing to the development of more reliable AI systems in healthcare.

## 3 Task Description

The MEDIQA-CORR 2024 shared task provides two distinct datasets: MS and UW (Ben Abacha et al., 2024b). The MS dataset consists of a Training Set containing 2,189 clinical texts and a Validation Set (#1) containing 574 clinical texts. The UW dataset, on the other hand, consists solely of a Validation Set (#2) containing 160 clinical texts. The test set for the shared task includes clinical texts from both the MS and UW collections.

The evaluation metrics for the MEDIQA-CORR 2024 shared task vary across the three subtasks:

- Subtask 1 (Error Flag Prediction): Evaluated using Accuracy.

- Subtask 2 (Error Sentence Detection): Evaluated using Accuracy.

- Subtask 3 (Sentence Correction): Evaluated using ROUGE (Lin, 2004), BERTScore (Zhang et al., 2020), BLEURT (Sellam et al., 2020), Aggregate-Score (mean of ROUGE-1-F, BERTScore, BLEURT-20), and Composite Scores.

The Composite Score for each text in Subtask 3 is calculated as follows:

1. Assign 1 point if both the system correction and the reference correction are "NA"

2. Assign 0 points if only one of the system correction or the reference correction is "NA"

3. Calculate the score based on metrics (ROUGE, BERTScore, BLEURT and the Aggregate-Score) within the range of [0, 1] if both the system correction and reference correction are non-"NA" sentences.

## 4 Approach

### 4.1 Overview

Upon reviewing the MS and UW datasets, it became apparent that these two datasets presented distinct challenges. The errors in the MS dataset were often extremely subtle, to the point that many errors did not actually seem like errors, and in fact, clinicians on our team often couldn't identify the presence of an error within the text. However, when reviewing corrected text from the training set, it became clear that corrections were often 'optimal' completions. For example, consider the following error and its correction:

> **Error sentence:** After reviewing imaging, the causal pathogen was determined to be Haemophilus influenzae. (Ben Abacha et al., 2024b)

> **Corrected sentence:** After reviewing imaging, the causal pathogen was determined to be Streptococcus pneumoniae. (Ben Abacha et al., 2024b)

These types of errors are subtle and seem akin to multiple-choice questions, where often multiple answers could independently be seen as correct completions, but only in the context of one another would you deem one answer wrong. On the other hand, the UW dataset appeared to reflect realistic clinical notes, and the errors were more apparent. For example, consider the following error and its correction:

> **Error sentence:** Hypokalemia - based on laboratory findings patient has hypervalinemia. (Ben Abacha et al., 2024b)

> **Corrected sentence:** Hypokalemia - based on laboratory findings patient has hypokalemia. (Ben Abacha et al., 2024b)

In this case, the error involves a nonsensical term (hypervalinemia, a rare metabolic condition) when the context makes it clear that the patient has hypokalemia (low potassium levels). These are errors that a clinician can identify from the text alone.

The distinct characteristics of the MS and UW datasets prompted us to develop a two-pronged approach to the MEDIQA-CORR 2024 shared task. For the MS dataset, we employed a retrieval-based system to identify similar questions from external medical question-answering datasets and leverage the knowledge contained in these datasets to detect

and correct errors. For the UW dataset, we created a series of modules to detect, localize, and correct errors in clinical text snippets. Both approaches were built on DSPy ([Khattab et al., 2023](#)), a novel framework for systematically optimizing prompts and few-shot examples in LLM based programs.

## 4.2 Approach for MS Dataset

Our approach to the MS dataset involves a multi-step process that leverages retrieval-based methods and the DSPy framework, as illustrated in Figures [1](#), [2](#), and [3](#). In all of our experiments, we utilized GPT-4-0125-preview as the underlying large language model, using default generation parameters (temperature of 1.0, top_p of 1) with the exception of a max tokens value of 4096.

### 4.2.1 Retrieval of Similar Questions

First, we employ a retrieval-based approach to identify similar questions from the MedQA dataset ([Jin et al., 2020](#)). MedQA is a medical question-answering dataset that contains multiple-choice questions, each with a set of answer options and a correct answer. By leveraging the knowledge contained in this external dataset, we aim to detect and correct errors in the MS dataset. We use TF-IDF ([Sparck Jones, 1972](#)) to calculate the similarity between the given question in the MS dataset and the questions in MedQA, retrieving the most similar questions along with their answer options and correct answers for further analysis.

### 4.2.2 Identifying Answer Choices within Query Text

To identify the implicit answer choice within the query text, we employ a two-step process using DSPy programs. First, we send both the query text and the identified similar multiple-choice question to a DSPy module that utilizes chain of thought ([Wei et al., 2023](#)) and the BootstrapFewShotWithRandomSearch teleprompter ([Khattab et al., 2023](#)). This teleprompter generates 20 few-shot examples by sampling from the training set and testing the module's performance on the validation set. The module aims to extract the answer choice that appears to be present in the query text.

The output from this module is then passed to a second DSPy module, which also leverages the BootstrapFewShotWithRandomSearch teleprompter. This module creates multiple few-shot examples that compare the extracted answer against the true answer from the multiple-choice
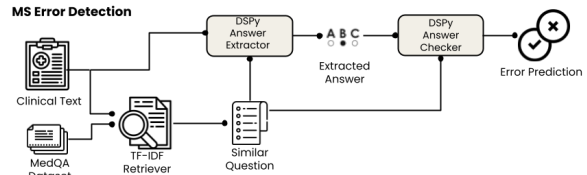


Figure 1: Predicting the presence of an error through a comparison to the retrieved question
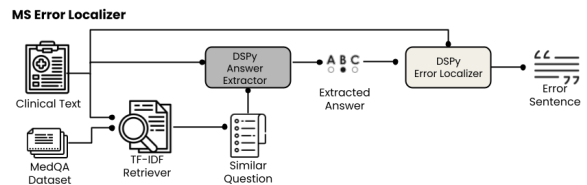


Figure 2: Identifying the error sentence

question, as shown in Figure [1](#). We simultaneously bootstrap these two steps, optimizing the entire pipeline based on the accuracy of the overall error flag prediction.

The result of this bootstrapping process is a compiled program with optimized multi-step chain of thought prompts based on the module's performance on error detection accuracy. This approach allows us to effectively identify the presence of errors in the query text by leveraging the knowledge from external medical question-answering datasets.

### 4.2.3 Localizing Errors within Query Text

After detecting an error in the query text, we use a DSPy module to identify the specific line containing the error, as illustrated in Figure [2](#). This module takes the extracted answer choice and the preprocessed query text as inputs and then an LLM call is done to determine which line most closely matches the erroneous answer choice.

Our experiments showed that GPT-4's performance was high enough that we did not need to compile the program or bootstrap few-shot prompts via a DSPy teleprompter.

The module outputs the line number where the error is located, which is crucial for the subsequent error correction step, as it allows for targeted correction of the relevant text.

### 4.2.4 Error Correction with DSPy

After identifying the error location within the query text, we use a final DSPy module to generate a corrected version of the text, as illustrated in Figure [3](#). This module takes three inputs: the error line, the extracted answer choice, and the correct answer
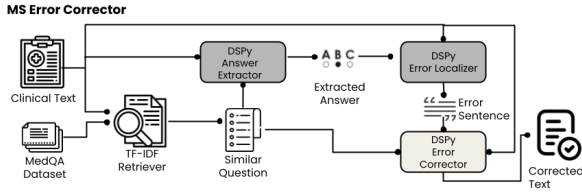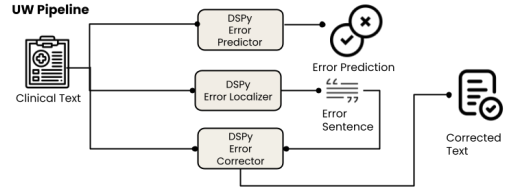
Figure 3: Generating the corrected sentence



Figure 4: Overview of the UW dataset pipeline, consisting of three main stages: error detection, error localization, and error correction. Each stage is implemented using a DSPy module optimized with the MIPRO teleprompter (Khattab et al., 2023) The pipeline also includes a quality control step based on the ROUGE-L score between the original erroneous text and the corrected version.

derived from the most similar retrieved multiple-choice question.

The error correction module utilizes a chain of thought prompt along with 20 few-shot examples generated by the BootstrapFewShotWithRandom-Search teleprompter. This teleprompter samples examples from the training set and generates intermediate labels, such as rationales for the chain of thought, to provide additional context and guidance for the language model during the error correction process. The teleprompter optimizes the selection of few-shot prompts based on their performance on the validation set, using the ROUGE-L score as the metric.

The selected few-shot examples, accompanied by the generated intermediate labels, demonstrate how to modify the error line based on the extracted answer choice and the correct answer, serving as a reference for the model to learn from and adapt to the specific error correction task.

The module outputs the corrected version of the query text, with the error line revised based on the correct answer derived from the most similar multiple-choice question. This corrected text represents the final output of our retrieval-based approach for the MS dataset, addressing the subtle errors present in the clinical text.

### 4.3 Approach for UW Dataset

Our approach for the UW dataset involves optimizing a series of DSPy modules to accomplish all three subtasks sequentially, as illustrated in Figure 4. In all of our experiments, we utilized GPT-4-0125-preview as the underlying large language model, using default generation parameters (temperature of 1.0, top_p of 1) with the exception of a max tokens value of 4096.

#### 4.3.1 Error Detection with DSPy

For the UW dataset, we first employ a DSPy program to identify whether an error exists in the given clinical text snippet. This program is optimized using the Multi-prompt Instruction Proposal Optimizer (MIPRO) teleprompter, which generates

and optimizes both the base prompts and few-shot examples. MIPRO optimizes the prompts and few-shot examples to maximize performance on the validation set, which we created by dividing the UW training collection (160 examples) into 80 training examples, 40 validation examples, and 40 test examples. The optimizer uses error flag accuracy as the metric to optimize and generates 20 examples. We also incorporate chain of thought reasoning into the DSPy module.

#### 4.3.2 Error Localization

If an error is detected in the clinical text snippet, we use another DSPy module to identify the specific line containing the error. This module is also optimized using MIPRO, which generates 20 bootstrap examples that include chain of thought rationales. Using a separate DSPy module for error localization allows us to precisely identify the source of the error and facilitate targeted corrections. The exact match of the error line is used as the metric for optimization, and this module is trained only on a subset of the training samples that contain errors.

#### 4.3.3 Error Correction

After identifying the error line, we use a third DSPy module to generate a corrected version of the erroneous text. This module is also optimized using MIPRO, following the same process as the previous modules. The error correction module takes the erroneous text as input and generates a corrected version based on the optimized prompts and weights. MIPRO uses the ROUGE-L score against the known correct sentence as the metric to optimize, and this module is trained only on a subset of the training samples that contain errors.

| Rank | Team | Error Flags Accuracy |
|------|------|---------------------|
| 1 | WangLab | 86.5% |
| 2 | MediFact | 73.7% |
| 3 | knowlab_AIMed | 69.4% |
| 4 | EM_Mixers | 68.0% |
| 5 | IKIM | 67.8% |
| 6 | IryoNLP | 67.1% |
| 7 | Edinburgh Clinical NLP | 66.9% |
| 8 | hyeonhwang | 63.5% |
| 9 | PromptMind | 62.2% |
| 10 | CLD-MEC | 56.6% |

Table 1: Top 10 teams' performance on Task 1 (Error Flags Accuracy)

| Rank | Team | Error Sentence Detection Accuracy |
|------|------|----------------------------------|
| 1 | WangLab | 83.6% |
| 2 | EM_Mixers | 64.0% |
| 3 | knowlab_AIMed | 61.9% |
| 4 | hyeonhwang | 61.5% |
| 5 | Edinburgh Clinical NLP | 61.1% |
| 6 | IryoNLP | 61.0% |
| 7 | PromptMind | 60.9% |
| 8 | MediFact | 60.0% |
| 9 | IKIM | 59.0% |
| 10 | HSE NLP | 52.0% |

Table 2: Top 10 teams' performance on Task 2 (Error Sentence Detection Accuracy)

### 4.3.4 Quality Control with ROUGE-L

To ensure the quality of the generated corrections, we calculate the ROUGE-L score between the original erroneous text and the corrected version. If the ROUGE-L score is below a threshold of 0.7, which we set as an arbitrary estimate for quality, we reject the correction and use the original erroneous text instead. This fallback mechanism is based on the observation that the ROUGE-L score of the erroneous text tends to be quite high since the error is only a small portion of the sentence. However, this fallback is more of a contest-metric-focused feature rather than something that significantly improves performance.

## 5 Results and Discussion

### 5.1 Overall Performance in the MEDIQA-CORR 2024 Shared Task

Our approach achieved top performance in the MEDIQA-CORR 2024 shared task across all three subtasks. Tables 1, 2, and 3 present the performance of the top 10 teams in each subtask.

### 5.2 Performance on Subtask 1 - Error Prediction

In the official contest results for binary error prediction, our approach achieved an accuracy of 86.5%, ranking first among all participating teams. Table 1 shows the top 10 teams' performance on Task 1.

### 5.3 Performance on Subtask 2 - Error Sentence Detection

For error sentence detection, we obtained an accuracy of 83.6%, ranking first among all teams. Table 2 presents the top 10 teams' performance.

These results demonstrate the effectiveness of our few-shot learning and CoT-based approach in detecting the presence of errors and localizing the specific sentences containing the errors.

### 5.4 Performance on Subtask 3 - Sentence Correction

For subtask C (Sentence Correction), the official contest results show that our approach achieved an Aggregate-Score of 0.789, which is the mean of ROUGE-1-F (0.776), BERTScore (0.809), and BLEURT (0.783). This was the highest score among the participating teams for the sentence correction task. Table 3 displays the top 10 teams' performance on Task 3.

The official contest results highlight the competitive performance of our approach across all three subtasks of the MEDIQA-CORR 2024 shared task, demonstrating its effectiveness in detecting, localizing, and correcting medical errors in clinical text for both the MS and UW datasets.

### 5.5 Implications and Limitations of the Approach

Our work contributes to the ongoing efforts in improving the accuracy and reliability of medical information in clinical text. The automated detection and correction of certain types of errors could ensure the quality and consistency of medical documentation, ultimately supporting patient safety and quality of care. The development and integration of more advanced systems could help alleviate the burden of manual error checking for the specific error types addressed, allowing healthcare providers to allocate more time and resources to delivering high-quality patient care.

However, it is important to acknowledge the limitations of our approach in the context of the diverse nature of errors in medical documentation. While our system demonstrates strong performance on the MS and UW datasets, it focuses on a specific subset of errors and has not been shown to be effec-

| Rank | Team | AggregateScore | R1F | BERTSCORE | BLEURT | AggregateCR |
|------|------|----------------|-----|-----------|--------|-------------|
| 1 | WangLab | 0.789 | 0.776 | 0.809 | 0.783 | 0.775 |
| 2 | PromptMind | 0.787 | 0.807 | 0.806 | 0.747 | 0.574 |
| 3 | HSE NLP | 0.781 | 0.779 | 0.806 | 0.756 | 0.512 |
| 4 | hyeonhwang | 0.734 | 0.729 | 0.767 | 0.705 | 0.571 |
| 5 | Maven | 0.733 | 0.703 | 0.744 | 0.752 | 0.524 |
| 6 | Edinburgh Clinical NLP | 0.711 | 0.678 | 0.744 | 0.711 | 0.563 |
| 7 | knowlab_AIMed | 0.658 | 0.643 | 0.677 | 0.654 | 0.573 |
| 8 | EM_Mixers | 0.587 | 0.571 | 0.595 | 0.596 | 0.548 |
| 9 | IryoNLP | 0.581 | 0.561 | 0.592 | 0.591 | 0.528 |
| 10 | IKIM | 0.559 | 0.523 | 0.564 | 0.588 | 0.550 |

Table 3: Top 10 teams' performance on Task 3 (Aggregate Score and its components)

tive in addressing the wide diversity of errors that can occur in medical documentation.

For instance, our approach does not currently address errors that are propagated through multiple notes when a physician references prior documents containing inaccuracies, such as incorrect medical history. Such errors can be particularly challenging to identify and correct, as they may require a comprehensive understanding of the patient's medical history, the context of the referenced documents, and the resolution of conflicting statements across documents. Our system has not been designed or evaluated for handling these types of errors.

Moreover, our approach does not cover errors that originate from sources beyond the scope of our training data, such as poor transcriptions, entries in the wrong medical record, or errors in decision making. These types of errors may necessitate different strategies and techniques for detection and correction, and our current approach has not been developed to handle them.

Additionally, the reliance on external datasets for the retrieval-based approach in the MS dataset limits the generalizability of our method to other medical domains or datasets. In fact, we believe that an approach used in the MS dataset might actually create further errors if used on real clinical text, as real clinical practice does not always reflect optimal or most likely completions. The effectiveness of our approach in detecting and correcting errors may vary depending on the specific characteristics and error types present in different medical contexts, and further evaluation would be necessary to assess its performance in diverse settings.

### 5.5.1 Impact of Different LLMs and Compilation

After the competition ended, we performed additional experiments to compare the performance of our approach when using GPT-4 and GPT-3.5 as the underlying language models for the DSPy modules, as well as the impact of using compiled and uncompiled DSPy programs.

Table 4 presents the results of the ablation study for error flag accuracy (Task 1), error sentence detection accuracy (Task 2), and various metrics for Task 3. The results show that using GPT-4 as the underlying LLM consistently yields better performance compared to GPT-3.5 across all tasks. For Task 1, the compiled GPT-4 model achieves the highest accuracy of 97.3% (0.1%), while for Task 2, it achieves an accuracy of 97.0% (0.1%). The compiled DSPy programs outperform their uncompiled counterparts for both GPT-3.5 and GPT-4.

In Task 3, the compiled GPT-4 model consistently outperforms the other models across all metrics, with the highest AggregateC score of 0.878 (0.002). Moreover, the results demonstrate that using compiled DSPy programs consistently outperforms the uncompiled approach across all tasks and datasets, emphasizing the significance of systematic optimization techniques in enhancing the performance of our error detection and correction system.

It is important to note that we did not isolate the impact of retrieval in our post-competition experiments, as it was a fundamental component of all the modules in our approach. Removing the retrieval component would require the development of a new solution. However, the strong performance of our uncompiled GPT-3.5 solution suggests that a significant portion of the performance could be attributed to the retrieval process itself. Future work should

| Error Flags Accuracy (Task 1) | | | | |
|---|---|---|---|---|
| | GPT-3.5 Compiled | GPT-3.5 Uncompiled | GPT-4 Compiled | GPT-4 Uncompiled |
| Error Flags Accuracy | 94.0% (0.4%) | 81.2% (0.7%) | 97.3% (0.1%) | 88.9% (0.5%) |
| Error Sentence Detection Accuracy (Task 2) | | | | |
| | GPT-3.5 Compiled | GPT-3.5 Uncompiled | GPT-4 Compiled | GPT-4 Uncompiled |
| Error Sentence Detection Accuracy | 92.8% (0.5%) | 78.5% (0.8%) | 97.0% (0.1%) | 88.0% (0.8%) |
| Task 3 Metrics | | | | |
| Metric | GPT-3.5 Compiled | GPT-3.5 Uncompiled | GPT-4 Compiled | GPT-4 Uncompiled |
| aggregate_subset_check | 0.853 (0.001) | 0.809 (0.011) | 0.824 (0.003) | 0.827 (0.003) |
| R1F_subset_check | 0.827 (0.003) | 0.778 (0.017) | 0.789 (0.003) | 0.792 (0.003) |
| BERTSCORE_subset_check | 0.874 (0.001) | 0.827 (0.013) | 0.856 (0.003) | 0.857 (0.002) |
| BLEURT_subset_check | 0.859 (0.000) | 0.824 (0.006) | 0.827 (0.002) | 0.832 (0.003) |
| AggregateC | 0.864 (0.004) | 0.736 (0.010) | 0.878 (0.002) | 0.792 (0.005) |

Table 4: Ablation studies for error flag accuracy (Task 1), error sentence detection accuracy (Task 2), and Task 3 metrics. Numbers in parentheses represent standard deviations.

explore the impact of different retrieval strategies on the performance of error detection and correction in clinical text.

### 5.6 Future Research Directions

Although our approach has demonstrated competitive performance in the MEDIQA-CORR 2024 shared task, there are several potential avenues for future research that could further improve the effectiveness and applicability of our system.

One area for future investigation is the fine-tuning of open access models specifically for clinical notes (Toma et al., 2023). While fine-tuning may lead to higher performance, we focused on working with DSPy in the current study and did not have the computational resources to maintain the necessary throughput and latency during initial experimentation. Future studies could examine the trade-offs between fine-tuning and using off-the-shelf models with prompt optimization techniques, taking into account factors such as performance, efficiency, and scalability.

Another direction for future research is the expansion of the benchmark dataset to include a broader range of errors, such as those spanning multiple documents or involving suboptimal clinical decisions. Broadening the scope of the dataset would enhance the robustness of error detection and correction systems and extend their applicability to more complex clinical scenarios.

Integrating domain-specific knowledge, such as medical ontologies or expert-curated rules, into our approach could improve the system's ability to handle complex medical cases and make more informed decisions. This would be particularly relevant if the errors include suboptimal clinical decisions, as the system could provide more comprehensive support to healthcare professionals.

Lastly, developing more comprehensive and robust methods for measuring and correcting errors is an area with significant potential. This could involve creating standardized evaluation metrics and datasets that better capture the intricacies of medical errors and developing more advanced error correction techniques that can handle a wider range of error types and contexts.

## 6 Conclusion

The approach presented in this paper, which combines retrieval-based methods, few-shot learning, and systematic prompt optimization, demonstrates the potential of AI-assisted tools for detecting and correcting medical errors in clinical text. The strong performance achieved across all three subtasks of the MEDIQA-CORR 2024 shared task highlights the effectiveness of our methods in addressing the specific challenges posed by different datasets and error types. However, further research is necessary to extend the applicability of our approach to a wider range of medical contexts, incorporate domain-specific knowledge, and integrate with existing clinical systems. As the field of AI-assisted medical error detection and correction continues to evolve, collaboration between AI researchers and healthcare professionals will be crucial to develop solutions that effectively augment and support clinical decision-making processes, ultimately contributing to improved patient safety and healthcare quality.

# References

Asma Ben Abacha, Wen wai Yim, Yujuan Fu, Zhaoyi Sun, Fei Xia, and Meliha Yetisgen. 2024a. Overview of the mediqa-corr 2024 shared task on medical error detection and correction. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Asma Ben Abacha, Wen wai Yim, Yujuan Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. 2024b. Medec: A benchmark for medical error detection and correction in clinical notes. *CoRR*.

Asma Ben Abacha, Wen-wai Yim, Griffin Adams, Neal Snider, and Meliha Yetisgen. 2023. Overview of the MEDIQA-chat 2023 shared tasks on the summarization & generation of doctor-patient conversations. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 503–513, Toronto, Canada. Association for Computational Linguistics.

John Giorgi, Augustin Toma, Ronald Xie, Sondra Chen, Kevin An, Grace Zheng, and Bo Wang. 2023. WangLab at MEDIQA-chat 2023: Clinical note generation from doctor-patient conversations using large language models. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 323–334, Toronto, Canada. Association for Computational Linguistics.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Preprint*, arXiv:2009.13081.

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *Preprint*, arXiv:2310.03714.

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.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. *Preprint*, arXiv:2303.13375.

Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. Bleurt: Learning robust metrics for text generation. *Preprint*, arXiv:2004.04696.

Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.

Augustin Toma, Patrick R. Lawler, Jimmy Ba, Rahul G. Krishnan, Barry B. Rubin, and Bo Wang. 2023. Clinical camel: An open expert-level medical language model with dialogue-based knowledge encoding. *Preprint*, arXiv:2305.12031.

Dave Van Veen, Cara Van Uden, Louis Blankemeier, 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*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. *Preprint*, arXiv:1904.09675.

# WangLab at MEDIQA-M3G 2024: Multimodal Medical Answer Generation using Large Language Models

**Ronald Xie**[1,3]    **Steven Palayew**[1,3]    **Augustin Toma**[1,3]
**Gary D. Bader**[1,4,5]    **Bo Wang**[1,2,3,6]

[1]University of Toronto    [2]Peter Munk Cardiac Centre, University Health Network
[3]Vector Institute    [4]Princess Margaret Cancer Centre, University Health Network
[5]Lunenfeld-Tanenbaum Research Institute, Sinai Health System
[6]AI Hub, University Health Network

{augustin.toma, ronald.xie, steven.palayew}@mail.utoronto.ca, gary.bader@utoronto.ca, bowang@vectorinstitute.ai

## Abstract

This paper outlines our submission to the MEDIQA2024 Multilingual and Multimodal Medical Answer Generation (M3G) shared task. We report results for two solutions under the English category of the task, the first involving two consecutive API calls to the Claude 3 Opus API and the second involving training an image-disease label joint embedding in the style of CLIP for image classification. These two solutions scored 1st and 2nd place respectively on the competition leaderboard, substantially outperforming the next best solution. Additionally, we discuss insights gained from post-competition experiments. While the performance of these two solutions have significant room for improvement due to the difficulty of the shared task and the challenging nature of medical visual question answering in general, we identify the multi-stage LLM approach and the CLIP image classification approach as promising avenues for further investigation.

## 1 Introduction

An increased demand for healthcare services and recent pandemic needs have accelerated the adoption of telehealth, which was previously underused and understudied (Shaver, 2022; wai Yim et al., 2024a). There has been significant recent interest in integrating artificial intelligence (AI) into telehealth Ma et al., 2024; Toma et al., 2023, as these technologies have the potential to enhance and expand its ability to address important healthcare needs (Sharma et al., 2023). The task of consumer health question answering, an important part of telehealth, has been explored actively in research. However, the focus of this existing research has been on text (Ben Abacha et al., 2019), which is limiting as medicine is inherently multimodal in nature, requiring clinicians to work not just with text but also with imaging among other modalities (Corrado and Matias, 2023).

To help address this gap, the MEDIQA-M3G shared task was proposed (wai Yim et al., 2024a). This task requires the automatic generation of clinical responses given relevant user generated text and images as input, with a specific focus on clinical dermatology (wai Yim et al., 2024a).

This work describes our submission to this task. We explored two standalone solutions, one involving two consecutive API calls to the recently released Claude 3 Opus model (Anthropic) and the other trains a joint image-disease label embedding model using CLIP (Radford et al., 2021) for image classification. These two strategies took 1st and 2nd place respectively during the competition. While our strategy's effectiveness relative to other submissions highlight that Claude 3 Opus and multi-stage LLM frameworks have potential value in the area of multi-modal medical AI, both our solutions' performance is limited despite their leaderboard success, highlighting the difficulty of the shared task and the unsolved challenge of medical visual question answering.

## 2 Shared task and provided dataset

The MEDIQA-M3G competition focuses on the problem of clinical dermatology multimodal query response generation. The inputs include text which give clinical context and queries, as well as one or more images associated with the case (wai Yim et al., 2024b). The task is to generate responses to these cases resembling those made by medical professionals in the field of dermatology. Participants have the option to generate these responses in three languages: Chinese (Simplified), English, and Spanish. (wai Yim et al., 2024a)

The dataset consists of 842 train, 56 validation, and 100 test cases. Each case consists of one or more images of skin conditions, their accompanying query text which may or may not include clinical context, patient queries, additional details

regarding the disease and in some cases possible diagnosis. Finally, for each case there are multiple responses made by one or more medical professionals, which are used as targets to score the model predictions. The cases also notably include metadata on the rank and validation level of the authors of content, which are used in evaluation (wai Yim et al., 2024b). For evaluation, the competition uses a version of the deltaBLEU (Galley et al., 2015) metric to allow a single score to be computed based on word matching, weighted by the consistency (most frequent response) and the seniority of the medical professional across all responses given for that particular case. (wai Yim et al., 2024a)

The query text and target responses are given in multiple languages, namely Chinese, English, and Spanish (wai Yim et al., 2024b). It's worth noting that while the test and validation sets were translated by medical professionals, the training set of 842 cases seems to be translated automatically with some potential room for errors. For our submission we focus on only providing the English solution.

## 3 Related Work

There has recently been a substantial amount of interest in medical applications of multimodal machine learning, and large multimodal models. Some notable examples of research in this area include the open source LLAVA-MED model (Li et al., 2023), and ELIXR, with the latter, similar to our work, exploring not only the application of large multimodal models, but also training a model using CLIP (Xu et al., 2023). However, while there has been significant focus on certain areas such as radiology, the area of dermatology has not been explored to the same extent. Cirone et al. notably found that GPT-4V could accurately differentiate between benign lesions and melanoma (Cirone et al., 2024). However, this is a much less challenging task than the one proposed in this shared task, as the problem space is much smaller in scope than responding to dermatology questions which are not necessarily in the train set, with even the conditions of interest not necessarily being in the train set. The limited performance of our solution, along with it being by far the best performing solution in this competition demonstrate the challenge of this task, and highlight the need for significant progress before deployment in a clinical setting. However, our work highlights potentially important directions for future research, including further investigation

| Rank | Team | dBLEU (English) |
|------|------|-----------------|
| 1 | WangLab | 12.855 |
| 2 | kiyoonyoo | 3.827 |
| 3 | amdada | 2.662 |
| 4 | romarcg | 2.133 |
| 5 | xiaolihaixiao | 1.758 |
| 6 | pvashisht | 0.923 |
| 7 | nadia | 0.717 |
| 8 | abrygo | 0.457 |

Table 1: Top 8 teams' performance on the English catgeory for the MEDIQA-M3G competition

of multi-stage LLM systems, and the importance of evaluation metrics in the benchmarking of the clinical efficacy of developed systems.

## 4 Results

Upon examination of the evaluation metric and competition data, we have determined that a short response focusing on disease diagnosis alone is the most advantageous. This is due to two reasons. First, we notice both the training and validation sets often contain short responses, and in many cases merely the skin condition presented in the associated images. Second, the evaluation metric's penalty for short responses is significantly smaller than a longer, partially correct response. Given these initial findings, we evaluated two methods as outlined in 1 which took 1st and 2nd place in the English category of the leaderboard during the MEDIQA-M3G challenge by a significant margin over the next best submitted solution, the latter of which received a deltaBLEU score of 3.827 during the competition. The methods will be elaborated in the following sections in detail.
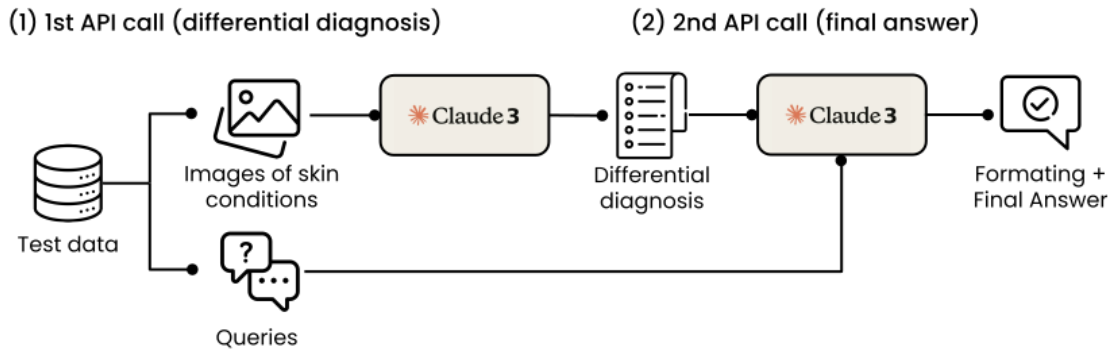
### 4.1 Claude 3 Opus API solution

The higher scoring of the two methods during the competition consists of two successive API calls to Claude 3 Opus (Anthropic). For each case in the test set, the first API call generates possible differential diagnosis for the given images, and the second API call further processes the response into the name of the most likely disease only, which is then returned.

This exact configuration was decided based on trial and error. Table 2 outlines the solutions tested. Notably, we observe that the disease diagnosis given by Claude 3 Opus was poorer quality when

## (A) Claude 3 Opus VQA

**(1) 1st API call (differential diagnosis)**   **(2) 2nd API call (final answer)**
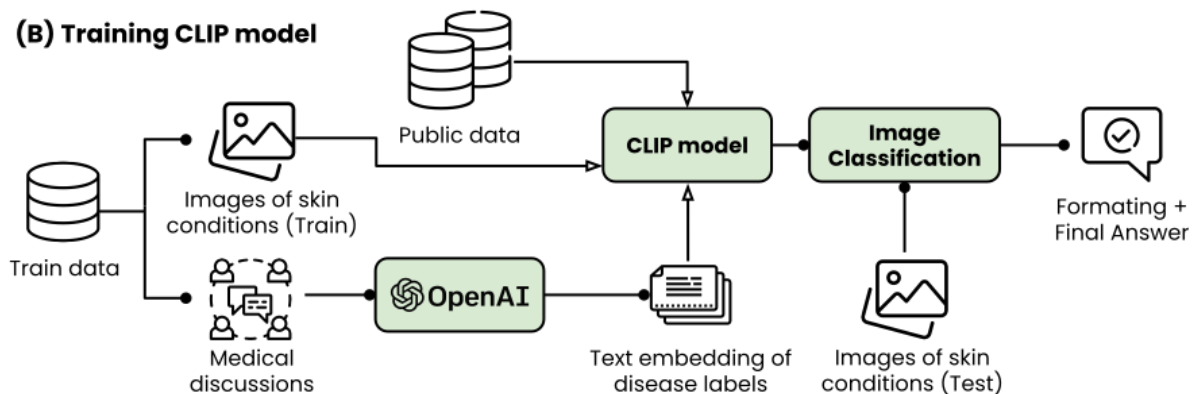
## (B) Training CLIP model

Figure 1: Overview of the two winning solutions. A) Test cases are directly submitted to the Claude 3 Opus API. The first of the two consecutive API calls generates differential diagnosis using only the images in the test cases and the second API call optionally includes the associated queries, specifies formatting, and generates final answer. The top performing Claude 3 Opus solution did not utilize test queries. B) The medical discussions included as a part of the training data is used to extract the most likely disease label for each case using GPT4-Turbo from OpenAI. The resulting image-disease label pair are used in conjunction with publicly available data to train a joint embedding in the style of CLIP. The disease labels are embedded using EmbeddingV3 from OpenAI and used to train the image encoder (ResNet50) and both the image and text projection layers. Finally, once the model is trained, the test images are classified inside the learned joint embedding which becomes the final output before performing post processing.

the prompt constrains the output format upon manual review. This was further confirmed by the inferior performance of the 1-call result. Therefore, we let the API generate differential responses with the provided images alone without any constraints on the format of the first response, and use a second API call to reformat that response into the desired form, which is just the name of the skin condition without any abbreviations.

Furthermore, we observe that including the accompanying query text for each case either in the 1st or 2nd pass was not able to outperform simply using the image alone to make the predictions. This finding may be attributed to the inconsistent information present in the query text, which may often harm the prediction from Claude 3 in some cases. It may also be a potential limitation of Claude's ability to reason with text and image simultane-

ously. Indeed the resulting predictions had substantial room for improvement even under the most favorable setting tested. All prompts used to produce the solutions in Table 2, including the winning solution are outlined in Appendix.

### 4.2 CLIP image classification solution

The second solution we've explored took second place during the MEDIQA-M3G challenge, and with subsequent tuning after the competition, was able to overtake the Claude 3 Opus API solution under the same evaluation setting used during the competition. The CLIP based solution involved learning a joint representation between the images of the skin conditions and their accompanying disease label. We achieved this by using a contrastive learning setup inspired by CLIP (Radford et al., 2021). We use a ResNet-50 (He et al., 2015) en-

Table 2: Performance of various Claude 3 Opus based solutions. 1 Call involves simply generating a response based on images, whereas 2 Calls involve first generating a differential diagnosis, then using a second API call to come up with a final diagnosis. For Img+text, both modalities are used in the first API call to generate the differential, whereas for Img then text the first API call uses only images, then the second API call uses text

| Scen. | dBLEU | BP | Ratio | Hyp_len | Ref_len |
|---|---|---|---|---|---|
| Img (1 Call) | 10.529 | 0.984 | 0.984 | 498 | 506 |
| Img (2 Calls) | 12.855 | 0.994 | 0.994 | 485 | 488 |
| Img then text (2 Calls) | 10.905 | 0.983 | 0.983 | 527 | 536 |
| Img + text (2 Calls) | 10.905 | 1.000 | 1.004 | 523 | 521 |

coder initialized with pretrained weights as the image encoder. Image augmentations include random flip, random rotations, random spatial cropping and random contrast adjustments to improve diversity and robustness of training. To obtain the disease label from the provided responses of medical professionals, we input all medical professional responses for each case in a GPT4-Turbo API call and prompt the GPT4-turbo model to return the most consistent disease diagnosis among all responses. We also curate additional image-disease pairs (n = 25528) in the domain of dermatology from publicly available sources. It's worth highlighting that there were 1245 unique disease labels among the image-disease label pairs curated. The label sparsity effectively makes training a traditional supervised classification model difficult. However, we make the observation that these labels were often the result of label inconsistency and frequently shared semantic meaning, which motivated our use of OpenAI's EmbeddingV3 (OpenAI, 2024) model to produce consistent, semantically meaningful word embeddings which effectively serve as the text encoder in our CLIP learning framework. We visualize the embeddings of the disease labels and verify that indeed many diseases with similar descriptions cluster together, as evident in Figure AS2. Specific hyperparameters used to produce the highest scoring CLIP solution are outlined in Table 6.

### 4.2.1 Image classification via nearest neighbour retrieval

Once the image encoder and the respective image and text projection layers are trained, the resulting joint embedding can be used to perform image classification via nearest neighbour retrieval. Specifically, we embed each image associated with a given case in the competition test set and find 5 nearest neighbours for each embedded image. We test 4 different conditions, namely retrieval between the

image embedding of the query (testing dataset) and either its nearest 5 text or image embeddings from the reference (training dataset), and whether the nearest neighbours are computed in PCA space (10 components) or as normal. We then pool the labels associated with the retrieved examples via majority voting and return the final predicted label for the case. The resulting scores are presented in Table 4. Of note, during the competition (1st row), random augmentations were mistakenly not turned off during inference when obtaining the image embeddings. This did not lead to better performance and was corrected after the competition concluded.

### 4.2.2 Importance of batch size

The CLIP loss heavily relies on a diverse source of positive and negative pairs to converge to a good solution. It's often the case that bigger batch sizes give more robust joint representations. However, under low data settings such as for this competition where the available labelled data is scarce, larger batch sizes may lead to overfitting which is destructive for generalization. We test 3 different batch sizes ranging from 128 to 512 and observe that a batch size of 256 is most suitable under prior evaluation scripts shared by the competition organizers. However, when using the updated evaluation script during test phase of the challenge, we observe that both batch size 256 and 512 exhibit comparable performance. The results under the updated evaluation script are presented in Table 3.

Table 3: Performance of the CLIP based solution across different batch sizes

| Model | dBLEU | BP | Ratio | Hyp_len | Ref_len |
|---|---|---|---|---|---|
| CLIP (batch 128) | 10.434 | 0.980 | 0.980 | 483 | 493 |
| CLIP (batch 256) | 12.080 | 0.966 | 0.966 | 461 | 477 |
| CLIP (batch 512) | 12.289 | 0.983 | 0.984 | 447 | 485 |

Table 4: Performance of CLIP with different retrieval related strategies, including retrieval in the PCA space (n=10), and retrieving based on either the image or the text embedding of the reference. The first row indicate the CLIP based solution submitted during the competition. Of note, the random image augmentations during inference were enabled unintentionally during the competition but disabled for all subsequent experiments.

| Random. Aug | PCA Space | Query-Reference | dBLEU |
|:---:|:---:|:---:|:---:|
| Yes* | Yes | Image-Image | 11.979 |
| No | No | Image-Text | 12.123 |
| No | Yes | Image-Text | 8.396 |
| No | No | Image-Image | 15.884 |
| No | Yes | Image-Image | 12.079 |

## 4.3 Post processing

Post processing is performed on both the Claude 3 Opus API solution and the CLIP based image classification solution in the same way. This includes putting the output disease name in predetermined sentence format to mimic the style of the given responses from medical professionals, specifically in the form of "It is [Disease name].". While a naive approach to the VQA task, we find this simple formatting allows our disease labels produced from images alone to score quite competitively under the deltaBLEU evaluation metric provided by the competition organizers compared to simply returning the disease name itself as evident in Table 5.

Furthermore, unlike other competitors' solutions based on finetuning existing VQA models (such as LLaVA-med) simultaneously using both the images and the associated query text, our solution does not take advantage of any potentially useful information included in the query text. As a naive way of overcoming this limitation, we compiled a dictionary of disease names present in the training data and do exact word matching with the query text. Cases where the query text matches with the dictionary will have their model predictions replaced with the matched disease condition. These matches constitute 15 cases out of 100 in the testing data. While this naive heuristic often times do not produce the correct diagnosis, considering the difficulty of the task this approach does confer some improvement in overall deltaBLEU score. The ablations of the post processing is outlined in Table 5.

| Solution | Word Matching | Sentence Structure | Both |
|:---|:---:|:---:|:---:|
| Claude Solution | 4.903 | 6.202 | 12.855 |
| CLIP Solution (competition) | 2.386 | 3.253 | 11.979 |
| CLIP Solution (batch 256) | 3.255 | 10.923 | 15.884 |

Table 5: Result of ablations on performance of top performing solutions. Sentence structure involves placing the predicted disease labels in predetermined sentence format, whereas word matching is a heuristic employed to utilize provided text via naively matching disease names with the given queries.

| HyperParameter | Value |
|:---|:---:|
| Image encoder | Resnet50 |
| Projection dim | 256 |
| Batch size | 256 |
| Text embedding dim | 3072 |
| Image embedding dim | 2048 |
| Num. projection layers | 1 |
| Augmentations | RandFlip, RandRotate, RandSpatialCrop, RandAdjustContrast |
| Weight decay | 0.001 |
| Learning rate | 0.001 |

Table 6: Hyperparameters corresponding to the highest performing CLIP solution

## 5 Discussion

We have presented two solutions to the MEDIQA2024-M3G competition, one involving API calls to an existing state of the art multimodal language model and the other involving the learning of an image-disease label joint embedding space for disease classification.

The superior performance of using two separate API calls to Claude 3 Opus over one pass was interesting to observe. The increase in performance is likely attributed to the reduced ability for the model to simultaneously reason with the images while adhering to the added difficulty of only returning the disease label without any additional textual generation. This finding is somewhat consistent with how chain of thought reasoning can improve model performance by asking the model to first consolidate evidence present in the given image followed by making several differential diagnoses. Further research such as (Zhang et al., 2023) also highlight the importance of using two-stage frameworks for multi-modal chain of thought that separate rationale generation and answer inference over one stage systems.

For the CLIP based solution, we find it extremely encouraging that a smaller scale model finetuned

on image-disease label pairs (n=25528) was able to outperform Claude 3 Opus (dBLEU of 15.884 vs 12.855). It perhaps demonstrates that smaller scale supervised training may sometimes outperform bigger more general purpose models for specific tasks of interest due to the advantage of training only on task specific examples. Furthermore, our additional experiments after the competition highlights the importance of proper selection of batch size and retrieval method. We observe that while CLIP effectively constructs a joint embedding space between images and their disease labels, the image embeddings and text embeddings remain as separate cluster in PCA space. As a result, we see that the nearest 5 neighbours in the text cluster for each embedded image (image-text) in the test set were much poorer in quality than those retrieved from the image cluster (image-image).

## 6 Limitations

While both the Claude 3 Opus API based solution and the CLIP based image classification solution achieved first and second place during the MEDIQA-M3G competition respectively, they have substantial room for improvement despite their leaderboard success.

First of all, the overall deltaBLEU score of both solutions are poor, mostly ranging from 10-15 dBLEU. The low absolute scores of the solutions really highlight the difficulty of the medical VQA task presented and the difficulty of such tasks in general. Upon examining the solutions, we observe that the models were seldom able to generate the exact name of the skin condition in question, although do a good job at identifying a disease similar in presentation or effect location (for example tinea scalp vs seborrheic dermatitis). Certainly both solutions require substantial improvements before they contribute meaningful benefits to the healthcare system in practice.

While the CLIP based solution was able to outperform our Claude 3 Opus API based solution with experiments conducted post competition, it is worth mentioning that such small scale finetuning may be less desirable as the model would have to be repurposed for new problems of interested each time. LLM based solutions have the advantage of being general purpose and do not have this issue. Furthermore, due to the tight schedules of the competition, both solutions were not explored to their full potential. We anticipate there are bigger up-

sides for the Claude 3 Opus API solution via more sophisticated prompting or compiling. Our rather simple implementation of the Claude based API solution may not represent the LLM's full capability but rather offers a competitive baseline for this task.

Next, both solutions while reproducible are not stable. The Claude API may be subject to randomness during generation due to the temperature parameter or the update of internal private model weights while the CLIP solutions observed inconsistencies during retrieval where the retrieved images' labels seldom agreed with each other despite relatively similar appearances, leading to low confidence in the final output. Retraining the CLIP model with the same experimental setup but initializing differently may yield completely different final disease label classification due to this inconsistency.

Lastly, the two solutions were formulated with the competition evaluation metric in mind as they are both framed as a disease label prediction task rather than a more usual VQA task which could cover a broader range of topics in their generated responses such as differential diagnoses, treatments and other recommendations as present in the actual ground truths for this competition. This is further reason to treat the performance of the presented solutions with a grain of salt. Specifically, upon our initial exploration, the deltaBLEU metric defined by the competition organizers favors short responses given the relatively heavy penalty incurred on incorrect k-mers present and relatively low penalty on a incomplete answer in comparison. This discourages model exploration during text generation and potentially penalizes model predictions that are correct semantically but are either too long or not containing the exact words present in the ground truths. This is highlighted in the ablation results in Table 5. Furthermore, the naive word matching often gave incorrect diagnosis as the patient writing the query does not have medical background, however the solution containing the disease label still scored well under the current metric as medical professionals respond with "not [disease label]" which has the opposite semantic meaning but similar k-mer composition. We recommend the organizers to slightly modify the existing metric to be more lenient with assessing the produced solutions and perhaps add a semantics component in addition to a k-mer based evaluation metric such as GPTscore (Fu et al., 2023), that can

provide more robustness in assessing the quality of generated responses.

Nevertheless, the competition serve as an important step towards the goal of automatically generating clinical responses given textual queries and associated images, and we sincerely thank the organizers for the work curating this dataset and organizing the competition.

## 7 Conclusion

We present two solutions to the English category of the MEDIQA2024-M3G shared task for Multilingual and Multimodal Medical Answer Generation. The Claude 3 Opus API based solution and the CLIP image classification based solution scored 1st and 2nd, respectively among all submissions. While there is still substantial room for improvement for these two solutions, we share and discuss our findings to contribute towards the important goal of automatically generating clinical responses given textual queries and associated images.

## 8 Acknowledgement

We extend our sincere thanks to the Digital Research Alliance of Canada for their support and computing resources. We also would like to express gratitude to both internal and external reviewers for their insightful feedback, which enhanced earlier versions of this paper. Finally, we would like to thank the organizers for all the work put into hosting this interesting and challenging competition.

# References

Anthropic. The claude 3 model family: Opus, sonnet, haiku.

Asma Ben Abacha, Chaitanya Shivade, and Dina Demner-Fushman. 2019. Overview of the MEDIQA 2019 shared task on textual inference, question entailment and question answering. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 370–379, Florence, Italy. Association for Computational Linguistics.

Katrina Cirone, Mohamed Akrout, Latif Abid, and Amanda Oakley. 2024. Assessing the utility of multimodal large language models (GPT-4 vision and large language and vision assistant) in identifying melanoma across different skin tones. *JMIR Dermatol*, 7:e55508.

Greg Corrado and Yossi Matias. 2023. Multimodal medical AI. https://blog.research.google/2023/08/multimodal-medical-ai.html. Accessed: 2023-12-4.

Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. GPTScore: Evaluate as you desire.

Michel Galley, Chris Brockett, Alessandro Sordoni, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 445–450, Beijing, China. Association for Computational Linguistics.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep residual learning for image recognition.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. LLaVA-Med: Training a large Language-and-Vision assistant for biomedicine in one day.

Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. 2024. Segment anything in medical images. *Nature Communications*, 15(1):654.

OpenAI. 2024. New embedding models and API updates. https://openai.com/blog/new-embedding-models-and-api-updates. Accessed: 2024-4-11.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision.

Sachin Sharma, Raj Rawal, and Dharmesh Shah. 2023. Addressing the challenges of AI-based telemedicine: Best practices and lessons learned. *J. Educ. Health Promot.*, 12:338.

Julia Shaver. 2022. The state of telehealth before and after the COVID-19 pandemic. *Prim. Care*, 49(4):517–530.

Augustin Toma, Patrick R Lawler, Jimmy Ba, Rahul G Krishnan, Barry B Rubin, and Bo Wang. 2023. Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding. *arXiv preprint arXiv:2305.12031*.

Wen wai Yim, Asma Ben Abacha, Velvin Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen, and Martin Krallinger. 2024a. Overview of the mediqa-m3g 2024 shared task on multilingual and multimodal medical answer generation. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wen wai Yim, Velvin Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Shawn Xu, Lin Yang, Christopher Kelly, Marcin Sieniek, Timo Kohlberger, Martin Ma, Wei-Hung Weng, Atilla Kiraly, Sahar Kazemzadeh, Zakkai Melamed, Jungyeon Park, Patricia Strachan, Yun Liu, Chuck Lau, Preeti Singh, Christina Chen, Mozziyar Etemadi, Sreenivasa Raju Kalidindi, Yossi Matias, Katherine Chou, Greg S Corrado, Shravya Shetty, Daniel Tse, Shruthi Prabhakara, Daniel Golden, Rory Pilgrim, Krish Eswaran, and Andrew Sellergren. 2023. ELIXR: Towards a general purpose x-ray artificial intelligence system through alignment of large language models and radiology vision encoders.

Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023. Multimodal Chain-of-Thought reasoning in language models.

## A  Claude 3 Opus API prompts

Example prompts used to perform API calling in the Claude 3 Opus solution and other tested variants.

### A.1  Image only 1-call

**System:** You are an expert assistant to a blind dermatology student, help him identify exactly what conditions would be included in the differential for this condition? Be concise. After brief description of the images and explanation of your choice, give the most commonly occuring skin disease out of the differentials at the end and nothing else, in the form of

Answer: [Disease Name]

**Content:** IMG_ENC00908_00001.jpg, IMG_ENC00908_00002.jpg
**Output:** Answer: Dyshidrotic eczema

### A.2  Image only 2-calls

**System:** You are an expert assistant to a dermatology student, help him identify exactly what skin conditions would be included in the differential for the images presented. Consider both resemblence and prevalence.
**Content:** IMG_ENC00908_00001.jpg, IMG_ENC00908_00002.jpg
**Output1:** Based on the images provided, the key skin findings are ... The differential diagnosis for these lesions would include:

1. Hand eczema (dyshidrotic eczema) ...
**System:** You are an expert assistant to a dermatology student. Given the following differentials, only return the name of the most likely diagnosis and nothing else. Do not include alternative names of the differential in brackets.
**Content:** Based on the images provided, the key skin findings are ... The differential diagnosis for these lesions would include:

1. Hand eczema (dyshidrotic eczema) ...
**Output2:** hand eczema

### A.3  Image then text 2-calls

Of note, the first API remains the same to the Image only 2-calls case, but the added Additional Information field contains the text query associated with each case in the test set.

**System:** You are an expert assistant to a dermatology student, help him identify exactly what skin conditions would be included in the differential for the images presented. Consider both resemblence and prevalence.
**Content:** IMG_ENC00908_00001.jpg, IMG_ENC00908_00002.jpg
**Output1:** Based on the images provided, the key skin findings are ... The differential diagnosis for these lesions would include:

1. Hand eczema (dyshidrotic eczema) ...
**System:** You are an expert assistant to a dermatology student, given the following differentials discussed and some additional information provided, only return the name of the most likely diagnosis and nothing else. Do not include alternative names of the differential in brackets.
textbfContent: Differentials:

Based on the images provided, the key skin findings are ... The differential diagnosis for these lesions would include:

1. Hand eczema (dyshidrotic eczema) ...
Additional information: Picture 1: On the outside of the thigh, there is a small circle of lump. Approximately 2 months.
Picture 2: Small red spots on the palm. There is slight numbness in the palm. **Output2:** hand eczema (dyshidrotic eczema)

### A.4  Image + text 2-calls

**System:** You are an expert assistant to a dermatology student, help him identify what skin conditions would be included in the differential for the presented images and additional information provided by the medical professional. If any skin conditions are mentioned in the additional information, include them as the most likely differential.
**Content:** Additional information: Picture 1: On the outside of the thigh, there is a small circle of lump. Approximately 2 months.
Picture 2: Small red spots on the palm. There is slight numbness in the palm.
IMG_ENC00908_00001.jpg,
IMG_ENC00908_00002.jpg
**Output1:** Based on the provided images and additional information, here are the potential skin conditions to consider in the differential

diagnosis: ... **System:** You are an expert assistant to a dermatology student. Given the following differentials, only return the name of the most likely diagnosis and nothing else. Do not include alternative names of the differential in brackets. textbfContent: Based on the provided images and additional information, here are the potential skin conditions to consider in the differential diagnosis: ...

**Output2:** picture 1: lipoma. picture 2: palmar erythema

Figure S1: Representative case example illustrating the images of the skin condition, their associated textual query and the predicted response given.



Figure S2: PCA visualization of all the training disease labels embedded by the EmbeddingV3 model. Skin conditions that are semantically similar are clustered together in this representation space.

# LG AI Research & KAIST at EHRSQL 2024: Self-Training Large Language Models with Pseudo-Labeled Unanswerable Questions for a Reliable Text-to-SQL System on EHRs

**Yongrae Jo**[1]*  **Seongyun Lee**[2]*  **Minju Seo**[2]*  **Sung Ju Hwang**[2]  **Moontae Lee**[1]

[1]LG AI Research, [2]KAIST

{yongrae.jo, moontae.lee}@lgresearch.ai  sjhwang82@kaist.ac.kr

## Abstract

Text-to-SQL models are pivotal for making Electronic Health Records (EHRs) accessible to healthcare professionals without SQL knowledge. With the advancements in large language models, these systems have become more adept at translating complex questions into SQL queries. Nonetheless, the critical need for reliability in healthcare necessitates these models to accurately identify unanswerable questions or uncertain predictions, preventing misinformation. To address this problem, we present a self-training strategy using pseudo-labeled unanswerable questions to enhance the reliability of text-to-SQL models for EHRs. This approach includes a two-stage training process followed by a filtering method based on the token entropy and query execution. Our methodology's effectiveness is validated by our top performance in the EHRSQL 2024 shared task, showcasing the potential to improve healthcare decision-making through more reliable text-to-SQL systems.

## 1 Introduction

Electronic Health Records (EHRs) are relational databases storing patients' medical histories within hospitals, covering details from admission to discharge. Common challenges with EHRs include difficulties in documenting and tracking health information, supporting team coordination, and sharing data (Cifuentes et al., 2015). Although ensuring the accurate capture of relevant information is crucial for addressing these challenges, accessing and querying these records often requires knowledge of SQL, making it challenging for healthcare provides in practical settings without technical expertise. A solution to this problem is developing a text-to-SQL model that can translate natural language questions into SQL queries to retrieve information from EHRs.

Recent advancements in Large Language Models (LLMs) have expanded their utility beyond natural language processing to include code generation, enabling them to interpret text for table manipulation and translate descriptions into code effectively (Lee et al., 2024b). These capabilities showcased by code-generating LLMs (Li et al., 2023; Roziere et al., 2023; Guo et al., 2024) demonstrate their potential in text-to-SQL applications. These developments suggest a promising horizon for leveraging LLMs to make EHR data more accessible to healthcare professionals, eliminating the prerequisite of SQL knowledge and significantly simplifying information retrieval (Hwang et al., 2019a; Lyu et al., 2020; Wang et al., 2020b; Park et al., 2021).

However, in the healthcare domain, the reliability of text-to-SQL models is crucial compared to other areas of NLP application. These models must not only generate precise SQL queries from natural language but also identify unanswerable questions—queries that cannot be solved with the available database—to avoid potentially harmful outcomes. The risk of providing answers to such questions underscores the need for these models to err on the side of caution, by preferring not to provide an answer rather than risking the provision of incorrect information (Lee et al., 2023). This approach underlines the unique challenges faced in healthcare NLP, emphasizing the need for accuracy and the ability to recognize when it cannot provide a reliable answer.

Our work introduces PLUQ, an approach leveraging the self-training paradigm to improve the reliability of text-to-SQL models for EHRs through training with **P**seudo-**L**abeled **U**nanswerable **Q**uestions. Self-training, a semi-supervised learning technique, involves re-training a model using its own predictions on unlabeled data to boost its performance. Our method adopts a two-stage version of self-training process, where we initially fine-tune a seed model using a given training

---

*These authors contributed equally to this work.

dataset. We then augment the training dataset by incorporating unanswerable questions that the fine-tuned model identifies from an unlabeled dataset. Subsequently, we fine-tune the model once more using this augmented training data to produce the final model.

Self-training is commonly used in scenarios where unlabeled data is abundant but obtaining labeled data is costly. Text-to-SQL for EHRs exemplifies such a scenario. In this context, a real-world service can collect users' natural language queries without much effort, but determining the correct SQL statement or verifying its answerability with the given database is time-consuming. We adopt the self-training approach to effectively address the issue of class imbalance between answerable and unanswerable questions, thus enhancing the robustness and reliability of the model's performance.

After the two-stage self-training process, we apply a filtering strategy to eliminate uncertain predictions. This strategy employs two types of filtering: one based on the maximum entropy of tokens and the other on the execution results of queries. Specifically, we assess the entropy in each token generated by the language model, designating the prediction's entropy as that of the token with the highest entropy. If a prediction's entropy exceeds a certain threshold, we consider it an unanswerable question, reflecting the model's lack of confidence in providing a correct answer. Additionally, we remove SQL queries that either produce errors or fail to retrieve valid values from the MIMIC-IV[1] dataset.

This approach was validated by our performance in the EHRSQL 2024 shared task (Lee et al., 2024a), where we achieved the top ranking, demonstrating the effectiveness of our method in improving the reliability of text-to-SQL systems in healthcare.

In summary, our study contributes a method that enhances the reliability of text-to-SQL systems for EHRs, addressing the shared task of handling unanswerable questions. This work supports better access to and utilization of EHRs, aiding in informed healthcare decision-making.

The main contributions of our paper are:

1. We propose a self-training method that uses pseudo-labeled unanswerable questions to train text-to-SQL models. This approach helps improve the model's ability to identify

queries it cannot answer accurately, thereby increasing reliability.

2. We detail the comprehensive strategy employed, from the initial prompting of the model to the filtering steps, to ensure the research can be reproduced. This clarity in methodology allows for the approach to be validated and applied by others in the field, enhancing text-to-SQL systems in healthcare.

3. Our method won the EHRSQL 2024 shared task, demonstrating its practical effectiveness in a competitive setting. This success showcases its potential to contribute to the healthcare field by improving access to EHRS through reliable text-to-SQL systems.

## 2 Related Work

In the field of Natural Language Processing (NLP), recent research has focused on text-to-SQL and applying large language models (LLMs) to Electronic Health Records (EHRs). These studies have advanced the handling of complex queries and the processing of healthcare data, setting the stage for our research on test-time data sample labeling and augmentation in EHRs.

**Text-to-SQL** In the evolving field of natural language processing, the development of text-to-SQL technologies represents a significant advancement. Pioneering efforts in this area, Hwang et al. (2019b) and Lyu et al. (2020), harnessed the power of BERT for column classification to tackle the Wiki SQL(Zhong et al., 2017) dataset, which is characterized by its simplistic select/where queries. For more complex scenarios, the SPIDER dataset(Yu et al., 2018), comprising Multi-Table questions, necessitated a understanding of relationship between different tables. Wang et al. (2020a) employed graph-based methods to integrating information, while Lin et al. (2020) introduced schema linking as an input. Moreover, fine-tuning pre-trained language models such as T5 (Raffel et al., 2019) has yielded substantial performance improvements in this field.

**Large Language Models in Text-to-SQL** The emergence of LLMs has inspired novel approaches to text-to-SQL tasks. Dong et al. (2023) introduced efficient zero-shot framework, which capitalize on the robust understanding capabilities of
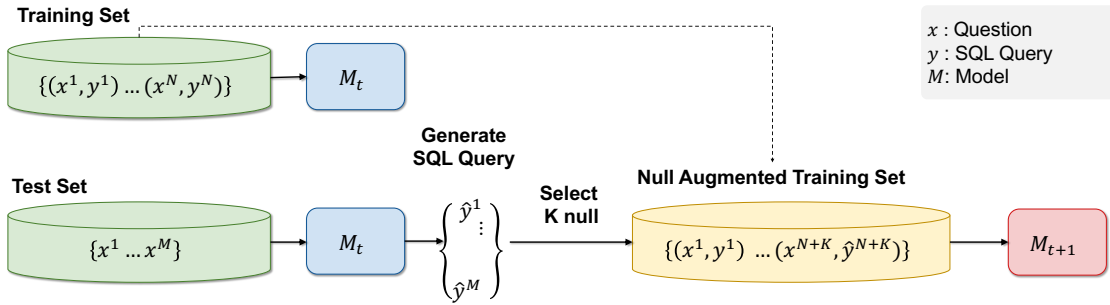
Figure 1: Training Process and SQL Query Generation. The model is initially trained using the training set. Then, a SQL query (or null) is generated for each sample in the test set using the trained model. Subsequently, we select $K$ null samples and add them to the training set, resulting in a null-augmented training set. This augmented dataset is then used to train the final model, denoted as $M_{t+1}$.

LLMs, with a particular emphasis on prompt-based techniques that demonstrates remarkable efficiency. Tai et al. (2023); Nan et al. (2023); Gao et al. (2024) explore optimal demonstrations based on methodologies like text dense similarity or query similarity selection. Pourreza and Rafiei (2023) enhances the integrity of generated SQL through decomposition of queries and self-correction strategies. Shi et al. (2024) proposed LLMs as an agent for generating code and executing it, leverage the few-shot learning capabilities of LLMs for solving the multi-tabular health record datasets.

**Enhancement of LLMs Through Data Augmentation and Self-Training** Amini et al. (2022) presents an extensive review of self-training methods, including consistency-based approaches and transductive learning. Post the rise of LLMs, the field of self-training methods has garnered considerable attention. To enhance the capabilities of LLMs, some studies have focused on autonomous data generation. Wang et al. (2023) stands out by generating synthetic data from a pool consisting not only of seed data but also data generated through an iterative process. Seo et al. (2024) employs a few-shot learning approach, drawing samples from external sources to align the seed data in low-resource settings. The line of autonomously augmenting data broadens and enhances the model's capabilities. Yuan et al. (2024) introduces using their own outputs to continuously improve both their instruction-following and reward-modeling abilities, demonstrating significant performance enhancements over traditional training methods.

**NLP in EHRs** The application of NLP techniques in EHRs has been extensively explored, utilizing texts and structured knowledge. Pampari et al. (2018) proposed a question-answering sys-

tem based on unstructured clinical notes. More recently, many works have been developed in the development of generation task based on structured EHRs. Wang et al. (2020b) construct the table-based QA datasets using MIMIC-III(Johnson et al., 2016) . Park et al. (2021) introduced a graph-based EHR QA system that leverages SPARQL queries from the MIMIC-SQL dataset(Wang et al., 2020b). Raghavan et al. (2021) focuses on QA tasks using the structured patient records in the MIMIC-III. Lee et al. (2023) datasets containing multi-table queries and those involving null values, reflecting real-world scenarios in healthcare domain.

In our research, we utilize a trained LLM on the training dataset for test-time data sample labeling and subsequent augmentation. This approach is particularly focused on addressing the imbalance in the 'null' class.

## 3 Method

We train a seed model using the original training dataset and then use this model for pseudo labeling on the test set. From this, we select only the samples labeled as unanswerable and augment them to the original training dataset to create the final dataset for self-training. Our self-trained model, PLUQ generates SQL queries. We apply post-processing and two stages of filtering to these queries to ensure their reliability and produce the final answers.

### 3.1 Seed model fine-tuning

In developing a model specialized for the text-to-SQL task, we initially fine-tuned seed model on the given training data. Because it is widely recognized that there exists a performance gap between open-source LLMs and proprietary LLMs in many benchmarks. In section 4.3, the results re-

garding performance corresponding to changes in the model substantiate it. Therefore, we utilized the Finetuning API provided by OpenAI to fine-tune the GPT-3.5-Turbo-0125 model.

The training dataset comprises a total of 5,124 samples, including both answerable and unanswerable questions. We employed all of these data samples in our training. Furthermore, to ensure that PLUQ accurately references the correct column names when generating SQL queries, we converted the table schema of the provided MIMIC-IV demo database into text format and incorporated it into the input for training. Additionally, to enhance PLUQ capability in distinguishing between answerable and unanswerable questions, we incorporated information about unanswerable questions into the input. This strategic inclusion aimed at refining PLUQ discernment, thus improving its overall accuracy in classifying questions.

### 3.2 Self-training

**Unanswerable Question Pseudo Labeling** Unanswerable questions refer to queries that either do not align with the given table schema or require external knowledge, rendering them unsolvable using only the MIMIC-IV demo database for SQL query generation. In our training dataset, the number of answerable questions is considerable, reaching 5,124, whereas unanswerable questions are limited to just 450. This disparity highlights a data imbalance issue within our training dataset, which may impede the model's ability to correctly respond to unanswerable questions during testing.

Moreover, there is a low similarity between the queries in the training data and those in the development/test sets. We found that the average cosine similarity between query embeddings in the train and development sets is only 0.36, and between the train and test sets is 0.34, measured using OpenAI's text-embedding-3-large embedding model. Such a disparity in dataset distribution could lead to significant performance declines for the model at test time. To address these issues, we initially perform pseudo labeling on the development/test set using PLUQ, which was originally trained solely on the original training dataset.

**Training With Augmented Data** Pseudo-labeling is one of the techniques used in semi-supervised learning, serving as a powerful tool for addressing issues of data scarcity and label imbalance. Particularly with the EHRSQL

dataset, a notable disparity exists: the quantity of unanswerable questions is significantly lower compared to answerable ones within the training data. Training a model with such data increases the likelihood of the model's inability to accurately respond to unanswerable questions. In tasks where reliability is crucial, especially compared to other domains, this could result in substantial penalties. Therefore, we choose to augment the original training set with those samples predicted as unanswerable. Finally, we fine-tune PLUQ using the augmented dataset.

### 3.3 Filtering

Despite the two-stage training process, including self-training, PLUQ still generates incorrect SQL queries. To enhance the reliability of our final predictions, we implemented a filtering process to sift out samples that were either inaccurately generated or produced with uncertainty by the model. This filtering stage plays a crucial role in ensuring the outputs of PLUQ are more dependable and accurate.

**Maximum Token Entropy Based Filtering** Tokens in a language model-generated output have higher entropy when the information is uncertain. Therefore, treating samples with high entropy as unanswerable questions aids in creating a more reliable system while incurring fewer penalties. We evaluate the entropy of each token produced by the language model, and define the entropy of the prediction based on the token exhibiting the maximum entropy. Then, in the entire set of predictions, samples exceeding a certain entropy level are considered as unanswerable questions and are filtered out. We have set a threshold for this filtering process, determined by the proportion of unanswerable questions in the dataset we aim to predict. This proportion of unanswerable questions is used as a hyperparameter to calibrate the threshold for filtering.

**Execution Based Filtering** Finally, we implement an additional process of filtering to ensure that the remaining SQL queries, after the initial filtering, can successfully access the MIMIC-IV demo database and retrieve valid values. Utilizing the sqlite3 library in Python, we test each SQL query. Queries that trigger errors, return empty values, or yield None are deemed unable to retrieve valid values. Consequently, we filter these queries as unanswerable questions. This step further en-

sures the accuracy and reliability of the system by only allowing queries that can effectively interact with the database.

## 4 Experiments

The experiments are conducted on the development and test sets provided by the EHRSQL 2024 shared tasks. All results presented are derived from runs on the official platform. Section 4.1 details the models, datasets, and metrics used for training and inference, while section 4.2 discusses the experimental results. Finally, in Section 4.3, we conduct ablation studies on various components of PLUQ to examine their individual contributions and impacts.

### 4.1 Settings

**Dataset & Model** We utilize the EHRSQL 2024 dataset for both training and evaluation. The dataset comprises 5,124 training, 1,163 development, and 1,167 test data entries. Notably, only the training dataset is accompanied by gold SQL queries and their corresponding executed gold answers. For questions deemed answerable, it's essential to generate the correct SQL query. For those classified as unanswerable, a null output is required. Database for SQL query generation is the MIMIC-IV demo database. We employ the GPT-3.5-Turbo-0125 model for fine-tuning purposes. Evaluation of our method is conducted on the codabench platform, where we submitted SQL queries predicted by PLUQ for the test set and obtained scores based on their performance.

$$\phi_c(x) = \begin{cases} 1 & \text{if } x \in \mathcal{Q}_{\text{ans}}; g(x) = 1; \text{Acc}(x) = 1 \\ 0 & \text{if } x \in \mathcal{Q}_{\text{ans}}; g(x) = 0, \\ -c & \text{if } x \in \mathcal{Q}_{\text{ans}}; g(x) = 1; \text{Acc}(x) = 0 \\ -c & \text{if } x \in \mathcal{Q}_{\text{una}}; g(x) = 1, \\ 1 & \text{if } x \in \mathcal{Q}_{\text{una}}; g(x) = 0. \end{cases}$$

Figure 2: **Formal Definition of RS** for a single data instance. $Q_{\text{una}}$ denotes unanswerable question, $Q_{\text{ans}}$ represents answerable question. $g(x) = 1$ means that model generates SQL query and $g(x) = 0$ denotes that model generates 'null'. $Acc(x) = 1$ signifies instances where the model's prediction is correct, while $Acc(x) = 0$ indicates cases where the prediction is incorrect. $c$ represents the penalty.

**Metrics** We utilize the Reliability Score (RS) as our primary metric (Lee et al., 2023). In figure 2, the RS aims to accomplish two main objectives: firstly, it provides rewards for correctly generating SQL for answerable questions $Q_{\text{ans}}$ and for not generating SQL for unanswerable questions $Q_{\text{una}}$; secondly, it imposes penalties for wrongly generating SQL for $Q_{\text{ans}}$ and for any attempts to create SQL for $Q_{\text{una}}$. However, the RS neither rewards nor penalizes for choosing not to answer $Q_{\text{ans}}$. The penalties are structured as 0, 5, 10, or N, where N corresponds to the total number of entries in the dataset. The final score is calculated by adding 1 point for each correct sample and deducting points based on the penalty for incorrect ones, followed by averaging these scores. Importantly, in the EHRSQL 2024 shared task, the primary metric for determining rankings is RS(10).

### 4.2 Results

**Development Set** In the development set, PLUQ exhibits the highest performance in RS(10), the primary metric, which positions it at the top of the official leaderboard when compared with other models. A notable aspect of PLUQ is the minimal difference between its RS(0) and RS(10) scores compared to other models. This indicates that PLUQ effectively reduces penalties by categorizing uncertain outcomes in answerable questions and unanswerable questions as 'unanswerable.' This strategy underscores our model's superior reliability, as it avoids the risk of incorrect answers where uncertainty exists, a feature that sets it apart from its counterparts.

**Test Set** In the final ranking phase of the shared task, which utilized the test set, PLUQ experienced a slight overall decrease in scores compared to its performance in the development set. Despite this dip, it maintained a higher score across all RS, including the pivotal RS(10), when compared with other models. This consistent performance across all metrics, even amidst a minor decline, ultimately led PLUQ to win the EHRSQL 2024 shared task.

### 4.3 Ablation Studies

**Model Ablation** We observe the performances across difference models. A total of three models were used, namely Flan-T5-base, Tulu-7b, GPT-3.5-Turbo-0125, and GPT-4-Turbo-Preview. Flan-T5-base, Tulu-7b, and GPT-3.5-Turbo-0125 is fine-tuned, while GPT-4-Turbo-Preview is applied with

Table 1: **Results of the Development and Test Phases** on the Official Codabench Leaderboard. The best results are highlighted in bold. Pivotal metric is RS(10) in this shared task. Note that Ours score for the development phase differs from the official leaderboard because we didn't add it to the leaderboard.

| Team | Development | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| | RS(0) | RS(5) | RS(10) | RS(N) | RS(0) | RS(5) | RS(10) | RS(N) |
| **PLUQ (Ours)** | 90.37 | **89.51** | **88.65** | **-109.6** | **88.17** | **84.75** | **81.32** | **-711.83** |
| PromptMind | 66.38 | 59.5 | 52.62 | -1533.62 | 82.6 | 78.75 | 74.89 | -817.4 |
| ProbGate | 84.18 | 79.45 | 74.72 | -1015.82 | 81.92 | 78.06 | 74.21 | -818.08 |
| KU-DMIS | **91.57** | 82.98 | 74.38 | -1908.43 | 72.07 | 65.64 | 59.21 | -1427.93 |
| oleg1996 | 47.03 | 34.14 | 21.24 | -2952.97 | 68.89 | 56.47 | 44.04 | -2831.11 |
| LTRC-IIITH | N/A | N/A | N/A | N/A | 66.84 | 55.27 | 43.7 | -2633.16 |
| Saama Technologies | 57.78 | 50.47 | 43.16 | -1642.22 | 53.21 | 44.64 | 36.08 | -1946.79 |
| TEAM_optimist | N/A | N/A | N/A | N/A | 14.14 | -349.61 | -713.37 | -84885.86 |



```
System prompt

You are helpful text-to-sql assistant.

--------------------------------------------------------------------------------------------------------------------------

User prompt

You are given SQL table schema and the question. Generate the SQL query for the following question. Note that you should generate 'null' if
the question cannot be converted to SQL query given information.

[SQL Table Schema]
table admissions, columns = [row_id, subject_id, hadm_id, admittime, dischtime, admission_type, admission_location, discharge_location,
insurance, language, marital_status, age]
table chartevents, columns = [row_id, subject_id, hadm_id, stay_id, itemid, charttime, valuenum, valueuom]
table cost, columns = [row_id, subject_id, hadm_id, event_type, event_id, chargetime, cost]
table d_icd_diagnoses, columns = [row_id, icd_code, long_title]
table d_icd_procedures, columns = [row_id, icd_code, long_title]
table d_items, columns = [row_id, itemid, label, abbreviation, linksto]
table d_labitems, columns = [row_id, itemid, label]
table diagnoses_icd, columns = [row_id, subject_id, hadm_id, icd_code, charttime]
table icustays, columns = [row_id, subject_id, hadm_id, stay_id, first_careunit, last_careunit, intime, outtime]
table inputevents, columns = [row_id, subject_id, hadm_id, stay_id, starttime, itemid, amount]
table labevents, columns = [row_id, subject_id, hadm_id, itemid, charttime, valuenum, valueuom]
table microbiologyevents, columns = [row_id, subject_id, hadm_id, charttime, spec_type_desc, test_name, org_name]
table outputevents, columns = [row_id, subject_id, hadm_id, stay_id, charttime, itemid, value]
table patients, columns = [row_id, subject_id, gender, dob, dod]
table prescriptions, columns = [row_id, subject_id, hadm_id, starttime, stoptime, drug, dose_val_rx, dose_unit_rx, route]
table procedures_icd, columns = [row_id, subject_id, hadm_id, icd_code, charttime]
table transfers, columns = [row_id, subject_id, hadm_id, transfer_id, eventtype, careunit, intime, outtime]
foreign_keys = [admissions.subject_id = patients.subject_id, diagnoses_icd.hadm_id = admissions.hadm_id, diagnoses_icd.icd_code =
d_icd_diagnoses.icd_code, procedures_icd.hadm_id = admissions.hadm_id, procedures_icd.icd_code = d_icd_procedures.icd_code,
labevents.hadm_id = admissions.hadm_id, labevents.itemid = d_labitems.itemid, prescriptions.hadm_id = admissions.hadm_id,
cost.hadm_id = admissions.hadm_id, cost.event_id = diagnoses_icd.row_id, cost.event_id = procedures_icd.row_id, cost.event_id =
labevents.row_id, cost.event_id = prescriptions.row_id, chartevents.hadm_id = admissions.hadm_id, chartevents.stay_id = icustays.stay_id,
chartevents.itemid = d_items.itemid, inputevents.hadm_id = admissions.hadm_id, inputevents.stay_id = icustays.stay_id, inputevents.itemid
= d_items.itemid, outputevents.hadm_id = admissions.hadm_id, outputevents.stay_id = icustays.stay_id, outputevents.itemid =
d_items.itemid, microbiologyevents.hadm_id = admissions.hadm_id, icustays.hadm_id = admissions.hadm_id, transfers.hadm_id =
admissions.hadm_id]

Question: {question}
SQL Query:
```

Figure 3: The system prompt and the user prompt template used in PLUQ. The prompt integrates instructions for handling unanswerable questions and the MIMIC-IV database schema.

in-context learning. All results are conducted on the development set.

Among the fine-tuned models, GPT-3.5-Turbo-0125 demonstrates the highest performance. This indicates that there is still a performance gap be-tween proprietary and open-source models. Furthermore, despite having more parameters, Tulu-7b shows lower performance compared to Flan-T5-base. Additionally, it is observed that GPT-4-Turbo-Preview, known for its high performance in nu-

Table 2: **Model Ablation Study** of the Development Set across the finetuned open-source LLMs and in-context learning, finetuned proprietary LLMs. FT denotes the fine-tuning of the model, while ICL represents in-context learning (Wei et al., 2022). In this work, the number of few-shot examples used for in-context learning is fixed at 4.

| Models | RS(0) | RS(5) | RS(10) | RS(N) |
|--------|-------|-------|--------|-------|
| Flan-T5-base FT | 82.11 | 76.53 | 70.94 | -1217.8 |
| Tulu-7b FT | 10.23 | -38.77 | -87.9 | -11389.7 |
| GPT-4-Turbo-Preview ICL | 63.52 | -118.85 | -301.22 | -186836.4 |
| GPT-3.5-Turbo-0125 FT (Ours) | **90.37** | **89.51** | **88.65** | **-109.6** |

Table 3: **Prompt Ablation Study** of the Development Set including the integration of table schema and the incorporation of unanswerable information to evaluate the impact of various prompts on model performance. The base model of fine-tuning is GPT-3.5-Turbo-0125 model.

| Models | RS(0) | RS(5) | RS(10) | RS(N) |
|--------|-------|-------|--------|-------|
| Fine-Tuning | 83.23 | 78.5 | 73.77 | -1016.7 |
| + Table Schema | 89.85 | 83.83 | 77.82 | -1310.1 |
| + Unans Info (Ours) | **90.37** | **89.51** | **88.65** | **-109.6** |

Table 4: **Filtering Ablation Study** of Development Set. For maximum token entropy based filtering, we filtered out SQL queries possessing high entropy within the top 7%, classifying them as unanswerable questions.

| Models | RS(0) | RS(5) | RS(10) | RS(N) |
|--------|-------|-------|--------|-------|
| No Filtering | 80.82 | 5.58 | -69.94 | -17419.1 |
| + Exec Filtering | **93.98** | **89.68** | 85.38 | -906.01 |
| + Ent Filtering (Ours) | 90.37 | 89.51 | **88.65** | **-109.6** |

merous benchmarks, scored lower than fine-tuned models when only in-context learning is applied.

**Prompt Ablation**  In the study, we compare the performance of models based on the information included in the input prompts during training. When table schema information is incorporated into the prompts, the models perform better than without it. This suggests that providing table schema information, such as column names, offers a valuable learning signal to the models.

Additionally, explicitly including information about unanswerable questions results in higher scores than when such information is omitted. By providing criteria for answerable and unanswerable questions, the models are aided in avoiding questions they could not answer and focusing on providing accurate responses to those that are answerable.

In the final version of the prompt, we incorporated the database schema of MIMIC-IV as well as

the instruction related to unanswerable questions. You can find the prompt in Figure 3.

**Filtering Ablation**  In table 4, by applying execution filtering, which treats invalid SQL queries that either do not execute or retrieve empty values as unanswerable questions, a significant performance improvement is observed, particularly in scenarios with substantial penalties such as RS(10) and RS(N). Additionally, by implementing entropy-based filtering, which filters out SQL queries with higher entropy than a set threshold among those with high maximum token entropy, performance is further enhanced by effectively eliminating SQL queries that, even when executed, return incorrect values.

## 5  Conclusion

In our work, we develop a self-training strategy designed to enhance the reliability of text-to-SQL models for Electronic Health Records (EHRs) through the inclusion of pseudo-labeled unanswerable questions. This approach is particularly valuable in scenarios where there is an abundance of unlabeled data and labeling is costly, thus providing substantial clinical utility in real-world applications. Our approach employs a two-stage training process alongside a filtering mechanism based on token entropy and query execution outcomes to improve the model's precision and its ability to identify unanswerable questions. The performance is validated by our leading performance in the EHRSQL 2024 shared task. Our method contributes towards rendering EHRs more accessible to healthcare professionals without SQL knowledge, addressing a critical need for reliable information retrieval in healthcare. Future research could explore how large language models facilitate the integration of unstructured medical texts into specific schemas, enhancing interoperability in varied healthcare settings.

## Limitations

Our method achieve the best score in this challenge, as we adopt various techniques to enhance reliability. However, there are some limitations to our approach. Since our model is fine-tuned using EHRs, its ability to generalize across the entire EHR dataset is limited. Additionally, the fine-tuning process requires training data, which poses a challenge due to the high costs and time associated with data collection. Furthermore, despite

achieving the highest score among all teams, our RS(N) score still remains negative, indicating that caution should be exercised when considering the application of our method in real-world scenarios.

## Acknowledgements

## References

Massih-Reza Amini, Vasilii Feofanov, Loïc Pauletto, Emilie Devijver, and Yury Maximov. 2022. Self-training: A survey. *CoRR*, abs/2202.12040.

Maribel Cifuentes, Melinda Davis, Doug Fernald, Rose Gunn, Perry Dickinson, and Deborah J Cohen. 2015. Electronic health record challenges, workarounds, and solutions observed in practices integrating behavioral health and primary care. *The Journal of the American Board of Family Medicine*, 28(Supplement 1):S63–S72.

Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, Lu Chen, Jinshu Lin, and Dongfang Lou. 2023. C3: zero-shot text-to-sql with chatgpt. *CoRR*, abs/2307.07306.

Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2024. Text-to-sql empowered by large language models: A benchmark evaluation. *Proc. VLDB Endow.*, 17(5):1132–1145.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. 2024. Deepseek-coder: When the large language model meets programming–the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.

Wonseok Hwang, Jinyeong Yim, Seunghyun Park, and Minjoon Seo. 2019a. A comprehensive exploration on wikisql with table-aware word contextualization. *arXiv preprint arXiv:1902.01069*.

Wonseok Hwang, Jinyeung Yim, Seunghyun Park, and Minjoon Seo. 2019b. A comprehensive exploration on wikisql with table-aware word contextualization. *CoRR*, abs/1902.01069.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2023. EHRSQL: A practical text-to-sql benchmark for electronic health records. *CoRR*, abs/2301.07695.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024a. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Younghun Lee, Sungchul Kim, Tong Yu, Ryan A Rossi, and Xiang Chen. 2024b. Learning to reduce: Optimal representations of structured data in prompting large language models. *arXiv preprint arXiv:2402.14195*.

Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.

Xi Victoria Lin, Richard Socher, and Caiming Xiong. 2020. Bridging textual and tabular data for cross-domain text-to-sql semantic parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 4870–4888. Association for Computational Linguistics.

Qin Lyu, Kaushik Chakrabarti, Shobhit Hathi, Souvik Kundu, Jianwen Zhang, and Zheng Chen. 2020. Hybrid ranking network for text-to-sql. *CoRR*, abs/2008.04759.

Linyong Nan, Yilun Zhao, Weijin Zou, Narutatsu Ri, Jaesung Tae, Ellen Zhang, Arman Cohan, and Dragomir Radev. 2023. Enhancing few-shot text-to-sql capabilities of large language models: A study on prompt design strategies. *CoRR*, abs/2305.12586.

Anusri Pampari, Preethi Raghavan, Jennifer J. Liang, and Jian Peng. 2018. emrqa: A large corpus for question answering on electronic medical records. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2357–2368. Association for Computational Linguistics.

Junwoo Park, Youngwoo Cho, Haneol Lee, Jaegul Choo, and Edward Choi. 2021. Knowledge graph-based question answering with electronic health records. In *Proceedings of the Machine Learning for Healthcare Conference, MLHC 2021, 6-7 August 2021, Virtual Event*, volume 149 of *Proceedings of Machine Learning Research*, pages 36–53. PMLR.

Mohammadreza Pourreza and Davood Rafiei. 2023. DIN-SQL: decomposed in-context learning of text-to-sql with self-correction. 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*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou,

Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *CoRR*, abs/1910.10683.

Preethi Raghavan, Jennifer J. Liang, Diwakar Mahajan, Rachita Chandra, and Peter Szolovits. 2021. emrkbqa: A clinical knowledge-base question answering dataset. In *Proceedings of the 20th Workshop on Biomedical Language Processing, BioNLP@NAACL-HLT 2021, Online, June 11, 2021*, pages 64–73. Association for Computational Linguistics.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.

Minju Seo, Jinheon Baek, James Thorne, and Sung Ju Hwang. 2024. Retrieval-augmented data augmentation for low-resource domain tasks. *CoRR*, abs/2402.13482.

Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce C. Ho, Carl Yang, and May D. Wang. 2024. Ehragent: Code empowers large language models for complex tabular reasoning on electronic health records. *CoRR*, abs/2401.07128.

Chang-Yu Tai, Ziru Chen, Tianshu Zhang, Xiang Deng, and Huan Sun. 2023. Exploring chain of thought style prompting for text-to-sql. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 5376–5393. Association for Computational Linguistics.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2020a. RAT-SQL: relation-aware schema encoding and linking for text-to-sql parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7567–7578. Association for Computational Linguistics.

Ping Wang, Tian Shi, and Chandan K. Reddy. 2020b. Text-to-sql generation for question answering on electronic medical records. In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*, pages 350–361. ACM / IW3C2.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 13484–13508. Association for Computational Linguistics.

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.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 3911–3921. Association for Computational Linguistics.

Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. *CoRR*, abs/2401.10020.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *CoRR*, abs/1709.00103.

# Overview of the EHRSQL 2024 Shared Task on
# Reliable Text-to-SQL Modeling on Electronic Health Records

**Gyubok Lee    Sunjun Kweon    Seongsu Bae    Edward Choi**
KAIST AI
{gyubok.lee, sean0042, seongsu, edwardchoi}@kaist.ac.kr

## Abstract

Electronic Health Records (EHRs) are relational databases that store the entire medical histories of patients within hospitals. They record numerous aspects of patients' medical care, from hospital admission and diagnosis to treatment and discharge. While EHRs are vital sources of clinical data, exploring them beyond a predefined set of queries requires skills in query languages like SQL. To make information retrieval more accessible, one strategy is to build a question-answering system, possibly leveraging text-to-SQL models that can automatically translate natural language questions into corresponding SQL queries and use these queries to retrieve the answers. The EHRSQL 2024 shared task aims to advance and promote research in developing a question-answering system for EHRs using text-to-SQL modeling, capable of reliably providing requested answers to various healthcare professionals to improve their clinical work processes and satisfy their needs. Among more than 100 participants who applied to the shared task, eight teams completed the entire shared task processes and demonstrated a wide range of methods to effectively solve this task. In this paper, we describe the task of reliable text-to-SQL modeling, the dataset, and the methods and results of the participants. We hope this shared task will spur further research and insights into developing reliable question-answering systems for EHRs.

## 1 Introduction

Electronic Health Records (EHRs) store all types of medical events that occur in the hospital, including hospital admissions, diagnoses, procedures, prescriptions, and discharges. They replace traditional paper-based records and provide a centralized repository for patient data. Over the years, the widespread adoption of EHRs in hospitals has been shown to improve patient care, increase efficiency, and enhance coordination among healthcare professionals (Upadhyay and Hu, 2022; Mullins et al.,

2020; Uslu et al., 2021). Although EHRs are a valuable source of patient data, the complexity of their data structures and the need for specialized skills, such as query languages like SQL, to extract and analyze the information, often hinder their effective utilization by healthcare professionals (Wang et al., 2020; Lee et al., 2022). These barriers lead to the underutilization of the full potential of EHRs in clinical practice and research.

An alternative way to utilize data stored in EHRs is to develop a question-answering (QA) system. QA systems provide a user-friendly interface that allows healthcare professionals to ask questions in natural language and receive relevant answers from the EHR data, without needing to know query languages or EHR database strctures. Specifically, text-to-SQL modeling is an effective approach for building QA systems for EHRs, which are typically relational databases. These models automatically convert natural language questions into their corresponding SQL queries, and then execute these queries on the database to obtain the final answer. With the impressive advances in large language models (LLMs), various high-performance text-to-SQL models have been introduced, which are accomplished through model fine-tuning (Scholak et al., 2021) or LLM prompting with demonstrations (Pourreza and Rafiei, 2024; Gao et al., 2023; Chang and Fosler-Lussier, 2023). If deployed with reliable performance, these models could significantly benefit healthcare professionals by allowing them to explore patient data more freely from the EHRs through natural language interactions.

Several datasets on question-answering for EHRs have been introduced, including MIMIC-SQL (Wang et al., 2020), emrKBQA (Raghavan et al., 2021), and EHRSQL (Lee et al., 2022). EHRSQL, in particular, poses unique challenges. It is the first dataset to compile a collection of questions that reflect the diverse needs of healthcare professionals, including physicians, nurses, and
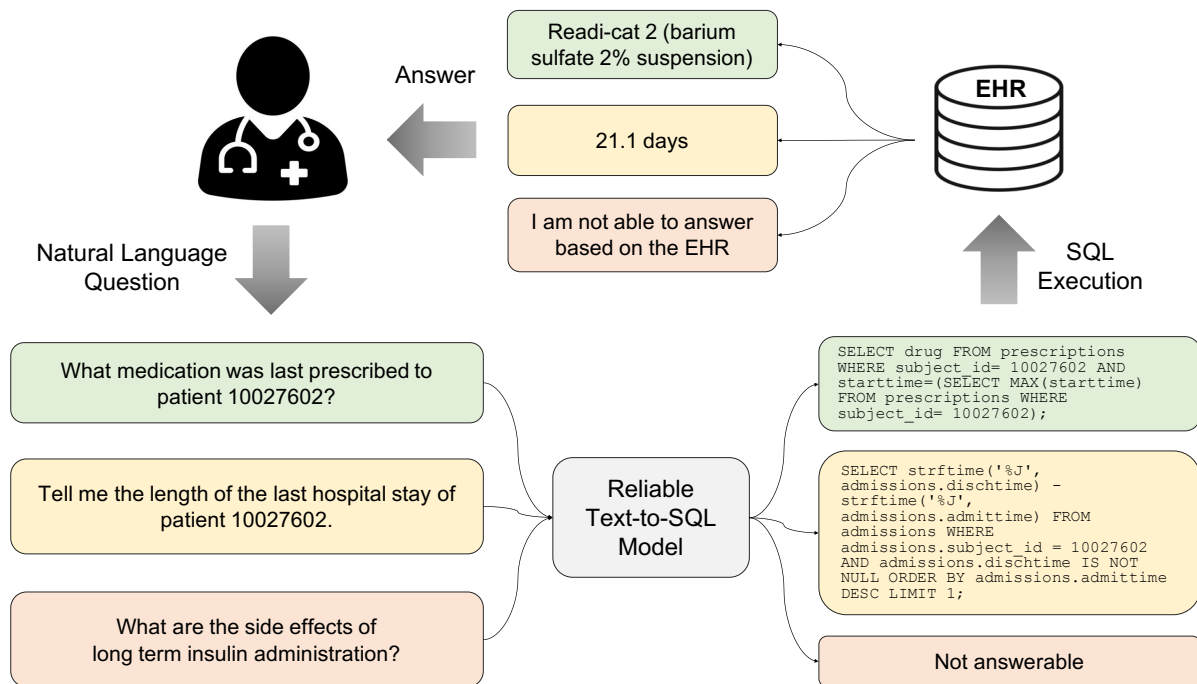
644

Figure 1: Overview of reliable text-to-SQL modeling on EHRs. For any input questions, a reliable text-to-SQL model should accurately predict SQL queries for what it can and abstain from what it cannot, such as for intrinsically unanswerable questions or ones that are likely incorrect by the model. Successfully developing such a model can serve as a valuable tool for healthcare professionals in hospitals for better accessibility of patient data and assistance of clinical decision-making.

hospital administrative staff. It contains extensive use of time expressions and includes SQL queries of increased complexity, which better reflect the real needs of a hospital setting. The SQL queries are linked to two open-source EHR databases[1], MIMIC-III (Johnson et al., 2016) and eICU (Pollard et al., 2018), retaining incompatible ones as unanswerable questions in the dataset (used to test a model's ability to abstain). Starting from their collected real-world questions, this shared task presents more up-to-date changes to the text-to-SQL modeling (use of MIMIC-IV and new paraphrases for questions) and more challenging problem settings (new data splitting and additional unanswerable questions). The dataset for this shared task is publicly available at https://github.com/glee4810/ehrsql-2024. The shared task platform is hosted on Codabench at https://www.codabench.org/competitions/1889/.

In this paper, we present the EHRSQL 2024

shared task and its dataset in Sections 2 and 3, respectively. Section 4 introduces the evaluation metric and baseline model for the task. Section 5 describes methods proposed by the participating teams and discusses interesting findings from the official results.

## 2 Task - Reliable Text-to-SQL Modeling

The goal of the EHRSQL 2024 shared task is to develop a reliable QA system for EHRs, specifically through text-to-SQL modeling. Reliability is crucial for the deployment of AI systems, especially in critical domains like hospitals, where incorrect predictions can have severe consequences. The term reliability in question answering refers to the system's preference for abstention over providing an incorrect answer (Whitehead et al., 2022; Chen et al., 2023; Lee et al., 2024). In this shared task, we adopt the definition of *reliable text-to-SQL* from TrustSQL (Lee et al., 2024), which first expands the scope of reliability to include unanswerable questions. A reliable text-to-SQL model should not only correctly generate SQL queries, providing utility, but also abstain from answering questions that are likely to be incorrect or are unanswerable, thereby

---

[1]SQL queries are database-dependent, meaning that even though a question attempts to retrieve the same information, the location of that information can vary across databases. For example, to list all drugs in MIMIC-III, you would use SELECT drug FROM prescriptions, whereas in eICU, it would be SELECT drugname FROM medication.

Figure 2: Data construction process of the EHRSQL shared task.

minimizing harm. This objective contrasts with most other text-to-SQL tasks, where the primary focus is to maximize SQL generation performance for answerable questions only. Further discussion on specific scenarios of measuring reliability for text-to-SQL is explained in Section 4.1.

## 3 Dataset Construction

In this section, we outline the key steps to generate data for the shared task. The overall data construction pipeline is illustrated in Figure 2. Each subsection provides a detailed explanation of each step.

### 3.1 Question Templates from EHRSQL

To construct the shared task data, we started from the pool of questions that reflect the real needs of diverse healthcare professionals in EHRSQL (Lee et al., 2022). This dataset is derived from the results of a poll participated in by more than 200 professionals at a university hospital in South Korea. The collected questions are those that the professionals would ask an AI speaker if it could access and synthesize structured information stored in EHRs (i.e., records in tabular form). The authors then translated the raw question utterances and removed duplicate ones to distill them into question templates. This shared task leverages the question templates collected in EHRSQL to generate diverse and realistic question-SQL pairs.

### 3.2 SQL Queries linked to MIMIC-IV Demo

Unlike the original EHRSQL dataset whose SQL queries are based on value-shuffled MIMIC-III and eICU[2], this shared task uses the demo version of MIMIC-IV[3] (Johnson et al., 2020) to construct question-SQL pairs. The demo version, containing

records of 100 patients from the full MIMIC-IV database, has the same database schema as the full MIMIC-IV and is openly-available for anyone who is interested in using the dataset without special training[4]. Since the demo database schema is identical to the full database, the same query can be used to retrieve information from both the full and demo versions.

### 3.3 New Question Paraphrases

We found that the style and naturalness of paraphrases generated by current LLMs, like ChatGPT, surpass the paraphrases in EHRSQL, which are produced through both human and machine efforts. To improve the quality of the paraphrases for each question template, we employed ChatGPT to generate new paraphrases that are more natural and conversational. We then manually reviewed all new paraphrases to ensure they maintain the intended meaning of the original question templates.

### 3.4 Challenging Unanswerable Questions

A recent study revealed that unanswerable questions in the EHRSQL dataset can mostly be filtered out using a combination of N-gram and beam search score filtering (Yang et al., 2024). This is primarily because the unanswerable questions in EHRSQL were collected erroneously due to human errors during the polling process[5], resulting in limited diversity. To increase the difficulty of the task, we combined the original unanswerable questions with those from the EHRSQL portion of TrustSQL (Lee et al., 2024), which contains adversarially created unanswerable questions, such as those referring to non-existing columns and requests that exceed SQL functionalities.

---

[2]This process was done to further de-identify the question-SQL pairs for public release. Please refer to more detailed reasons in the original paper.

[3]https://physionet.org/content/mimic-iv-demo/2.2/

[4]The full MIMIC-IV dataset requires researchers to complete the Collaborative Institutional Training Initiative (CITI) training before accessing the data.

[5]The poll participants were initially provided with examples of inappropriate questions for the system, including those requiring external knowledge, ambiguous or qualitative statements, and questions about the reasons behind certain clinical decisions.

| | Dev Phase | | Test Phase |
|---|---|---|---|
| | Train | Valid | Test |
| Answerable question template | 100 (100 seen) | 134 (100 seen + 34 unseen) | 134 (100 seen + 34 unseen) |
| Answerable samples | 4674 | 931 | 934 |
| Unanswerable samples | 450 | 232 | 233 |
| Total samples | 5124 | 1163 | 1167 |

Table 1: Data statistics for the shared task. All text-to-SQL data used in the shared task is based on MIMIC-IV.

## 3.5 New Data Split

In real-world scenarios, text-to-SQL models can encounter questions that are answerable based on the EHR schema but have not been seen in the training set (unseen SQL with respect to the training set). This situation can lead to increased confusion for the model in distinguishing between answerable and unanswerable questions. Unlike the original EHRSQL, where answerable questions were split in an identically distributed (IID) manner, we split the shared task data to include both seen and unseen question templates (or SQL structures) in the validation and test sets. For unanswerable questions, the original unanswerable questions from EHRSQL were distributed across all splits (training, validation, and test), while new unanswerable questions were added exclusively to the validation and test sets to increase the task's difficulty. Each of these splits has a 20% proportion of unanswerable questions. Table 1 shows the number of question templates and the size of each data split[6]. The training and validation sets were made available during the development phase (Jan 29, 2024 - Mar 26, 2024), and the test set was made available for the three-day test phase (Mar 26, 2024 - Mar 28, 2024).

---

[6]Even if the MIMIC-IV demo includes only 100 patients, a wide variety of question templates can exist. Consider patient ID 100 and two question templates: 'What is patient 100's gender?' and 'What is patient 100's last blood pressure?' The data splitting in text-to-SQL for EHRs does not have to be done by patient, such as 'What is patient 100's gender?' in the training set and 'Tell me patient 200's sex?' in the validation set, because the task could become relatively easy. Instead, it might include 'What is patient 100's gender?' in the training set and 'What is patient 100's last blood pressure?' in the validation set. A more challenging and realistic goal of text-to-SQL is to assess how well the model can generate SQL queries for both question templates (or SQL structures) that it has seen and those it has not seen. In this example, we show four question samples with two question templates.

## 4 Evaluation

### 4.1 Evaluation Metric

We chose the evaluation metric that best aligns with the purpose of our shared task: to build reliable text-to-SQL models aimed at accurately predicting correct SQL queries and identifying unanswerable questions, while minimizing incorrect SQL predictions and the wrongly classifying unanswerable questions as answerable. More concretely, we adopt the Reliability Score (RS) for reliable text-to-SQL (Lee et al., 2024), formally written as follows:

$$RS(c)(x) = \begin{cases} 1 & \text{if } x \in \mathcal{Q}_{ans}; g(x) = 1; Acc(x) = 1, \\ 0 & \text{if } x \in \mathcal{Q}_{ans}; g(x) = 0, \\ -c & \text{if } x \in \mathcal{Q}_{ans}; g(x) = 1; Acc(x) = 0, \\ -c & \text{if } x \in \mathcal{Q}_{una}; g(x) = 1, \\ 1 & \text{if } x \in \mathcal{Q}_{una}; g(x) = 0, \end{cases}$$
(1)

where $\mathcal{Q}_{ans}$ and $\mathcal{Q}_{una}$ denote answerable and unanswerable questions, respectively. $g(x) = 1$ implies that the model selects its SQL generation as the final answer, whereas $g(x) = 0$ implies that the model abstains. $Acc(x)$ indicates the accuracy of the generated SQL, based on execution accuracy, which is determined by whether the answers returned by the ground-truth and predicted SQL queries match.

The RS has five different cases for assigning the score:

- A score of $1$ is assigned if SQL is correctly generated by the model for answerable questions.

- A score of $0$ is assigned if the model abstains from generating SQL for answerable questions.

- A score of $-c$ is assigned if the model predicts incorrect SQL for answerable questions.

- A score of $-c$ is assigned if the model attempts to predict SQL for unanswerable questions.

- A score of $1$ is assigned if the model accurately detects unanswerable questions by abstaining.

The overall RS is calculated by taking the average of sample-level scores, represented in percentages. The penalty of c is chosen depending on the reliability requirements of the model. A penalty of 0 (RS_0) means no punishment for incorrect

| | Team | Affiliation | Paper | Code |
|---|---|---|---|---|
| 1 | LG AI Research & KAIST | LG AI Research & KAIST, South Korea | Jo et al. (2024) | [1] |
| 2 | PromptMind | - | Gundabathula and Kolar (2024) | [2] |
| 3 | ProbGate | KAIST, South Korea | Kim et al. (2024b) | [3] |
| 4 | KU-DMIS | Korea university, South Korea | Kim et al. (2024a) | [4] |
| 5 | AIRI NLP | AIRI, Russia | Somov et al. (2024) | [5] |
| 6 | LTRC-IIITH | IIIT Hyderabad, India | Thomas et al. (2024) | [6] |
| 7 | Saama Technologies | Saama Technologies, USA | Jabir et al. (2024) | [7] |
| 8 | TEAM_optimist | SUST, Bangladesh | Joy et al. (2024) | [8] |

[1] `https://github.com/sylee0520/ehrsql-2024` (private)
[2] `https://github.com/satyakesav/ehrsql-clinicalnlp-2024` (private)
[3] `https://github.com/venzino-han/probgate_ehrsql`
[4] `https://github.com/Chanwhistle/EHRSQL_NACCL`
[5] `https://github.com/runnerup96/EHRSQL-text2sql-solution`
[6] `https://github.com/jr-john/ehrsql_2024` (private)
[7] `https://github.com/upjabir/ehrsql_2024`
[7] `https://github.com/joy-2019331037/nlpConference`

Table 2: Participating teams, affiliation, paper, and code.

predictions, a penalty of 10 (RS_10) represents a moderately rigorous scenario, and a penalty of N (RS_N), where N refers to the evaluation data size, is the most rigorous scenario in which even a single mistake outweighs all correct predictions and abstentions. The maximum possible RS is $100\%$, and the minimum possible scores vary depending on the penalties: 0 for $c = 0$; $-1000\%$ for $c = 10$; $-100N\%$ for $c = N$. The main evaluation metric for the shared task is RS(10), where every ten accurate predictions weigh the same as one incorrect prediction.

## 4.2 Code Verification and Fact Sheet

The participants shared their code and the fact sheet following the instructions reported in Appendix A. The purpose of the fact sheet was to collect a brief summary of participants' methods, including any use of pre-trained models or external data. For code verification, participants had the option to submit their code either via email or through GitHub repositories. These repositories could be public or private, as long as access was granted to the task organizers. Upon receiving the submissions, we conducted a careful review to ensure that the provided code and the methods described in the fact sheets are consistent.

## 4.3 Baseline Model

For the baseline, we employ the simplest method, denoted as ABSTAIN-ALL, which abstains from answering all questions. Evaluating in the RS, abstaining from all questions results in an overall score of 20%. This score is earned by correctly abstaining from answering unanswerable questions. This is not a trivial score, particularly as the penalty for incorrect predictions increases, which can severely harm the overall score.

## 5 Official Results

### 5.1 Participating Teams

The EHRSQL shared task attracted over 100 participants from both academia and industry. Of these, 8 teams submitted their code and fact sheet. Table 2 lists the participating teams, their affiliations, the code submission status (not all of which is publicly available), and their working papers.

### 5.2 Methods and Results

Table 3 presents the official results for each team, along with short descriptions of their methods. The proposed methods can be categorized into two types: unified and pipeline-based approaches ('Modeling Type' in Table 3). The unified approach leverages LLMs to perform both SQL generation and abstention, while the pipeline-based approach

| | Team | RS_0 | RS_10 | RS_N | Modeling Type | Ensemble | Fine-tuned | Model Used |
|---|---|---|---|---|---|---|---|---|
| 1 | LG AI Research & KAIST | 88.17 | 81.32 | -711.83 | Unified | Yes | No | ChatGPT |
| 2 | PromptMind | 82.6 | 74.89 | -817.4 | Unified | Yes | Yes | GPT-4, ChatGPT, Claude Opus |
| 3 | ProbGate | 81.92 | 74.21 | -818.08 | Unified | No | Yes | ChatGPT |
| 4 | KU-DMIS | 72.07 | 59.21 | -1427.93 | Unified | No | Yes | ChatGPT |
| 5 | AIRI NLP | 68.89 | 44.04 | -2831.11 | Pipeline | No | Yes | T5-3B, Logistic Regression |
| 6 | LTRC-IIITH | 66.84 | 43.7 | -2633.16 | Pipeline | No | Yes | SQLCoder-7b-2 |
| 7 | Saama Technologies | 53.21 | 36.08 | -1946.79 | Pipeline | Yes | Yes | Decision Trees, CodeLlama-7b, ChatGPT |
| 8 | TEAM_optimist | 14.14 | -713.37 | -84.9K | Unified | No | No | SQLCoder-7b-2 |
| - | ABSTAIN-ALL | 20.0 | 20.0 | 20.0 | No | No | No | - |

Table 3: Official results. ABSTAIN-ALL is the baseline for the shared task, explained in Section 4.3. 'Ensemble' denotes the use of any ensemble methods. 'Fine-tuned' indicates whether any pre-trained models were further trained for SQL generation or abstention purposes. 'Pipeline-based' means the use of multiple methods in a sequence, such as a pipeline that consists of an answerability detector, an SQL generator, and subsequently an SQL error detector.

involves building a series of specialized, smaller models, such as SQLCoder or T5-3B, to ensure reliability as one system. The overall observation is that 1) methods that fall under the unified approach tend to outperform those in the pipeline-based approaches; 2) most teams chose to fine-tune LLMs on the training data, either general-purpose (e.g., ChatGPT) or code-specialized models (e.g., CodeLLama), highlighting the importance of domain-specific fine-tuning for adapting LLMs to this task; 3) teams with smaller discrepancies between the RS with different penalties (e.g., the gap between RS(0) and RS(10)) tend to rank higher, indicating that minimizing incorrect SQL predictions through effective abstention mechanisms is crucial for this task. Detailed discussions of each method by category are provided in the following paragraphs.

**Unified approach.** Five teams utilized methods under the unified approach. The LG AI Research & KAIST team achieved the best results, scoring 81.32 in RS(10) by using self-training LLMs (Amini et al., 2022; Yuan et al., 2024) with pseudo-labeling for unanswerable questions. The PromptMind team implemented an ensemble of LLMs, including fine-tuned ChatGPT, GPT-4, and Claude Opus. They selected SQL generation as the final prediction only if all three models unanimously agreed; otherwise, they would abstain. For SQL generation, they employed two retrievers (one for the general domain and another for the medi-

cal domain) to retrieve similar question-SQL pairs from the training set. The ProbGate team employed fine-tuned ChatGPT with log-probability thresholding and error handling for abstention, where the threshold was set heuristically based on the ratio of unanswerable questions in the validation set. The KU-DMIS team took a two-stage method. First, they generated question-SQL pairs to align the test set distribution with the training set using question templates from the original EHRSQL. Then, they fine-tuned ChatGPT on this newly generated dataset. Abstention was achieved by sampling multiple SQL predictions for each input question and abstaining if the outputs were not consistent. Lastly, the TEAM_optimist team used SQLCoder-7b-2 for direct generation of SQL and abstention labels (null) through in-context learning.

**Pipeline-based approach.** Alternatively, three teams adopted the pipeline-based approach. The AIRI NLP team used a two-stage method: initially using logistic regression to detect unanswerable questions, then generating SQL with a fine-tuned T5-3B (Raffel et al., 2020), and finally checking the executability of the generated SQL for final abstention. The LTRC-IIITH team used two different SQLCoder-7b-2 models, one for detecting unanswerable questions and the other for SQL generation. For final abstention, they utilized the log-probabilities from the SQL generator to detect potential errors in SQL generation, followed by an executability check of the SQL. The Saama Technolo-

gies team began with an ensemble of unanswerable question detectors, including multinomial naive bayes, SGD classifier, CatBoost (Prokhorenkova et al., 2018), and CodeLlama-7b (Roziere et al., 2023). They then generated SQL using CodeLlama-7b, and finally used a ChatGPT-based answer selector for final abstention.

# 6 Conclusion

With the increasing volume of data stored in EHRs and the impressive advances in LLMs, the EHRSQL 2024 shared task offered an opportunity to develop and test participants' creative methods to building reliable QA systems on EHRs using text-to-SQL modeling. The dataset for this shared task presents unique challenges, including questions that extensively use time expressions and the increased complexity of SQL queries, which more accurately reflect the real needs of a hospital setting. It also includes challenging unanswerable questions that should be avoided. This distinguishes the task from most other text-to-SQL challenges, as reliable text-to-SQL models must not only generate correct SQL queries, providing utility, but also abstain from answering questions that are likely incorrect or unanswerable, thereby minimizing harm.

The shared task attracted over 100 participants from academia and industry, with 8 teams ultimately submitting their code and fact sheets. As a novel task at the intersection of the NLP and clinical domains, it inspired a variety of proposed methods. These included self-training LLMs through pseudo-labeling, ensembling of different LLMs, generating synthetic question-SQL pairs to handle distribution shifts from training to test sets, leveraging log-probabilities for abstention, and pipeline-based approaches with specialized models for correct SQL generation and abstention. We hope that this shared task, emphasizing reliability, will encourage further research into building QA systems for EHRs that can truly serve as valuable tools for healthcare professionals in hospitals, improving clinical decision-making, facilitating research, and enhancing patient care quality. Future research directions include expanding reliable question answering for EHRs to multimodal settings by incorporating clinical notes, X-ray images, and ECG signals.

# Limitations

This shared task does not represent all types of answerable and unanswerable questions encountered in hospital settings. Additionally, this shared task employs MIMIC-IV as the EHR database, which is not a universally accepted EHR schema, and the databases are preprocessed for the QA task by eliminating duplicate values across different tables to reduce ambiguity. Lastly, further experiments are necessary for newly proposed LLMs, since most methods, including text-to-SQL generation and abstention, depend heavily on the underlying LLMs.

# References

Massih-Reza Amini, Vasilii Feofanov, Loic Pauletto, Emilie Devijver, and Yury Maximov. 2022. Self-training: A survey. *arXiv preprint arXiv:2202.12040.*

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. In *NeurIPS 2023 Second Table Representation Learning Workshop.*

Jiefeng Chen, Jinsung Yoon, Sayna Ebrahimi, Sercan Arik, Tomas Pfister, and Somesh Jha. 2023. Adaptation with self-evaluation to improve selective prediction in llms. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5190–5213.

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. *arXiv preprint arXiv:2308.15363.*

Satya Gundabathula and Sriram Kolar. 2024. Promptmind team at ehrsql-2024: Improving reliability of

sql generation using ensemble llms. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Mohammed Jabir, Kamal Kanakarajan, and Malaikannan Sankarasubbu. 2024. Saama technologies at ehrsql 2024: Sql generation through classification answer selector by llm. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Yongrae Jo, Seongyun Lee, and Minju Seo. 2024. Lg ai research kaist at ehrsql 2024: Self-training large language models with pseudo-labeled unanswerable questions for a reliable text-to-sql system on ehrs. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2020. Mimic-iv. *PhysioNet. Available online at: https://physionet. org/content/mimiciv/1.0/(accessed August 23, 2021)*, pages 49–55.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

Sourav Bhowmik Joy, Rohan Redwan, Argha Pratim Saha, Minhaj Ahmed, Utsho Das, and Partha Sarothi Bhowmik. 2024. Team optimist at ehrsql 2024: Text-to-sql generation using large language model for ehr analysis. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Chanhwi Kim, Hajung Kim, Hoonick Lee, Jiwoo Lee, Kyochul Jang, Kyungjae Lee, Gangwoo Kim, and Jaewoo Kang. 2024a. Ku-dmis: Generating sql query via question templatization in ehr. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Sangryul Kim, Donghee Han, and Sehyun Kim. 2024b. Probgate at ehrsql 2024: Enhancing sql query generation accuracy through probabilistic threshold filtering and error handling. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Gyubok Lee, Woosog Chay, Seonhee Cho, and Edward Choi. 2024. Trustsql: A reliability benchmark for text-to-sql models with diverse unanswerable questions. *arXiv preprint arXiv:2403.15879*.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Alexandra Mullins, Renee O'Donnell, Mariam Mousa, David Rankin, Michael Ben-Meir, Christopher Boyd-Skinner, and Helen Skouteris. 2020. Health outcomes and healthcare efficiencies associated with the use of electronic health records in hospital emergency departments: a systematic review. *Journal of Medical Systems*, 44(12):200.

Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. 2018. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data*, 5(1):1–13.

Mohammadreza Pourreza and Davood Rafiei. 2024. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Advances in Neural Information Processing Systems*, 36.

Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. 2018. Catboost: unbiased boosting with categorical features. *Advances in neural information processing systems*, 31.

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.

Preethi Raghavan, Jennifer J Liang, Diwakar Mahajan, Rachita Chandra, and Peter Szolovits. 2021. emrK-BQA: A clinical knowledge-base question answering dataset. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 64–73, Online. Association for Computational Linguistics.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.

Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. Picard: Parsing incrementally for constrained auto-regressive decoding from language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9895–9901.

Oleg Somov, Elena Tutubalina, and Alexei Dontsov. 2024. Airi nlp team at ehrsql 2024: T5 and logistic regression to the rescue. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Jerrin John Thomas, Pruthwik Mishra, Dipti Sharma, and Parameswari Krishnamurthy. 2024. Ltrc-iiith at ehrsql 2024: Enhancing reliability of text-to-sql systems through abstention and confidence thresholding.

In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Soumya Upadhyay and Han-fen Hu. 2022. A qualitative analysis of the impact of electronic health records (ehr) on healthcare quality and safety: Clinicians' lived experiences. *Health Services Insights*, 15:11786329211070722.

Aykut Uslu, Jürgen Stausberg, et al. 2021. Value of the electronic medical record for hospital care: update from the literature. *Journal of medical Internet research*, 23(12):e26323.

Ping Wang, Tian Shi, and Chandan K Reddy. 2020. Text-to-sql generation for question answering on electronic medical records. In *Proceedings of The Web Conference 2020*, pages 350–361.

Spencer Whitehead, Suzanne Petryk, Vedaad Shakib, Joseph Gonzalez, Trevor Darrell, Anna Rohrbach, and Marcus Rohrbach. 2022. Reliable visual question answering: Abstain rather than answer incorrectly. In *European Conference on Computer Vision*, pages 148–166. Springer.

Yongjin Yang, Sihyeon Kim, SangMook Kim, Gyubok Lee, Se-Young Yun, and Edward Choi. 2024. Towards unbiased evaluation of detecting unanswerable questions in EHRSQL. In *ICLR 2024 Workshop on Navigating and Addressing Data Problems for Foundation Models*.

Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*.

## A  Code Submission and Fact Sheet Template

# Fact Sheet for EHRSQL-2024 Shared Task:

### I. Team name

- Username on Codalab:

- Team leader affiliation:

- Team leader email:

- Name of other team members (and affiliation):

- Team website URL (if any):

### II. Contribution

- **Title of the contribution**
    - Provide a brief summary of the method and contributions.

- **Representative image / workflow diagram of the method**
    - An image (or several images) to support method description to better understand the approach and model pipeline. You can refer to these images in the method description part.

- **Detailed method description**
    - Provide a technical and detailed description of the method and contributions. The explanations must be self-contained and one must be able to reproduce the approach by reading this section.

- **Shared task results**
    - $RS_0$:
    - $RS_5$:
    - $RS_{10}$:
    - $RS_N$:

- **Final Remarks**
    - Please identify the pros and cons (if any) of the proposed approach.

### III. Additional method details

- Did you use any pre-trained model?

- Did you use external data?

- Did you perform any data augmentation?

- At the test phase, did you use the provided validation set as part of your training set?

- Did you use any regularization strategies/terms?

- Did you use handcrafted features?

- Did you use any domain adaptation strategy?

**IV. Code Repository**

- Link to a code repository with complete and detailed instructions so that the results obtained on Codabench can be reproduced.

- If private repo, share the repo with `glee4810`

# Saama Technologies at EHRSQL 2024: SQL Generation through Classification Answer Selector by LLM

**Mohammed Jabir, Kamal Raj Kanakarajan, Malaikannan Sankarasubbu**
Saama Technologies
{mohammed.jabir, kamal.raj, malaikannan.sankarasubbu}@saama.com

## Abstract

The EHRSQL task aims to develop a dependable text-to-SQL model for Electronic Health Records (EHR) databases, which are crucial sources of clinical data that store patients' medical histories in hospitals. Large language models (LLM) have been proven to exhibit state-of-the-art performance for text-to-SQL tasks across various domains. To this end, we have developed a framework, SQL Generation through Classification Answer Selector by LLM (SCAS), which comprises two modules. The CAS module determines the answerability of the question, while the SG model generates the SQL query exclusively for answerable questions. Our system ranked 7th on the leaderboard with a Reliability Score of 53.21 on the official test set.

## 1 Introduction

Electronic Health Records (EHRs) are an essential component of modern healthcare. They store a patient's complete medical history, allowing hospital staff to make better clinical decisions (Wang et al., 2020; Bardhan et al., 2022) by quickly accessing relevant patient information. However, accessing this information can be time-consuming, especially when complex queries are involved. The traditional way of accessing EHRs involves using a predefined rule conversion system to convert user queries to SQL and retrieve the relevant information. This process can become a bottleneck for users who need to build custom queries or deal with complex queries. To address this issue, the EHRSQL (Lee et al., 2022) task aims to develop a system that can automatically translate user questions into corresponding SQL queries, making retrieving the information they need easier and quicker. The system's objective is to build a text-to-SQL system that converts natural language queries to SQL and informs users whether their queries are answerable.

Text-to-SQL tasks (Katsogiannis-Meimarakis and Koutrika, 2023) involve mapping natural language questions onto a given relational database into SQL queries. Early studies (Dong and Lapata, 2016; Wang et al., 2019) tackled this task with pre-defined rules or as a sequence-to-sequence task. However, recent advancements in large language models (LLMs) (Brown et al., 2020; OpenAI et al., 2024; Touvron et al., 2023) have become a milestone for natural language processing and machine learning. LLMs are pre-trained on massive text corpus, which enables them to perform various natural language tasks, and their ability to do in-context learning (Liu et al., 2021) makes them most suitable for text-to-SQL generation.

In this paper, we present our approach to tackling the EHRSQL 2024 (Lee et al., 2024) shared task, which involves a complex dataset of electronic health records. Our proposed framework, the SQL Generation through Classification Answer Selector by LLM (SCAS), helps avoiding incorrect SQL generation by using the Classification Answer Selector (CAS) module. The CAS module uses an LLM prompting method that incorporates the output of other classification models to generate the final classification output, thereby abstaining from incorrect responses. The SCAS framework generates SQL queries only for necessary questions by utilizing other LLM models. Our system achieved a 7th position on the leaderboard, with a Reliability Score of 53.21 on the official test set. The code to reproduce the experiments mentioned in this paper is publicly available[1].

## 2 Background

### 2.1 Task and Dataset Description

EHRSQL is a text-to-SQL task aiming to convert natural language queries into corresponding SQL queries while identifying untranslatable ones. The

---

[1] https://github.com/upjabir/ehrsql_2024

original EHRSQL (Lee et al., 2022) task was built on MIMIC-III (Johnson et al., 2016) and EICU (Pollard et al., 2018) datasets, which are available from Physionet (Goldberger et al., 2000). The EHRSQL dataset contains a wide range of questions across various domains in EHR, including Demographics, Prescription, Vital signs, and more, as well as Time Sensitive questions. The dataset for the task comprises 5124 training, 1163 validation, and 1162 testing samples.The competition dataset is derived from the MIMIC-IV (Johnson et al., 2023) open-access database demo subset and includes both answerable and unanswerable questions in the training set. The system should output corresponding SQL queries for answerable questions and null for unanswerable questions.

The Reliability Score (RS) is a new evaluation metric for text-to-SQL models used in the EHRSQL task. It rewards accurate SQL generation for certain types of questions while penalizing incorrect SQL generation for others. It does not assign any reward or penalty for abstaining from answering certain questions. The competition uses several scoring systems, including RS(0), RS(5), RS(10), and RS(N), with RS(10) being the primary metric for the leaderboard. In RS(10), correct predictions receive one positive point, while incorrect predictions receive -10 points. N in RS(N) represents the size of the test set.

## 2.2 Related Works

The text to SQL task poses a significant challenge and has previously been approached as a sequence-to-sequence task. (Brunner and Stockinger, 2021) utilized the BERT (Devlin et al., 2019) model as an encoder architecture to achieve state-of-the-art results in this task. They incorporated user questions and employed a neural network architecture to extract values and generate SQL queries. Another study of (Qi et al., 2022) demonstrates an innovative approach by incorporating various types of existing relations and co-references, thereby introducing new parameters to the encoder-decoder (Sutskever et al., 2014) architecture model.

Researchers have been utilizing the Language Model LLM for text-to-SQL since its emergence. Downstream tasks for LLM can be achieved through in-context learning and fine-tuning methods. (Wei et al., 2023) proposed a Chain of Thought style prompting technique to enhance the capabilities of LLM. (Pourreza and Rafiei, 2023) proposed a decomposed in-context learning method where

the text-to-SQL task is divided into subtasks. On the other hand (Tai et al., 2023) introduced a new CoT-style prompting method specifically for text-to-SQL parsing, which showed significant improvements compared to standard prompting methods and the least-to-most prompting method. Additionally, a new prompt engineering method called DIAL SQL was proposed by (Gao et al., 2023), demonstrating the potential of fine-tuning LLMs for Text-to-SQL while highlighting the degeneracy of in-context learning capability after fine-tuning.

Text classification is a crucial task in machine learning, and it can be accomplished using classical machine learning models such as Random Forest and Deep learning models like Transformer (Vaswani et al., 2023), which is also effective in handling complex language tasks. With the emergence of LLM, which is trained on large text corpus, (Wang et al., 2023) suggests that the efficiency of text classification has been increased by using LLM as a zero-shot classifier. Although using LLM for downstream tasks is quite challenging. (Sun et al., 2023) addresses the difficulty of using LLM for downstream tasks by implementing effective prompting techniques, thereby improving the efficiency of LLM in text classification. On the other hand, (Zhang et al., 2024) overcomes this challenge by fine-tuning LLM, resulting in impressive performance surpassing in-context zero-shot learning capabilities of pre-trained LLM models like GPT 4 (OpenAI et al., 2024) in the healthcare domain.
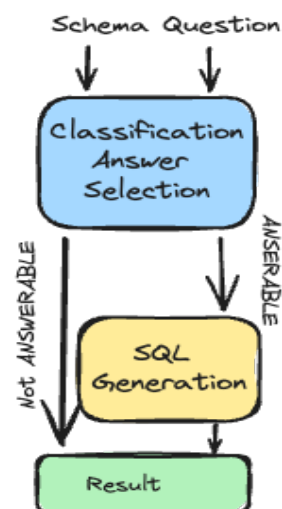
## 3 System Overview



Figure 1: SCAS framework input and output flow.

Our research paper presents a novel framework

comprising two modules: The Classification Answer Selector (CAS) module and the SQL generation (SG) module. The CAS module is responsible for determining whether a question can be answered, while the SG module is designed to generate an SQL query for questions.

## 3.1 Classification Answer Selector Module

The CAS (Classification Answer Selector) module is a powerful tool that includes two distinct classification models and a selector to generate the final classification answer. The selector utilizes the advanced capabilities of Azure's OpenAI GPT-3.5-turbo (Brown et al., 2020) LLM to ensure the accuracy and comprehensiveness of the final answer.

The classification model uses a classical machine learning approach and methodology for the classification task. Feature selection is used to reduce the dataset's dimensionality by eliminating irrelevant features, with TF-IDF (Sparck Jones, 1972) as an effective methodology for text classification. Multiple classification models were employed, including MultinomialNB (Lewis, 1998), SGD (Robbins, 1951; Kiefer and Wolfowitz, 1952), and CatBoost (Dorogush et al., 2018), for ensemble classification using a weighted ensemble approach. Predictions were based on probabilities, with a threshold of 0.4 established to convert predicted probabilities into class labels.

$$ClassLabel = \begin{cases} 1 & \text{if probability} > 0.4 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

The Second classification model is a fine-tuned LLM specifically for classification tasks. We utilize a pre-trained Large language model from the Codellama family (Rozière et al., 2024), specifically CodeLlama-7b-Instruct-hf[2], for this task. In this task, we have a large language model M and a training dataset $D=\{x_i,y_i\}$, where $x_i$ represents the input prompt and $y_i$ is the class label. The goal is to minimize the weighted cross-entropy loss, which is calculated by dividing the total number of instances in the training data by twice the count of positive or negative target values:

$$l_n = -w_{y_n} \log \frac{\exp(x_{n,y_n})}{\sum_{c=1}^{C} \exp(x_{n,c})} \quad (2)$$

The selector takes the outputs from both the first and second classification models and utilizes an in-context learning method to determine the final classification result. The prompt used for the selector is shown in figure 2. A full listing of examples is available in Appendix E

---

Based on the database schema and table description, determine which AI assistant's answer accurately identifies whether the given question can generate an SQL query or not.

### Database Table Description
The table name and its corresponding description are as follows:
{table description}

### Database Schema
This query will run on a database whose schema is represented in this string:
{schema}

{few shots}
Question: "{question}"
Ai Assitant 1's Answer: {model1 answer}
Ai Assitant 2's Answer: {model2 answer}
Answer: Let's think step by step.

---

Figure 2: Prompt for Classification Answer Selector.

## 3.2 SQL Generation Module

We employ the same pre-trained large language model of the Codellama family, which is used in classifier tasks for SQL generation. We perform instruction tuning only by considering the answerable questions from the EHRSQL 2024 dataset. The dataset, denoted as $D=\{x_i\}$, consists of input prompts where $x_i$ represents the input prompt. The training objective is Causal language modeling. A full listing of prompts and examples are shown in Appendix D

## 4 Finetuning

We fine-tuned the model using the efficient parameter tuning method LoRA (Hu et al., 2022) and the HuggingFace library (Mangrulkar et al., 2022). The finetuning process for both the CAS and SQL generation modules involved using the AdamW optimizer and a cosine learning rate scheduler, targeting all linear layers within the model. We em-

---

[2]https://huggingface.co/meta-llama/CodeLlama-7b-Instruct-hf

| Exp No | Development Phase | RS(0) | RS(5) | RS(10) | RS(N) |
|--------|-------------------|-------|-------|--------|-------|
| 1 | gpt-3.5-turbo | 38.60 | -261.47 | -561.56 | -69761.39 |
| 2 | codellama FD | 40.84 | -250.21 | -541.27 | -67659.15 |
| 3 | codellama SA + CL1 | **74.97** | 46.17 | 17.36 | -6625.02 |
| 4 | codellama SA + CL1$^*$ | 42.73 | 40.15 | 37.57 | **−557.26** |
| 5 | codellama SA + CL2 | 43.59 | 39.29 | 34.99 | 956.40 |
| 6 | codellama SA + CAS | 55.97 | **49.52** | **43.07** | -1444.02 |
| **Exp No** | **Test Phase** | | | | |
| 1 | codellama SA + CAS | **77.03** | -36.07 | -149.18 | -26322.96 |
| 2 | codellama SA + CAS$^*$ | 53.21 | **44.64** | **36.07** | **-1946.786** |

Table 1: Experimental results during development phase and test phase. FD means Full data set used for training, SA means only Selected Answerable Questions. CL1 is the classical model, CL2 is codellama model, and CAS is the Classical Answer Selector Module. $^*$ Adjusted threshold.

ployed a maximum sequence length of 4096 tokens for training and inference for the SQL generation task, using beam-search decoding strategies with 4 beams during inference. The hyperparameters utilized for the finetuning process are outlined in Appendix A, while Appendix B and C detail the dataset preprocess and postprocess methods employed for fine-tuning. Additionally, every fine-tuning process was done using a 4× Quadro RTX 8000 (48GB VRAM) card.

# 5 Results

We tried both fine-tuning and in-context learning of LLM for this task. We established a baseline for our experiment using the gpt-3.5-turbo model, which received a RS(10) score of -561.56 points. To enable in-context learning for the model, we used few-shot prompting. In the second experiment, we fine-tuned the codellama model, and despite having only 7B parameters, it outperformed the gpt-3.5-turbo model. Notably, both the question classifier and SQL generation in both 1 & 2 experiments used the same model. Experiment 3 showcased the robustness of our SQL generation module, as it achieved an impressive RS(0) score of 74.97. This indicates that our system can correctly generate executable queries 74.97% of the time. Experiment 4 focused on improving the question classification model by adjusting the threshold to identify unanswerable questions better. While this enhanced the RS(10) score by 23.48 points, it caused a decrease in the RS(0) score due to misclassification. Experiment 5 evaluated the performance of the codellama-based question classifier, which showed no significant improvement over classical models. Finally, in experiment 6, we used the CAS mod-

ule that combines the result of two classification models to enhance the gpt-3.5-turbo model's performance. The input the CAS module is detailed in section 2. After completing the development phase, we submitted our top-performing model for the testing phase, scoring RS(0) of 53.21 and RS(10) of 36.07 points. The SQL generator module achieved an impressive RS(0) score of 77.03 points during the test phase, without any adjustments made to the Question Classifier threshold, which shows the capabilities of the SQL generator module.

After performing an error analysis on the CAS module, it was discovered that false negatives were higher. This indicates that some answerable questions were incorrectly classified as unanswerable. Since our system is designed as a pipeline model, only the questions classified as answerable will advance to the SQL generation model. This resulted in a decrease in the RS(0) score. Notably, most of the false negative predictions were observed in queries related to test procedures, hospitals and departments.

# 6 Conclusion

We developed a sophisticated system that can generate SQL queries from user queries in the EHR dataset, provided they are convertible to SQL. Our system was able to achieve an impressive rank of 7 on the EHRSQL task. Our experiments have shown that fine-tuning an LLM for task-specific SQL generation significantly enhances its performance compared to in-context learning. However, we acknowledge that our system needs improvement in identifying which user queries can be successfully converted to SQL. This is crucial for ensuring the reliability of our SQL converter system.

To facilitate the reproducibility of our work, we have made available instruction templates, code, and pre-trained models as open-source resources.

## Limitations

Our research specifically focused on Codellama models, and we discovered that models fine-tuned on text-to-SQL tasks, such as SQLCoder[3], did not perform well in the EHRSQL task. Additionally, there is a limited amount of data available to train for unanswerable questions within the provided training data, with only 450 out of 5124 questions being unanswerable. In future work, generating synthetic data for unanswerable questions using models like GPT-4 could potentially improve performance. It is important to note that all experiments were conducted using Codellama 7B models.

## References

Jayetri Bardhan, Anthony Colas, Kirk Roberts, and Daisy Zhe Wang. 2022. Drugehrqa: A question answering dataset on structured and unstructured electronic health records for medicine related queries. *arXiv preprint arXiv:2205.01290*.

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. *Preprint*, arXiv:2005.14165.

Ursin Brunner and Kurt Stockinger. 2021. Valuenet: A natural language-to-sql system that learns from database information. *Preprint*, arXiv:2006.00888.

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.

Li Dong and Mirella Lapata. 2016. Language to logical form with neural attention. *Preprint*, arXiv:1601.01280.

Anna Veronika Dorogush, Vasily Ershov, and Andrey Gulin. 2018. Catboost: gradient boosting with categorical features support. *Preprint*, arXiv:1810.11363.

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.

Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. 2000. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220.

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*.

Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2023. Mimic-iv clinical database demo.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9.

George Katsogiannis-Meimarakis and Georgia Koutrika. 2023. A survey on deep learning approaches for text-to-sql. *The VLDB Journal*, 32(4):905–936.

Jack Kiefer and Jacob Wolfowitz. 1952. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*, pages 462–466.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

David D Lewis. 1998. Naive (bayes) at forty: The independence assumption in information retrieval. In *European conference on machine learning*, pages 4–15. Springer.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Preprint*, arXiv:2107.13586.

Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul. 2022. PEFT: State-of-the-art parameter-efficient fine-tuning methods. https://github.com/huggingface/peft.

---

[3]https://huggingface.co/defog/sqlcoder-7b-2

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gompani, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. 2018. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data*, 5(1):1–13.

Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Preprint*, arXiv:2304.11015.

Jiexing Qi, Jingyao Tang, Ziwei He, Xiangpeng Wan, Yu Cheng, Chenghu Zhou, Xinbing Wang, Quanshi Zhang, and Zhouhan Lin. 2022. Rasat: Integrating relational structures into pretrained seq2seq model for text-to-sql. *Preprint*, arXiv:2205.06983.

Herbert E. Robbins. 1951. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407.

Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code. *Preprint*, arXiv:2308.12950.

Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.

Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023. Text classification via large language models. *Preprint*, arXiv:2305.08377.

Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. *Preprint*, arXiv:1409.3215.

Chang-You Tai, Ziru Chen, Tianshu Zhang, Xiang Deng, and Huan Sun. 2023. Exploring chain-of-thought style prompting for text-to-sql. *Preprint*, arXiv:2305.14215.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2019. Rat-sql: Relation-aware schema encoding and linking for text-to-sql parsers. *arXiv preprint arXiv:1911.04942*.

Ping Wang, Tian Shi, and Chandan K Reddy. 2020. Text-to-sql generation for question answering on electronic medical records. In *Proceedings of The Web Conference 2020*, pages 350–361.

Zhiqiang Wang, Yiran Pang, and Yanbin Lin. 2023. Large language models are zero-shot text classifiers. *Preprint*, arXiv:2312.01044.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

Xiaodan Zhang, Nabasmita Talukdar, Sandeep Vemulapalli, Sumyeong Ahn, Jiankun Wang, Han Meng, Sardar Mehtab Bin Murtaza, Dmitry Leshchiner, Aakash Ajay Dave, Dimitri F. Joseph, Martin Witteveen-Lane, Dave Chesla, Jiayu Zhou, and Bin Chen. 2024. Comparison of prompt engineering and fine-tuning strategies in large language models in the classification of clinical notes. *medRxiv*.

## A  Hyperparameter

| Hyperparameter | Codellama SA | CL2 |
|---|---|---|
| Learning rate | 2.5e-5 | 1e-4 |
| Batch size | 4 | 4 |
| Epochs | 3 | 3 |
| Weight decay | 0.01 | 0.01 |
| Weight for loss | - | 0.54, 5.69 |
| Lora rank (r) | 16 | 16 |
| Lora rank alpha ($\alpha$) | 32 | 32 |
| Lora rank dropout | 0.05 | 0.1 |

Table 2: Hyperparameter used for the best performing model.

Hyperparameters used by the best performing pre-trained language model are listed in Table 2, and the total hyperparameter search space is listed in Table 3. Also, the hyperparameter for the classification model is listed in the Table 4.

| Hyperparameter | Value |
|---|---|
| Learning rate | 2.5e-5, 5e-5, 1e-4, 2e-4 |
| Batch size | 4, 8, 16, 32 |
| Epochs | 1-5 |
| Weight decay | 0.01, 0.02, 0.05, 0.1 |
| Lora rank (r) | 8, 16, 32, 64 |
| Lora rank alpha ($\alpha$) | 16, 24, 32 |
| Lora rank dropout | 0.05, 0.08, 0.1 |

Table 3: Full list of hyperparameter search space for finetuning LLM

## B  Dataset Preprocess

For classification, we derive a binary dataset from the raw EHRSQL dataset, which contains only two classes based on the question's answerability. The class label is one if the question is answerable; otherwise, it is zero. We use the raw question and its answerability for the first classification model to create a dataset of questions and their respective class labels. For the second classification model, we map the input to the format of figure 3, which involves providing the SQL schema with foreign keys and the question itself. The class label determination remains the same as in the first classification model. Figure 2 showcases the prompt for the selector in the CAS module, which includes $< dti, dsi, qi, fi, cii, cij >$. We select a few-shot example based on the cosine similarity between the

| Model | HyperParameters | values |
|---|---|---|
| MultinomialNB Classifier | $alpha$ | 0, 0.01, **0.02**, 1.0 |
| | $fit\_prior$ | **True**,False |
| SGD classifier | $loss$ | **modified_huber**, log_loss,huber |
| | $max\_iter$ | 1000, 5000, **8000** |
| | $tol$ | **1e-4**, 1e-5, 2e-5 |
| | $penalty$ | **l2**, l1 |
| CatBoost classifier | $learning\_rate$ | 0.01, **0.0056**, 0.01, 0.2 |
| | $depth$ | 4,5,**6**,8 |
| | $l2\_leaf\_reg$ | 1, 4, **6.5**, 8.5, 10 |
| | $subsample$ | 0.1, **0.3**, 0.5, 1 |
| | $loss$ | LogLoss, **CrossEntropy** |

Table 4: Hyperparameter space for the classification experiments. Hyperparameters in bold are what we used for the our classification models

given question ($qi$) and the entire set of training questions using a pre-trained sentence transformer called all-mpnet-base-v2 [4]. From this process, we identify the four most similar training questions along with their corresponding SQL query, which will serve as our few-shot examples.For SQL generation, we formatted the raw data into the format of figure 2. The prompt includes $< qi, dti, dsi, qsi >$, where $qi$ is the question, $dti$ is the database table information, $dsi$ is the schema with foreign key details, and $qsi$ is the SQL query for the corresponding question.

## C Dataset Postprocess

In order to guarantee that the classification models' outputs are effectively conveyed to the selector within the CAS module, a postprocessing step must be incorporated. This step entails modifying the classification output: if the output is 1, it is transformed to "Able to generate answer"; otherwise, it is converted to "Unable to generate answer". Furthermore, the generated output in the SQL generation module is trimmed to include only the section between the [SQL] and [SQL] keywords.

---

[INST] ### Task
Generate a SQL query to answer [QUESTION]question[/QUESTION].

### Database Table Description
The table name and its corresponding description are as follows:
{table description}

### Database Schema
This query will run on a database whose schema is represented in this string:
{schema}

### Answer
Given the database schema, here is the SQL query that answers [QUESTION]question[/QUESTION]
[SQL]{sql}[/SQL]

Figure 3: Prompt for SQL generation.

## D Prompt and Examples for SQL generation module

Prompt in the figure 3 is used to train and inference pre-trained large language model for the SQL generation task. Given below is a full-fledged example for SQL generation prompt.

**Example 1**

Based on the database schema and table description, determine which AI assistant's answer accurately identifies whether the given question can generate an SQL query or not.
### Database Table Description
The table name and its corresponding description are as follows:
ADMISSIONS – Every unique hospitalization for each patient in the database
PATIENTS – Every unique patient in the database
D_ICD_DIAGNOSES – International Statistical Classification of Diseases and Related Health Problems (ICD-9) codes relating to diagnoses
D_ICD_PROCEDURES – International Statistical Classification of Diseases and Related Health Problems (ICD-9) codes relating to procedures
D_LABITEMS – Local codes ('ITEMIDs') appearing in the database that relate to laboratory tests
D_ITEMS – Local codes ('ITEMIDs') appearing in the database, except those that relate to laboratory tests
DIAGNOSES_ICD – Hospital assigned diagnoses, coded using the International Statistical Classification of Diseases and Related Health Problems (ICD) system
PROCEDURES_ICD – Patient procedures, coded using the International Statistical Classification of Diseases and Related Health Problems (ICD) system
LABEVENTS – Laboratory measurements for patients both within the hospital and in outpatient clinics
PRESCRIPTIONS – Medications ordered for a given patient
COST – All patients events cost
CHARTEVENTS – All charted observations for patients
INPUTEVENTS – Intake for patients monitored while in the ICU
OUTPUTEVENTS – Output information for patients while in the ICU
MICROBIOLOGYEVENTS – Microbiology culture results and antibiotic sensitivities from the hospital database

ICUSTAYS – Every unique ICU stay in the database
TRANSFERS – Patient movement from bed to bed within the hospital

Database Schema
This query will run on a database whose schema is represented in this string:
CREATE TABLE patients
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of the patient
subject_id INT NOT NULL UNIQUE, – Unique subject id of the patient
gender VARCHAR(5) NOT NULL, – Gender of the patient
dob TIMESTAMP(0) NOT NULL, – Date of birth of the patient
dod TIMESTAMP(0) – Date of death of the patient
);
CREATE TABLE admissions
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of the admission
subject_id INT NOT NULL, – Subject id of the admission
hadm_id INT NOT NULL UNIQUE, – Unique hospital admission id of the admission
admittime TIMESTAMP(0) NOT NULL, – Admit time of the admission
dischtime TIMESTAMP(0), – Discharge time of the admission
admission_type VARCHAR(50) NOT NULL, – Admission type of the admission
admission_location VARCHAR(50) NOT NULL, – Admission location of the admission
discharge_location VARCHAR(50), – Discharge location of the admission
insurance VARCHAR(255) NOT NULL, – Insurance of the admission
language VARCHAR(10), – Langauge of the admission
marital_status VARCHAR(50), – Marital status of the admission
age INT NOT NULL, – Age of the admission
);

663

CREATE TABLE d_icd_diagnoses
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of the icd diagnose
icd_code VARCHAR(10) NOT NULL
UNIQUE, – Unique icd code of the icd
diagnose
long_title VARCHAR(255) NOT NULL –
Title of the icd diagnose
);
CREATE TABLE d_icd_procedures
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of icd procedure
icd_code VARCHAR(10) NOT NULL
UNIQUE, – Unique icd code of the icd
procedure
long_title VARCHAR(255) NOT NULL –
Title of the icd procedure
);
CREATE TABLE d_labitems
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of the item relate to laboratory
tests
itemid INT NOT NULL UNIQUE, –
Unique item id of the item relate to
laboratory tests
label VARCHAR(200) – Label of the item
relate to laboratory tests
);
CREATE TABLE d_items
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of the item excepts item relate
to laboratory tests
itemid INT NOT NULL UNIQUE, –
Unique item id of the item excepts item
relate to laboratory tests
label VARCHAR(200) NOT NULL, –
Label of item excepts item relate to
laboratory tests
abbreviation VARCHAR(200) NOT NULL,
– Abbreviation of item excepts item relate
to laboratory tests
linksto VARCHAR(50) NOT NULL –
Event linked to item excepts item relate to
laboratory tests
);
CREATE TABLE diagnoses_icd

(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of diagnose
subject_id INT NOT NULL, – Subject id
of diagnose
hadm_id INT NOT NULL, – Hospital
admission id of diagnose
icd_code VARCHAR(10) NOT NULL, –
ICD code of diagnose
charttime TIMESTAMP(0) NOT NULL, –
Chart time of diagnose
);
CREATE TABLE procedures_icd
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of procedures
subject_id INT NOT NULL, – Subject id
of procedures
hadm_id INT NOT NULL, – Hospital
admission id of procedures
icd_code VARCHAR(10) NOT NULL, –
ICD code of procedures
charttime TIMESTAMP(0) NOT NULL, –
Chart time of procedures
);
CREATE TABLE labevents
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of laboratory event
subject_id INT NOT NULL, – Subject id
of laboratory event
hadm_id INT NOT NULL, – Hospital
admission id of laboratory event
itemid INT NOT NULL, – Item id of
laboratory event
charttime TIMESTAMP(0), – Chart time of
laboratory event
valuenum DOUBLE PRECISION, – Nu-
merical value measured of laboratory event
valueuom VARCHAR(20), – Unit of
numerical value of laboratory event
);
CREATE TABLE prescriptions
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of prescription
subject_id INT NOT NULL, – Subject id
of prescription
hadm_id INT NOT NULL, – Hospital
admission id of prescription

starttime TIMESTAMP(0) NOT NULL, –
Start time of prescription
stoptime TIMESTAMP(0), – Stop time of
prescription
drug VARCHAR(255) NOT NULL, – Drug
name of prescription
dose_val_rx VARCHAR(100) NOT NULL,
– Dosage value of prescription
dose_unit_rx VARCHAR(50) NOT NULL,
– Dosage unit of prescription
route VARCHAR(50) NOT NULL, – Intake
method of prescription
);
CREATE TABLE cost
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of cost event
subject_id INT NOT NULL, – Subject id
of cost event
hadm_id INT NOT NULL, – Hospital
admission id of cost event
event_type VARCHAR(20) NOT NULL, –
Event type of cost event
event_id INT NOT NULL, – Event id of
cost event
chargetime TIMESTAMP(0) NOT NULL, –
Charge time of cost event
cost DOUBLE PRECISION NOT NULL, –
Cost of cost event
);
CREATE TABLE chartevents
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of chart event
subject_id INT NOT NULL, – Subject id
of chart event
hadm_id INT NOT NULL, – Hospital
admission id of chart event
stay_id INT NOT NULL, – Stay ID of
chart event
itemid INT NOT NULL, – Item ID of chart
event
charttime TIMESTAMP(0) NOT NULL, –
Chart time of chart event
valuenum DOUBLE PRECISION, – Nu-
merical value measured of chart event
valueuom VARCHAR(50), – Unit of
numerical value of chart event
);
CREATE TABLE inputevents

(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of input event
subject_id INT NOT NULL, – Subject id
of input event
hadm_id INT NOT NULL, – Hospital
admission id of input event
stay_id INT NOT NULL, – Stay id of input
event
starttime TIMESTAMP(0) NOT NULL, –
Start time of input event
itemid INT NOT NULL, – Item id of input
event
amount DOUBLE PRECISION, – Amount
of input event
);
CREATE TABLE outputevents
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of output event
subject_id INT NOT NULL, – Subject id
of output event
hadm_id INT NOT NULL, – Hospital
admission id of output event
stay_id INT NOT NULL, – Stay id of
output event
charttime TIMESTAMP(0) NOT NULL, –
Chart time of output event
itemid INT NOT NULL, – Item id of output
event
value DOUBLE PRECISION, – Value of
output event
);
CREATE TABLE microbiologyevents
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of microbiologyevent
subject_id INT NOT NULL, – Subject id
of microbiologyevent
hadm_id INT NOT NULL, – Hospital
admission id of microbiologyevent
charttime TIMESTAMP(0) NOT NULL, –
Chart time of microbiologyevent
spec_type_desc VARCHAR(100), – Speci-
men name of microbiologyevent
test_name VARCHAR(100), – Test name
of microbiologyevent
org_name VARCHAR(100), – Organism
name of microbiologyevent
);

CREATE TABLE icustays
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of icu stay
subject_id INT NOT NULL, – Subject id
of icu stay
hadm_id INT NOT NULL, – Hospital
admission id of icu stay
stay_id INT NOT NULL UNIQUE, – Stay
id of icu stay
first_careunit VARCHAR(20) NOT NULL,
– first care unit of icu stay
last_careunit VARCHAR(20) NOT NULL,
– Last care unit of icu stay
intime TIMESTAMP(0) NOT NULL, – In
time of icu stay
outtime TIMESTAMP(0), – Out time of icu
stay
);
CREATE TABLE transfers
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of transfer
subject_id INT NOT NULL, – Subject Id
of transfer
hadm_id INT NOT NULL, – Hospital
admission id of transfer
transfer_id INT NOT NULL, – Transfer Id
of transfer
eventtype VARCHAR(20) NOT NULL, –
Event type of transfer
careunit VARCHAR(20), – Care unit of
transfer
intime TIMESTAMP(0) NOT NULL, – In
time of transfer
outtime TIMESTAMP(0), – Out time of
transfer
);
– admissions.subject_id can be joined with
patients.subject_id
– diagnoses_icd.hadm_id can be joined with
admissions.hadm_id
– diagnoses_icd.icd_code can be joined with
d_icd_diagnoses.icd_code
– procedures_icd.hadm_id can be joined
with admissions.hadm_id
– procedures_icd.icd_code can be joined
with d_icd_procedures.icd_code
– labevents.hadm_id can be joined with
admissions.hadm_id

– labevents.itemid can be joined with
d_labitems.itemid
– prescriptions.hadm_id can be joined with
admissions.hadm_id
– cost.hadm_id can be joined with admissions.hadm_id
– cost.event_id can be joined with diagnoses_icd.row_id
– cost.event_id can be joined with procedures_icd.row_id
– cost.event_id can be joined with labevents.row_id
– cost.event_id can be joined with prescriptions.row_id
– chartevents.hadm_id can be joined with
admissions.hadm_id
– chartevents.stay_id can be joined with
icustays.stay_id
– chartevents.itemid can be joined with
d_items.itemid
– inputevents.hadm_id can be joined with
admissions.hadm_id
– inputevents.stay_id can be joined with
icustays.stay_id
– inputevents.itemid can be joined with
d_items.itemid
– outputevents.hadm_id can be joined with
admissions.hadm_id
– outputevents.stay_id can be joined with
icustays.stay_id
– outputevents.itemid can be joined with
d_items.itemid
– microbiologyevents.hadm_id can be
joined with admissions.hadm_id
– icustays.hadm_id can be joined with
admissions.hadm_id
– transfers.hadm_id can be joined with
admissions.hadm_id

### Answer
Given the database schema, here is the SQL
query that answers [QUESTION]What was
the drug that patient 10015931 was prescribed with within the same hospital visit
after the replacement of aortic valve with
zooplastic tissue, percutaneous approach
since 5 months ago?[/QUESTION]
[SQL] SELECT admissions.subject_id,
prescriptions.drug,prescriptions.starttime,
admissions.hadm_id FROM prescrip-

666

tions JOIN admissions ON prescriptions.hadm_id = admissions.hadm_id WHERE admissions.subject_id = 10015931 [/SQL]

# E   Prompt and Examples for CAS Module

Following is the full example for the prompt in the figure 2

### Example 2

Based on the database schema and table description, determine which AI assistant's answer accurately identifies whether the given question can generate an SQL query or not.
### Database Table Description
The table name and its corresponding description are as follows:
ADMISSIONS – Every unique hospitalization for each patient in the database
PATIENTS – Every unique patient in the database
D_ICD_DIAGNOSES – International Statistical Classification of Diseases and Related Health Problems (ICD-9) codes relating to diagnoses
D_ICD_PROCEDURES – International Statistical Classification of Diseases and Related Health Problems (ICD-9) codes relating to procedures
D_LABITEMS – Local codes ('ITEMIDs') appearing in the database that relate to laboratory tests
D_ITEMS – Local codes ('ITEMIDs') appearing in the database, except those that relate to laboratory tests
DIAGNOSES_ICD – Hospital assigned diagnoses, coded using the International Statistical Classification of Diseases and Related Health Problems (ICD) system
PROCEDURES_ICD – Patient procedures, coded using the International Statistical Classification of Diseases and Related Health Problems (ICD) system
LABEVENTS – Laboratory measurements for patients both within the hospital and in outpatient clinics
PRESCRIPTIONS – Medications ordered for a given patient

COST – All patients events cost
CHARTEVENTS – All charted observations for patients
INPUTEVENTS – Intake for patients monitored while in the ICU
OUTPUTEVENTS – Output information for patients while in the ICU
MICROBIOLOGYEVENTS – Microbiology culture results and antibiotic sensitivities from the hospital database
ICUSTAYS – Every unique ICU stay in the database
TRANSFERS – Patient movement from bed to bed within the hospital

Database Schema
This query will run on a database whose schema is represented in this string:
CREATE TABLE patients
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of the patient
subject_id INT NOT NULL UNIQUE, – Unique subject id of the patient
gender VARCHAR(5) NOT NULL, – Gender of the patient
dob TIMESTAMP(0) NOT NULL, – Date of birth of the patient
dod TIMESTAMP(0) – Date of death of the patient
);
CREATE TABLE admissions
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of the admission
subject_id INT NOT NULL, – Subject id of the admission
hadm_id INT NOT NULL UNIQUE, – Unique hospital admission id of the admission
admittime TIMESTAMP(0) NOT NULL, – Admit time of the admission
dischtime TIMESTAMP(0), – Discharge time of the admission
admission_type VARCHAR(50) NOT NULL, – Admission type of the admission
admission_location VARCHAR(50) NOT NULL, – Admission location of the admission
discharge_location VARCHAR(50), –

Discharge location of the admission
insurance VARCHAR(255) NOT NULL, –
Insurance of the admission
language VARCHAR(10), – Langauge of
the admission
marital_status VARCHAR(50), – Marital
status of the admission
age INT NOT NULL, – Age of the
admission
);
CREATE TABLE d_icd_diagnoses
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of the icd diagnose
icd_code VARCHAR(10) NOT NULL
UNIQUE, – Unique icd code of the icd
diagnose
long_title VARCHAR(255) NOT NULL –
Title of the icd diagnose
);
CREATE TABLE d_icd_procedures
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of icd procedure
icd_code VARCHAR(10) NOT NULL
UNIQUE, – Unique icd code of the icd
procedure
long_title VARCHAR(255) NOT NULL –
Title of the icd procedure
);
CREATE TABLE d_labitems
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of the item relate to laboratory
tests
itemid INT NOT NULL UNIQUE, –
Unique item id of the item relate to
laboratory tests
label VARCHAR(200) – Label of the item
relate to laboratory tests
);
CREATE TABLE d_items
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of the item excepts item relate
to laboratory tests
itemid INT NOT NULL UNIQUE, –
Unique item id of the item excepts item
relate to laboratory tests
label VARCHAR(200) NOT NULL, –

Label of item excepts item relate to
laboratory tests
abbreviation VARCHAR(200) NOT NULL,
– Abbreviation of item excepts item relate
to laboratory tests
linksto VARCHAR(50) NOT NULL –
Event linked to item excepts item relate to
laboratory tests
);
CREATE TABLE diagnoses_icd
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of diagnose
subject_id INT NOT NULL, – Subject id
of diagnose
hadm_id INT NOT NULL, – Hospital
admission id of diagnose
icd_code VARCHAR(10) NOT NULL, –
ICD code of diagnose
charttime TIMESTAMP(0) NOT NULL, –
Chart time of diagnose
);
CREATE TABLE procedures_icd
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of procedures
subject_id INT NOT NULL, – Subject id
of procedures
hadm_id INT NOT NULL, – Hospital
admission id of procedures
icd_code VARCHAR(10) NOT NULL, –
ICD code of procedures
charttime TIMESTAMP(0) NOT NULL, –
Chart time of procedures
);
CREATE TABLE labevents
(
row_id INT NOT NULL PRIMARY KEY, –
Unique ID of laboratory event
subject_id INT NOT NULL, – Subject id
of laboratory event
hadm_id INT NOT NULL, – Hospital
admission id of laboratory event
itemid INT NOT NULL, – Item id of
laboratory event
charttime TIMESTAMP(0), – Chart time of
laboratory event
valuenum DOUBLE PRECISION, – Nu-
merical value measured of laboratory event
valueuom VARCHAR(20), – Unit of

numerical value of laboratory event
);
CREATE TABLE prescriptions
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of prescription
subject_id INT NOT NULL, – Subject id of prescription
hadm_id INT NOT NULL, – Hospital admission id of prescription
starttime TIMESTAMP(0) NOT NULL, – Start time of prescription
stoptime TIMESTAMP(0), – Stop time of prescription
drug VARCHAR(255) NOT NULL, – Drug name of prescription
dose_val_rx VARCHAR(100) NOT NULL, – Dosage value of prescription
dose_unit_rx VARCHAR(50) NOT NULL, – Dosage unit of prescription
route VARCHAR(50) NOT NULL, – Intake method of prescription
);
CREATE TABLE cost
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of cost event
subject_id INT NOT NULL, – Subject id of cost event
hadm_id INT NOT NULL, – Hospital admission id of cost event
event_type VARCHAR(20) NOT NULL, – Event type of cost event
event_id INT NOT NULL, – Event id of cost event
chargetime TIMESTAMP(0) NOT NULL, – Charge time of cost event
cost DOUBLE PRECISION NOT NULL, – Cost of cost event
);
CREATE TABLE chartevents
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of chart event
subject_id INT NOT NULL, – Subject id of chart event
hadm_id INT NOT NULL, – Hospital admission id of chart event
stay_id INT NOT NULL, – Stay ID of chart event

itemid INT NOT NULL, – Item ID of chart event
charttime TIMESTAMP(0) NOT NULL, – Chart time of chart event
valuenum DOUBLE PRECISION, – Numerical value measured of chart event
valueuom VARCHAR(50), – Unit of numerical value of chart event
);
CREATE TABLE inputevents
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of input event
subject_id INT NOT NULL, – Subject id of input event
hadm_id INT NOT NULL, – Hospital admission id of input event
stay_id INT NOT NULL, – Stay id of input event
starttime TIMESTAMP(0) NOT NULL, – Start time of input event
itemid INT NOT NULL, – Item id of input event
amount DOUBLE PRECISION, – Amount of input event
);
CREATE TABLE outputevents
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of output event
subject_id INT NOT NULL, – Subject id of output event
hadm_id INT NOT NULL, – Hospital admission id of output event
stay_id INT NOT NULL, – Stay id of output event
charttime TIMESTAMP(0) NOT NULL, – Chart time of output event
itemid INT NOT NULL, – Item id of output event
value DOUBLE PRECISION, – Value of output event
);
CREATE TABLE microbiologyevents
(
row_id INT NOT NULL PRIMARY KEY, – Unique ID of microbiologyevent
subject_id INT NOT NULL, – Subject id of microbiologyevent
hadm_id INT NOT NULL, – Hospital

admission id of microbiologyevent

charttime TIMESTAMP(0) NOT NULL, – Chart time of microbiologyevent

spec_type_desc VARCHAR(100), – Specimen name of microbiologyevent

test_name VARCHAR(100), – Test name of microbiologyevent

org_name VARCHAR(100), – Organism name of microbiologyevent

);

CREATE TABLE icustays

(

row_id INT NOT NULL PRIMARY KEY, – Unique ID of icu stay

subject_id INT NOT NULL, – Subject id of icu stay

hadm_id INT NOT NULL, – Hospital admission id of icu stay

stay_id INT NOT NULL UNIQUE, – Stay id of icu stay

first_careunit VARCHAR(20) NOT NULL, – first care unit of icu stay

last_careunit VARCHAR(20) NOT NULL, – Last care unit of icu stay

intime TIMESTAMP(0) NOT NULL, – In time of icu stay

outtime TIMESTAMP(0), – Out time of icu stay

);

CREATE TABLE transfers

(

row_id INT NOT NULL PRIMARY KEY, – Unique ID of transfer

subject_id INT NOT NULL, – Subject Id of transfer

hadm_id INT NOT NULL, – Hospital admission id of transfer

transfer_id INT NOT NULL, – Transfer Id of transfer

eventtype VARCHAR(20) NOT NULL, – Event type of transfer

careunit VARCHAR(20), – Care unit of transfer

intime TIMESTAMP(0) NOT NULL, – In time of transfer

outtime TIMESTAMP(0), – Out time of transfer

);

– admissions.subject_id can be joined with patients.subject_id

– diagnoses_icd.hadm_id can be joined with admissions.hadm_id

– diagnoses_icd.icd_code can be joined with d_icd_diagnoses.icd_code

– procedures_icd.hadm_id can be joined with admissions.hadm_id

– procedures_icd.icd_code can be joined with d_icd_procedures.icd_code

– labevents.hadm_id can be joined with admissions.hadm_id

– labevents.itemid can be joined with d_labitems.itemid

– prescriptions.hadm_id can be joined with admissions.hadm_id

– cost.hadm_id can be joined with admissions.hadm_id

– cost.event_id can be joined with diagnoses_icd.row_id

– cost.event_id can be joined with procedures_icd.row_id

– cost.event_id can be joined with labevents.row_id

– cost.event_id can be joined with prescriptions.row_id

– chartevents.hadm_id can be joined with admissions.hadm_id

– chartevents.stay_id can be joined with icustays.stay_id

– chartevents.itemid can be joined with d_items.itemid

– inputevents.hadm_id can be joined with admissions.hadm_id

– inputevents.stay_id can be joined with icustays.stay_id

– inputevents.itemid can be joined with d_items.itemid

– outputevents.hadm_id can be joined with admissions.hadm_id

– outputevents.stay_id can be joined with icustays.stay_id

– outputevents.itemid can be joined with d_items.itemid

– microbiologyevents.hadm_id can be joined with admissions.hadm_id

– icustays.hadm_id can be joined with admissions.hadm_id

– transfers.hadm_id can be joined with admissions.hadm_id

Question: 'What was the drug that patient 10015931 was prescribed with within the same hospital visit after the replacement of aortic valve with zooplastic tissue, percutaneous approach since 5 months ago?'
Answer: Able to generate SQL Query.

Question: 'Tell me the name of the prescription drug that patient 10015931 was prescribed in the same day after having a replacement of aortic valve with zooplastic tissue, percutaneous approach procedure since 4 months ago?'
Answer: Able to generate SQL Query.

Question: 'What was prescribed to patient 10015931 during the same hospital visit following their replacement of aortic valve with zooplastic tissue, percutaneous approach during this month?'
Answer: Able to generate SQL Query.

Question: 'What was the drug that patient 10025463 was prescribed for during the same hospital encounter after the procedure of excision or destruction of other lesion or tissue of heart, endovascular approach?'
Answer: Able to generate SQL Query.

Question: "What was the drug that patient 10015931 was prescribed with within the same hospital visit after the replacement of aortic valve with zooplastic tissue, percutaneous approach since 5 months ago?"
Ai Assitant 1's Answer: Able to generate SQL Query.
Ai Assitant 2's Answer: Able to generate SQL Query.
Answer: Let's think step by step."

# KU-DMIS at EHRSQL 2024:
# Generating SQL query via question templatization in EHR

**Hajung Kim[1*], Chanhwi Kim[1*], Hoonick Lee[1], Kyochul Jang[1], Jiwoo Lee[1],**
**Kyungjae Lee[2], Gangwoo Kim[1], Jaewoo Kang[1,3†]**

[1]Korea University, [2]LG AI Research, [3]AIGEN Sciences
{hajungk, chanhwi_kim, hoonick, gcj0125, hijiwoo7}@korea.ac.kr
kyungjae.lee@lgresearch.ai, {gangwoo_kim, kangj}@korea.ac.kr

## Abstract

Transforming natural language questions into SQL queries is crucial for precise data retrieval from electronic health record (EHR) databases. A significant challenge in this process is detecting and rejecting unanswerable questions that request information beyond the database's scope or exceed the system's capabilities. In this paper, we introduce a novel text-to-SQL framework that robustly handles out-of-domain questions and verifies the generated queries with query execution. Our framework begins by standardizing the structure of questions into a templated format. We use a powerful large language model (LLM), fine-tuned GPT-3.5 with detailed prompts involving the table schemas of the EHR database system. Our experimental results demonstrate the effectiveness of our framework on the EHRSQL-2024 benchmark benchmark, a shared task in the ClinicalNLP workshop. Although a straightforward fine-tuning of GPT shows promising results on the development set, it struggled with the out-of-domain questions in the test set. With our framework, we improve our system's adaptability and achieve competitive performances in the official leaderboard of the EHRSQL-2024 challenge.

## 1 Introduction

Electronic Health Records (EHRs) are crucial elements of the contemporary healthcare system, storing patients' medical histories in relational databases. However, retrieving data from EHRs can be challenging, requiring specialized training in Structured Query Language (SQL). To bridge this gap, the previous works build AI-powered systems that parse the user's question (Yin et al., 2020; Brunner and Stockinger, 2021) or convert it into an SQL query that the database can process. Lee et al.



Figure 1: In the proposed Text-to-SQL framework, when a query is presented in natural language, the model generates SQL code to retrieve the required information from the database. If the query requires information absent from the database, the Text-to-SQL model returns a 'null' response.

(2022) identify an essential component in this text-to-SQL task; recognizing and adequately handling unanswerable questions that seek information beyond what the database contains. Hence, to ensure reliability and trustworthiness, the systems should be able to refrain from answering unanswerable questions.

To further encourage research in this field, the Clinical NLP 2024 workshop has introduced a new shared task called EHRSQL-2024 (Lee et al., 2024) to motivate the development of more reliable question-answering (QA) systems. The EHRSQL-2024 dataset involves the real-world needs of medical personnel, incorporating templates of their most

---

* Equal contribution, † Corresponding author

common questions. In this challenge, systems are tasked to generate SQL queries that accurately return the desired information from tables from the MIMIC-IV (Johnson et al., 2016a). Additionally, the dataset includes inherently unanswerable questions, either due to the restrictions of the database schema or the request for information not contained within the databases. On the other hand, the test set presents distracting question types that contain noisy words, further testing the robustness of participants' systems.

In this paper, we introduce a novel framework created to convert natural language questions to SQL queries for EHR databases. This framework transforms free-form questions into a templated format to handle distracting questions. We fine-tune GPT-3.5-turbo (Brown et al., 2020), one of the most performant large language models (LLMs), optimizing it to effectively interpret intricate medical queries and produce the corresponding SQL queries. We also provide detailed prompts that describe the tables in the EHR database. For SQL generation, given the task's complexity and the relationships between tables, we break it down into two steps: selecting relevant tables and then generating SQL by reflecting in-depth on the selected tables. We enhance the accuracy and reliability of the generated SQL queries by correcting any errors in table names and applying ensemble techniques with majority voting.

Our empirical results of fine-tuned GPT-3.5 on the EHRSQL-2024 benchmark highlight its capability, achieving third place on the development set. However, it revealed a limitation in generalizing to questions in the test set that diverged from the predefined templates. By using our framework, we successfully address this gap between free-form questions, resulting in a notable improvement of 26.5 in the RS(10) metric in the test set. Additionally, we find that decomposing the task into two steps contributed to this success, with a significant improvement in RS (10) in the test set. Furthermore, by employing further verification and ensemble techniques, we attain fourth place in the EHRSQL-2024 challenge's official leaderboard.

We conduct in-depth analyses of the questions to uncover disparities in each split. In particular, we apply $N$-gram counting of the questions to highlight the distribution gaps. This variation emphasizes the need to develop a resilient model capable of adapting to and performing consistently across datasets with diverse word distributions. Ad-

ditionally, we manually categorize the unanswerable questions into three distinct types.

## 2 Related Works

### 2.1 Text-to-SQL Generation

Text-to-SQL conversion requires interpreting natural language questions, matching them with the database schema, and producing accurate SQL queries that reflect the question's intent. This task is particularly challenging for individuals unfamiliar with database structures, highlighting the need for methods that translate natural language into SQL queries—a focus of ongoing research due to real-world applications. However, accurately generating SQL code from natural language is complex, mainly because of the challenges in integrating precise database knowledge into the model (Qin et al., 2022; Katsogiannis-Meimarakis and Koutrika, 2023).

Initially, efforts to address Text-to-SQL employed predefined rules (Sen et al., 2020) to handle existing difficulties. The field has evolved since then to explore encoder-decoder models (Cai et al., 2017; Popescu et al., 2022), and Text-to-SQL is tested on sequence-to-sequence approaches (Qi et al., 2022). With the rapid advancement in deep learning research, methodologies incorporating graph representation (Xu et al., 2018; Wang et al., 2019; Brock et al., 2021) and attention mechanisms (Liu et al., 2023b) have been extensively applied to Text-to-SQL tasks. Additionally, the Text-to-SQL task, tailored to real-world data, has been conducted on datasets such as WikiSQL (Zhong et al., 2017), Spider (Yu et al., 2018), KaggleD-BQA (Lee et al., 2021), and BIRD (Li et al., 2023).

With the emergence of LLMs like GPT (Brown et al., 2020) and Llama (Touvron et al., 2023), research leveraging these models has proliferated. Their comprehensive pretraining on massive text corpora enables them to show promising results using techniques like prompt engineering and in-context learning (Trummer, 2022; Liu et al., 2023a; Chang and Fosler-Lussier, 2023; Dong et al., 2023; Sun et al., 2023). Despite these advancements, exploring supervised fine-tuning has led to even greater enhancements in their performance (Gao et al., 2023).

### 2.2 Text-to-SQL in EHR database

The MIMIC-III (Johnson et al., 2016b) is a prominent EHR database in the healthcare domain. MIM-

ICSQL (Tarbell et al., 2023) is the first dataset constructed based on the MIMIC-III database, designing questions generated from pre-formatted templates. Similarly, emrKBQA (Raghavan et al., 2021) derived from the MIMIC-III database and the emrQA (Yue et al., 2020) dataset focused on clinical reading comprehension expands the research scope. EHRSQL, introduced by Lee et al. (2022), is an extensive text-to-SQL dataset that is linked to the two open-source EHR databases, MIMIC-III and eICU (Pollard et al., 2018). Created based on feedback from 222 professionals with varied experience levels, EHRSQL covers a wide range of real-world scenarios. This dataset includes time-sensitive questions to highlight the critical importance of time in the healthcare domain. Additionally, it incorporates unanswerable questions to evaluate the system's capability to recognize and handle such inquiries effectively.

## 2.3 Discriminating Unanswerable Questions

The distinction between answerable and unanswerable questions is crucial in NLP tasks, especially in domains where accuracy and reliability are critical, such as healthcare. Discriminating between these types of questions is complex due to the subtle differences in what a question may require for a satisfactory answer. The language models often exhibit overconfidence in their ability to accurately respond to a given question. To address this, the specialized datasets have been designed through various methodologies, such as rule-based editing (Jia and Liang, 2017), distant supervision (Joshi et al., 2017), and crowdsourcing (Rajpurkar et al., 2018), each method offering its own set of benefits and challenges for identifying unanswerable questions. This advancement facilitates more reliable and accurate question-answering capabilities, which is crucial for applications where the cost of misinformation can be high.

## 3 Methods

Figure 2 presents an outline of our proposed methodology. Our process starts with the templatization of questions, transforming free-form inquiries into a standardized format to ensure consistency in how queries are represented. Additionally, we enrich the model's understanding by supplying detailed information about the database tables, thereby improving its capacity to formulate precise queries. To further elevate the accuracy of the

generated SQL queries, we introduce a verification phase to confirm that the queries accurately correspond to the intended data retrieval objectives.

## 3.1 Question Templatization

We introduce question templatization to handle the diverse forms of question presentation. This approach addresses the challenge of questions deviating from a standard template by employing a reverse engineering strategy. By converting free-form questions into a templated format, we aim to align them more closely with similar patterns, thus bridging the gap between the varied question formats in real-world contexts. Specifically, we prompt GPT-4-turbo to rewrite questions to match the structure of pre-defined templates more closely.

Identifying semantically close questions involves searching for questions similar to the input question. This similarity is quantified by calculating the Euclidean distance between the question embeddings and comparing input questions to potential neighbors. We mask identification information to ensure that specific table values do not skew this comparison. For example, a question like *"Count how many times in the first hospital visit patient 10004457 had coronary arteriography using two catheters."* is transformed into *"Count how many times patient <patient number> had <procedure> during their first hospital visit ."* By adopting this method, we achieve a uniform question format, effectively standardizing free-form queries and reducing discrepancies in dataset distribution. The templatized question is utilized as the input question.

## 3.2 SQL generation

Considering the complexity of the text-to-SQL task and the intricate relationships among more than ten tables in the database, we propose a two-stage approach that involves a table selection phase followed by a self-reflection phase.

**Table Selection** We task the GPT model with converting natural language questions into SQL queries. The construction of prompts for the model involves three essential components: (1) outlining the text-to-SQL task guides the model to convert a natural language question into an SQL query for data retrieval from the EHR database. We clarify
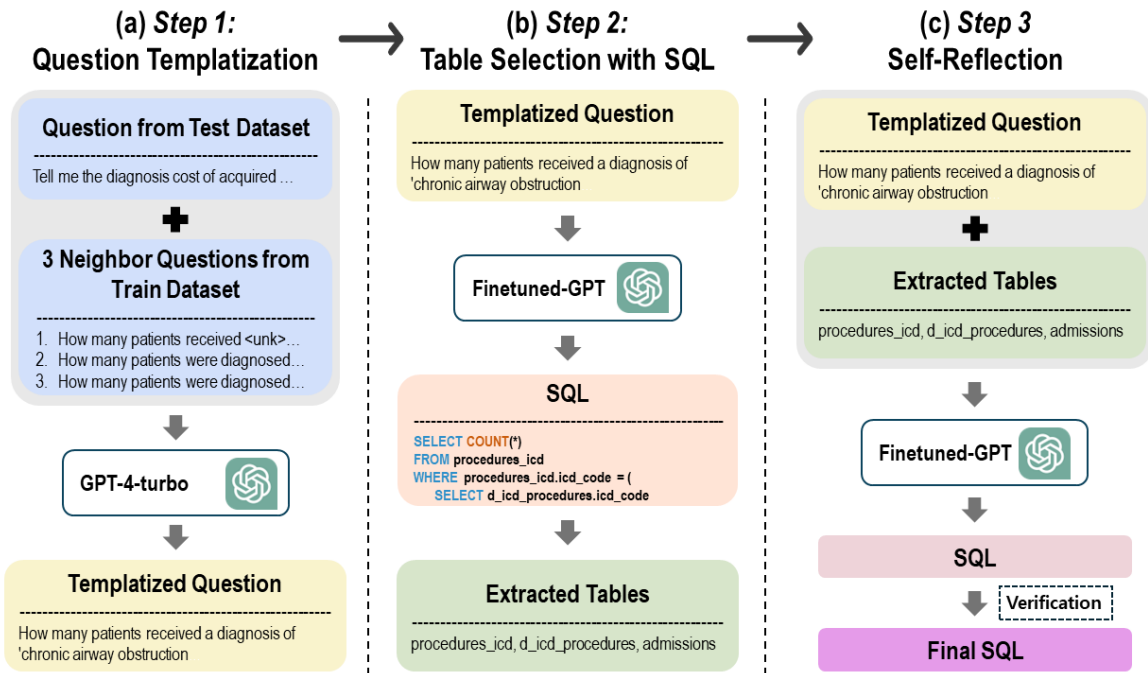
Figure 2: Overview of our framework. (a) Question Templatization (Sec. 3.1). Implementing question templatization to convert free-form questions into a structured format. (b) SQL Generation (Sec. 3.2). Providing task outlines and table information to aid in precise query generation. (c) Self-Reflection and Verification (Sec. 3.2, 3.3). Providing detailed table information identified in the initial SQL generation and then finalizing the process.

that the database uses SQLite and highlight the syntactical nuances between SQLite and other SQL dialects to guide the model's syntax choice. (2) By detailing the database tables, we describe the database's complex structure, listing over ten tables with brief descriptions and their respective columns. This detail is crucial since it aids the model in identifying the relevant tables and navigating their relational schema without direct access to the database values. We follow the format introduced in DAIL-SQL (Gao et al., 2023) for table schema details, which allows both natural language and SQL representations. (3) Presenting the question for conversion is the natural language question to be transformed. By using this prompt, we prompt the model to produce an SQL query that matches the question and subsequently identifies the table name mentioned within the SQL query.

**Self-Reflection** The prompt for the self-reflection stage is similar to Table Selection, except for detailing table information. In this stage, the prompt is augmented with detailed descriptions for each table column identified in the initial SQL query. This refinement aims to enhance the SQL query formulation by providing a more comprehensive understanding of the selected

table's specifics, enabling the model to generate a more accurate and targeted SQL query.

## 3.3 SQL Verification

We implement a verification step on the generated SQL queries to address two specific scenarios. In the first scenario, some questions can be technically converted into SQL queries but remain unanswerable due to the absence of required information in the dataset. To avoid providing incorrect SQL results and improper answers, which are unanswerable, we verify the validity of each SQL query by executing it against the database. If the execution results in an error, indicating the SQL query cannot retrieve the correct answer, we replace the SQL query with *null* instead. This adjustment ensures the query is considered valid but unanswerable, optimizing score outcomes.

The second scenario addresses instances where the generated SQL query includes incorrect references to table names or column names. In such cases, we identify the correct table name and associated column names based on the table values mentioned in the SQL query. We then modify the SQL query to accurately reflect the proper table name and column names to which the

| | Development | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| Team | RS(0) | RS(5) | RS(10) | RS(N) | RS(0) | RS(5) | RS(10) | RS(N) |
| LG AI Research & KAIST | 90.37 | **89.51** | **88.65** | **-109.6** | **88.17** | **84.75** | **81.32** | **-711.83** |
| PromptMind | 66.38 | 59.5 | 52.62 | $-1533.62$ | 82.6 | 78.75 | 74.89 | $-817.4$ |
| ProbGate | 84.18 | 79.45 | 74.72 | $-1015.82$ | 81.92 | 78.06 | 74.21 | $-818.08$ |
| **KU-DMIS (Ours)** | **91.57** | 82.98 | 74.38 | $-1908.43$ | 72.07 | 65.64 | 59.21 | $-1427.93$ |
| oleg1996 | 47.03 | 34.14 | 21.24 | $-2952.97$ | 68.89 | 56.47 | 44.04 | $-2831.11$ |
| LTRC-IIITH | N/A | N/A | N/A | N/A | 66.84 | 55.27 | 43.7 | $-2633.16$ |
| Saama Technologies | 57.78 | 50.47 | 43.16 | $-1642.22$ | 53.21 | 44.64 | 36.08 | $-1946.79$ |
| TEAM_optimist | N/A | N/A | N/A | N/A | 14.14 | $-349.61$ | $-713.37$ | $-84885.86$ |

Table 1: Official results of the leaderboard on EHRSQL-2024 dataset. The teams are ranked based on Reliability Score RS(10).

### 3.4 Ensemble with Majority Voting

We incorporate an ensemble method to determine the final SQL query. We first instruct GPT-4-turbo to evaluate whether the generated SQL query accurately captures the intent of the original natural language question. This alignment check ensures that the model prioritizes the core intent of the query, such as using the 'COUNT' function in SQL queries asking for a count of patients. To finalize the SQL query or its resulting answer from the database execution, we adopt a majority voting system. This ensemble strategy mitigates the variability inherent in the fine-tuned model and improves the robustness. Using majority voting to select the SQL query or derive its answer aims to improve performance metrics by effectively managing *null* responses.

## 4 Experiments

### 4.1 Experimental Setup

**Dataset** We evaluate our frameworkusing the EHRSQL-2024 challenge benchmark dataset (Lee et al., 2024). This large-scale Text-to-SQL dataset contains 5,124 instances in the train set, 1,163 instances in the development set, and 1,167 instances in the test set, spanning 17 tables. The train dataset comprises natural language questions paired with their corresponding SQL queries. However, the SQL queries associated with the questions in the development and test sets are not provided.

**Metric** Following Lee et al. (2024), we use the Reliability Score (RS). RS is unique because it rewards correct SQL queries for answerable questions ($Q_{ans}$) and the choice to abstain from answering unanswerable questions ($Q_{una}$). At the same time, it penalizes incorrect SQL generation for $Q_{ans}$ and any attempt to answer $Q_{una}$. Moreover, RS includes a penalty factor 'c' to adjust the evaluation's strictness according to specific safety requirements. The corresponding formula is as follows.

$$\phi_c(x) = \begin{cases} 1 & \text{if } x \in Q_{\text{ans}}, g(x) = 1, \text{Acc}(x) = 1 \\ 0 & \text{if } x \in Q_{\text{ans}}, g(x) = 0 \\ -c & \text{if } x \in Q_{\text{ans}}, g(x) = 1, \text{Acc}(x) = 0 \\ -c & \text{if } x \in Q_{\text{una}}, g(x) = 1 \\ 1 & \text{if } x \in Q_{\text{una}}, g(x) = 0. \end{cases}$$

The adaptability of RS is demonstrated by evaluating models under four different scenarios, which vary based on the severity of the penalty term: RS(0), RS(10), and RS(N). For this challenge, the primary metric is RS(10), emphasizing the importance of accurately assessing answerable questions and refraining from generating SQL for unanswerable questions.

### 4.2 Implementation Details

We utilize GPT, one of the most performant Large Language Models (LLMs), to enhance the translation from text to SQL. We investigate the effectiveness of in-context learning and supervised fine-tuning methods. We employ GPT-3.5-turbo, GPT-4-turbo, and GPT-4 models for in-context learning, augmenting the prompt with three more examples. These examples consist of pairs of semantically related questions, including the input question and their corresponding SQL queries. The semantic relatedness is determined

| Model | Few-shot | Table info. | RS(0) | RS(10) | RS(N) |
|---|---|---|---|---|---|
| GPT-3.5-turbo | 0 | O | 29.53 | -250.19 | -28670.47 |
| | 3 | O | 48.54 | 15.40 | -3351.46 |
| | 3 | - | 75.34 | 29.53 | -4624.66 |
| GPT-4-turbo | 0 | O | 36.74 | -196.20 | -23863.26 |
| | 3 | O | 67.84 | -18.91 | -8832.16 |
| | 3 | - | 85.87 | 42.98 | -4314.13 |
| GPT-4 | 0 | O | 38.40 | -215.98 | -26061.60 |
| | 3 | O | 79.82 | 15.50 | -6520.18 |
| | 3 | - | 90.45 | 62.18 | -2809.55 |
| Finetuned-GPT | - | O | **98.05** | **91.23** | **-601.95** |

Table 2: Training set performance. Comparison of GPT models.

| Model | RS(0) | RS(10) | RS(N) |
|---|---|---|---|
| GPT-3.5-turbo | 70.34 | 13.59 | -6529.66 |
| GPT-4-turbo | 76.53 | -6.02 | -9523.47 |
| GPT-4 | 79.28 | -17.88 | -11220.72 |
| Finetuned-GPT | **93.12** | 50.99 | -4806.88 |
|    w/ table info. in SQL form | 83.23 | 17.02 | -7616.77 |
|    w/ Self-Reflection | 83.15 | 62.51 | -2316.85 |
|    w/ Ensemble | 91.57 | **74.38** | **-1908.43** |

Table 3: Ablation study conducted on the development set showcases the performance of in-context learning with few examples using GPT-3.5-turbo, GPT-4-turbo, and GPT-4, alongside fine-tuning performed with GPT-3.5-turbo using various additional techniques.

by calculating the Euclidean distance between question embeddings derived from the training dataset and the input question embedding. For supervised fine-tuning, we focus on the GPT-3.5-turbo model, the primary model available for fine-tuning. The model is prompted without including neighboring examples. Based on the evaluation results, it is clear that the supervised fine-tuning methodology is particularly effective in addressing the challenges inherent in text-to-SQL tasks. Further details are provided in section 4.4. Consequently, the fine-tuned GPT-3.5-turbo model is selected for further detailed experiments.

### 4.3 Leaderboard Results

Table 1 presents the scores of the participants' systems, ranked according to the RS(10) score. We secure the fourth position in the test set rankings. All participating teams utilized Large Language Models (LLMs), with the top four teams, including ours, primarily employing a fine-tuned GPT model and incorporating various other techniques. This table underscores the efficacy of LLMs in addressing Text-to-SQL tasks.

### 4.4 In-Context Learning and Fine-tuning

To evaluate the effectiveness of various GPT models for Text-to-SQL tasks, we conduct experiments with GPT-3.5-turbo, GPT-4-turbo, and GPT-4 for in-context learning and a fine-tuned version of GPT-3.5-turbo for supervised fine-tuning. Due to submission limitations, we assessed the GPT models using the training set. We adopt a k-fold cross-validation method with $k = 5$, training on four folds and evaluating the remaining fold. To maintain the balance of answerable and unanswerable questions in the training dataset, we divide unanswerable questions into three categories. When partitioning the training dataset into five folds, we ensured that the proportions of these categories were reflected in each fold. A detailed analysis of these categorized groups is discussed in section 5.2.

Table 2 presents the comparison results of the GPT models. We experimented with variations by providing few-shot examples and including table information. The fine-tuned GPT model demonstrates superior performance across all metrics, making it our model of choice. Interestingly, the inclusion of table information slightly reduces performance in all in-context learning scenarios. We speculate that the table information in our experiment, which merely lists table names and column names, lacks detailed relational data like primary and foreign keys. Consequently, this minimal and potentially uninformative text might have acted as a distraction.

### 4.5 Table Information Format

The prompt includes table information to accurately identify the table and column names. Following the DAIL-SQL approach (Gao et al., 2023), we explore different formats of presenting table information, in both natural language and SQL format, within the same prompt framework. Our experiments, detailed in table 3, reveal that presenting table information in SQL format results in a decrease in the RS (10) score from 50.99 to 17.02. This suggests that natural language formats are more readily interpretable by the language model such as GPT.

### 4.6 Table Selection Results

Considering the complexity of the more than ten tables and the resulting SQL queries that reference multiple tables, we hypothesize that a self-

| Model | Inclusion | Jaccard | Exact Match |
|---|---|---|---|
| GPT-3.5-turbo | 0.7933 | 0.7930 | 0.7836 |
| GPT-4-turbo | 0.8912 | 0.8908 | 0.8723 |
| GPT-4 | 0.9250 | 0.9244 | 0.9123 |
| **Finetuned-GPT** | **0.9857** | **0.9855** | **0.9844** |
| Table Selector (GPT-3.5-turbo) | 0.8976 | 0.8488 | 0.7115 |

Table 4: Table selection performance.

| Model | RS(0) | RS(10) | RS(N) |
|---|---|---|---|
| Finetuned-GPT (Ensemble) | 72.07 | **59.21** | **-1427.93** |
| Finetuned-GPT (Single) | 78.06 | 13.80 | -7421.94 |
| w/ Self-Reflection | 77.55 | 27.85 | -5722.45 |
| w/ Question Templatization | **80.55** | 40.27 | -4619.45 |

Table 5: Ablation study on the test set. We provide the performance of ensembled and single results. Every component, including SQL regeneration and question templatization, plays a key role in enhancing overall performance.

reflection incorporating selective, detailed table information could enhance the accuracy of the generated SQL queries. In preparation for this self-reflection process, we assess the accuracy of the tables retrieved in the generated SQL queries. This assessment involves calculating the accuracy between correct tables and the tables extracted from the generated SQL queries. We use three metrics as an accuracy score: 1) inclusion score (indicating the presence of the correct tables within the generated SQL), 2) the Jaccard similarity score (comparing the intersection to the union of correct and extracted tables), 3) and the exact match score.

Table 4 suggests that the fine-tuned GPT model effectively identifies the relevant tables without a dedicated table selection model. We extract tables from the initially generated SQL queries and use prompts augmented with detailed information, such as descriptions of each column and examples of values, for the fine-tuned GPT model. The comparison between the initially generated SQL and the outcomes after the self-reflection stage, table 3 shows an increase in the RS(10) score from 50.99 to 62.51 in the development set, and table 5 also illustrates an improvement in the RS(10) score from 13.80 to 27.85 in the test set. This improvement indicates that the regenerated SQL queries provide more reliable and accurate outputs.

### 4.7 Question Templatization

Our analysis focuses on the characteristics of the questions across each dataset. It reveals a decline in the fine-tuned GPT model's scores from the train-

ing set to the development and test sets. This pattern highlights substantial variations among the training, development, and test datasets. To mitigate these discrepancies, we employ the technique that reverses the deviation of questions from templates. We utilize GPT-4-turbo to rephrase the original question. By prompting GPT-4-turbo with the original question and semantically similar questions from the training set and template from (Lee et al., 2022), we aim to achieve consistency with related queries. This approach significantly reduces the distribution gaps between the training and test sets, as demonstrated in Table 5. The improvement in the RS(10) score from 13.80 to 40.27 highlights the effectiveness of question templatization by comparing the performance of a single model before and after its application.

## 5 Analysis

In this section, we analyze the word distribution of questions for each dataset split: training, development, and test sets. The objective is to identify variations in question composition among these datasets. Furthermore, we investigate the distribution of unanswerable questions in the training set to better understand questions that yield an *null* response.

To focus solely on word analysis and minimize noise, we eliminate punctuation marks such as ".", ",", and "?", remove stop words such as "the", "a", and "an" from the questions, and convert all letters to lowercase. After eliminating these elements, we analyze the processed questions using N-grams. The analysis is limited to 1 to 3-grams, which is sufficient for understanding the context of questions while excluding the aforementioned noise. Appendix 9 details the ten most prevalent words alongside their respective frequencies within each dataset arranged in non-increasing order, including the collection of unanswerable queries labeled as *Unanswerable Train set*.

### 5.1 N-gram Distribution

The initial three columns of appendix 9 enumerate the top ten most frequent words in each dataset alongside their respective frequencies. Analysis of appendix 9 indicated that words with high frequency within one dataset tended to be frequent across other datasets as well, suggesting a pattern of similarity. However, it was observed that words with lower frequency, which were not in-

cluded in the table, often did not appear in other datasets. This discrepancy became particularly evident within the context of 3-gram sets, highlighting a distinct distribution among the datasets.

This disparity underscores the necessity of developing a robust model that can adapt and perform well across datasets with different word distributions.

## 5.2 Category of Unanswerable Questions

We analyzed the training set's *null* distribution, identifying 450 unanswerable questions. Our initial qualitative analysis involved categorizing these *null*-labeled questions into three distinct groups through a detailed manual review: (1) Incorrect Patient Number, (2) Require External Knowledge, (3) Out of EHR Database.

In the first case, based on the MIMIC-IV dataset's criteria, a legitimate patient number is identified by its 8-digit configuration; thus, questions featuring a patient number with fewer or more than 8 digits invariably resulted in a *null* response. Regarding the second case, specific questions, for example, *"I am curious what the protocols for the drugs that work to treat cancer."* could potentially be answered by a knowledgeable individual or through QA tasks using external information resources. The third group, while seemingly akin to the second, differed in that the questions could technically be converted into SQL queries; however, they remained unanswerable due to the absence of the required information in the dataset. Example questions include: *"Has patient 23224 an appointment in another hospital department?"*. Further examples for each category, along with their respective frequencies, are detailed in appendix **??**.

Additionally, a quantitative analysis of unanswerable questions was also conducted using N-grams. By examining the differences in word distribution between 'answerable' and 'unanswerable' questions, as highlighted by the contrast between the first and last columns of appendix 9, significant disparities were noted. For instance, an examination of the 1-gram columns for both the training set and the Unanswerable Training set reveals that the only overlapping words are "patient" and "last." This indicates a significantly different distribution between the two datasets.

Based on both qualitative and quantitative analysis, we were able to refine our framework to avoid generating SQL queries for questions that solely comprise words found in the unanswerable questions of the training set.

## 6 Conclusion

Throughout the challenge, we noticed that differences in the way data is distributed across training, development, and test sets can make it hard for our model to determine which questions are answerable or not. To tackle this issue, we templatized questions to make the word distribution of development and test data more similar to the training data. This method aimed to bridge the gap between the datasets, helping the model better understand the features of unanswerable questions within the test dataset.

Although we did not address this in this study, we anticipate that future research could see performance improvements by augmenting the training dataset to more closely match the distribution of unanswerable questions in the development and test sets. Focusing on refining the test data to align more closely with the characteristics observed in the training datasets, we expect to increase model performance in identifying unanswerable questions. Such data augmentation strategies could bridge the remaining gaps between datasets and ensure a more robust model performance across varied datasets. Also, we utilized finetuned gpt-3.5-turbo, which is expensive and unusable for other researchers. Thus, further study should be done with open sourced models, like llama or gemma.

## Acknowledgments

## References

Andrew Brock, J. Donahue, Karen Simonyan, Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Donghyun Choi, M. Shin, EungGyun Kim, Xiang Deng, Ahmed Hassan Awadallah, Oleksandr Meek, Huan Polozov, Sun Matthew, Yujian Gan, Xinyun Chen, Qiuping Huang, John R Purver, Jinxia Woodward, Xie Peng-301, Amol Kelkar, Rohan Relan, V. Bhardwaj, and Saurabh Vaichal. 2021. S 2 sql: Injecting syntax to question-schema interaction graph encoder for text-to-sql parsers.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Ursin Brunner and Kurt Stockinger. 2021. Valuenet: A natural language-to-sql system that learns from database information. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pages 2177–2182. IEEE.

Ruichu Cai, Boyan Xu, Zhenjie Zhang, Xiaoyan Yang, Zijian Li, and Zhihao Liang. 2017. An encoder-decoder framework translating natural language to database queries. In *International Joint Conference on Artificial Intelligence*.

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. *ArXiv*, abs/2305.11853.

Xuemei Dong, C. Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, Lu Chen, Jinshu Lin, and Dongfang Lou. 2023. C3: Zero-shot text-to-sql with chatgpt. *ArXiv*, abs/2307.07306.

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. *ArXiv*, abs/2308.15363.

Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. *ArXiv*, abs/1707.07328.

Alistair Johnson, Tom Pollard, and Roger Mark. 2016a. Mimic-iii clinical database (version 1.4). *PhysioNet*, 10(C2XW26):2.

Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Li wei H. Lehman, Mengling Feng, Mohammad Mahdi Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. 2016b. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3.

Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *ArXiv*, abs/1705.03551.

George Katsogiannis-Meimarakis and Georgia Koutrika. 2023. A survey on deep learning approaches for text-to-sql. *The VLDB Journal*, 32:905–936.

Chia-Hsuan Lee, Oleksandr Polozov, and Matthew Richardson. 2021. Kaggledbqa: Realistic evaluation of text-to-sql parsers. In *Annual Meeting of the Association for Computational Linguistics*.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiaxi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, Nan Huo, Chenhao Ma, Kevin C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2023. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *ArXiv*, abs/2305.03111.

Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S. Yu. 2023a. A comprehensive evaluation of chatgpt's zero-shot text-to-sql capability. *ArXiv*, abs/2303.13547.

Hu Liu, Yuliang Shi, Jianlin Zhang, Xinjun Wang, Hui Li, and Fanyu Kong. 2023b. Multi-hop relational graph attention network for text-to-sql parsing. *2023 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8.

Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. 2018. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data*, 5(1):1–13.

Octavian Popescu, Irene Manotas, Ngoc Phuoc An Vo, Hangu Yeo, Elahe Khorashani, and Vadim Sheinin. 2022. Addressing limitations of encoder-decoder based approach to text-to-sql. In *International Conference on Computational Linguistics*.

Jiexing Qi, Jingyao Tang, Ziwei He, Xiangpeng Wan, Chenghu Zhou, Xinbing Wang, Quanshi Zhang, and Zhouhan Lin. 2022. Rasat: Integrating relational structures into pretrained seq2seq model for text-to-sql. *ArXiv*, abs/2205.06983.

Bowen Qin, Binyuan Hui, Lihan Wang, Min Yang, Jinyang Li, Binhua Li, Ruiying Geng, Rongyu Cao, Jian Sun, Luo Si, Fei Huang, and Yongbin Li. 2022. A survey on text-to-sql parsing: Concepts, methods, and future directions. *ArXiv*, abs/2208.13629.

Preethi Raghavan, Jennifer J. Liang, Diwakar Mahajan, Rachita Chandra, and Peter Szolovits. 2021. emrkbqa: A clinical knowledge-base question answering dataset. In *Workshop on Biomedical Natural Language Processing*.

Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for squad. *ArXiv*, abs/1806.03822.

Jaydeep Sen, Chuan Lei, Abdul Quamar, Fatma Özcan, Vasilis Efthymiou, Ayushi Dalmia, Greg Stager, Ashish R. Mittal, Diptikalyan Saha, and Karthik Sankaranarayanan. 2020. Athena++. *Proceedings of the VLDB Endowment*, 13:2747 – 2759.

Ruoxi Sun, Sercan Ö. Arik, Hootan Nakhost, Hanjun Dai, Rajarishi Sinha, Pengcheng Yin, and Tomas Pfister. 2023. Sql-palm: Improved large language model adaptation for text-to-sql. *ArXiv*, abs/2306.00739.

Richard Tarbell, Kim-Kwang Raymond Choo, Glenn Dietrich, and Anthony Rios. 2023. Towards understanding the generalization of medical text-to-sql models and datasets. *AMIA ... Annual Symposium proceedings. AMIA Symposium*, 2023:669–678.

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*, abs/2302.13971.

Immanuel Trummer. 2022. Codexdb: Synthesizing code for query processing from natural language instructions using gpt-3 codex. *Proc. VLDB Endow.*, 15:2921–2928.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2019. Rat-sql: Relation-aware schema encoding and linking for text-to-sql parsers. In *Annual Meeting of the Association for Computational Linguistics*.

Kun Xu, Lingfei Wu, Zhiguo Wang, Mo Yu, Liwei Chen, and Vadim Sheinin. 2018. Sql-to-text generation with graph-to-sequence model. *ArXiv*, abs/1809.05255.

Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Sebastian Riedel. 2020. Tabert: Pretraining for joint understanding of textual and tabular data. *arXiv preprint arXiv:2005.08314*.

Tao Yu, Rui Zhang, Kai-Chou Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Z Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. *ArXiv*, abs/1809.08887.

Xiang Yue, Bernal Jimenez Gutierrez, and Huan Sun. 2020. Clinical reading comprehension: A thorough analysis of the emrqa dataset. *ArXiv*, abs/2005.00574.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *ArXiv*, abs/1709.00103.

# A Examples of unanswerable questions

| Category | Example | Frequency |
|---|---|---|
| Incorrect Patient Number | Will they have any urine test done for patient 24628? | 252 (56%) |
| | Is patient 21074 subject to tests involving covid-19? | |
| | Do you know what type of blood patient 1903 has? | |
| Require External Knowledge | What is a checklist before lumb/lmbosac fus ant/ant? | 83 (18.4%) |
| | What is the protocol used for the anticancer drugs? | |
| | So tell me what to do before you go for hemodialysis. | |
| Out of EHR Knowledge Base | What kind of blood patient 18866 has. | 115 (25.6%) |
| | List the single rooms that are available now? | |
| | When are dr. oneill's rounds and procedures? | |

Table 6: Examples of Unanswerable Questions with Respective Frequencies

## B  Prompt

| Task Description |
|---|
| I have a database related to healthcare that consists of 17 tables, each holding various pieces of data about hospital operations. I need to convert a natural language question into an SQL query to retrieve specific information from this database. Could you construct an SQL query that accurately reflects the question, considering the structure and details of my database? My database uses SQLite, and I'm looking for a query that's optimized for accuracy and efficiency. The detailed description of 17 tables is given below. Note that the patient number is an eight-digit number and current year is 2100 and all table values are in lower case. |

| Table Information – <Step 2> Initial SQL Generation | |
|---|---|
| **Natural Language Form** | **SQL Form** |
| 1.admissions:<br>  Documents each hospitalization event.<br>  • admissions [row_id, subject_id, hadm_id (hospital admission ID), admittime (admission time), dischtime (discharge time), admission_type, admission_location, discharge_location, insurance, language, marital_status, age];<br>2. d_icd_diagnoses:<br>  A reference for ICD-9 diagnosis codes<br>  • d_icd_diagnoses [row_id, icd_code, long_title];<br>3. d_icd_procedures:<br>  A reference for ICD-9 procedure codes.<br>  • d_icd_procedures [row_id, icd_code, long_title];<br>4. d_labitems:<br>  Acts as a dictionary for lab test ITEMIDs.<br>  • d_labitems [row_id, itemid, label]<br><br><br><br><br><br>... | 1.admissions:<br>  Documents each hospitalization event.<br>  CREATE TABLE admissions (<br>    row_id INTEGER,<br>    subject_id INTEGER REFERENCES patients(subject_id),<br>    hadm_id INTEGER,<br>    admittime TIMESTAMP,<br>    dischtime TIMESTAMP,<br>    admission_type TEXT,<br>    admission_location TEXT,<br>    discharge_location TEXT,<br>    insurance TEXT,<br>    language TEXT,<br>    marital_status TEXT,<br>    age INTEGER<br>  );<br>2. d_icd_diagnoses:<br>  A reference for ICD-9 diagnosis codes.<br>  CREATE TABLE d_icd_diagnoses (<br>    row_id INTEGER,<br>    icd_code TEXT,<br>    long_title TEXT<br>  );<br>3. d_icd_procedures:<br>  A reference for ICD-9 procedure codes.<br>  CREATE TABLE d_icd_procedures (<br>    row_id INTEGER,<br>    icd_code TEXT,<br>    long_title TEXT<br>  );<br>4. d_labitems:<br>  Acts as a dictionary for lab test ITEMIDs.<br>  CREATE TABLE d_labitems (<br>    row_id INTEGER,<br>    itemid INTEGER,<br>    label TEXT<br>  );<br>… |
| **Templatized Question** | |
| Count how many times patient 10004457 had 'coronary arteriography using two catheters' during their first hospital visit. | |

Table 7: The prompt used in the step 2 initial SQL generation.

## C  Example Appendix

| **Task Description** |
|---|
| I have a database related to healthcare that consists of 17 tables, each holding various pieces of data about hospital operations. I need to convert a natural language question into an SQL query to retrieve specific information from this database. Could you construct an SQL query that accurately reflects the question, considering the structure and details of my database? My database uses SQLite, and I'm looking for a query that's optimized for accuracy and efficiency. The detailed description of 17 tables is given below. Note that the patient number is an eight-digit number and current year is 2100 and all table values are in lower case. |

| **Table Information – <Step 3> SQL Regeneration** |
|---|

1. "procedures_icd:
   - Records procedures using ICD codes.
   - procedures_icd (row_id, subject_id, hadm_id, icd_code, charttime)
   - Description of Columns:
       row_id: Unique record identifier;
       subject_id: Unique identifier assigned to each individual patient;
       hadm_id (hospital admission ID): Unique identifier assigned to each separate hospital admission of a patient;
       icd_code: ICD code for the procedure performed;
       charttime: Timestamp of when the procedure was documented in the patient's chart",
2. "d_icd_procedures:
   - A reference for ICD-9 procedure codes.
   - d_icd_procedures (row_id, icd_code, long_title)
   - Description of Columns:
       row_id: Unique record identifier;
       icd_code: Unique ICD-9 procedure code;
       long_title: Detailed description of the procedure.",
3. "admissions:
   - Documents each hospitalization event.
   - admissions (row_id, subject_id, hadm_id, admittime, dischtime, admission_type, admission_location, discharge_location, insurance, language, marital_status, age)
   - Description of Columns:
       row_id: Unique record identifier;
       subject_id: Unique identifier assigned to each individual patient;
       hadm_id (hospital admission ID): Unique identifier assigned to each separate hospital admission of a patient;
       admittime: Admission time to the hospital;
       dischtime: Discharge time from the hospital;
       admission_type: Type of hospital admission;
       admission_location: Location from where the patient was admitted;
       discharge_location: Location to where the patient was discharged;
       insurance: Patient's insurance type;
       language: Patient's primary language;
       marital_status: Patient's marital status;
       age: Patient's age at admission.",

| **Templatized Question** |
|---|
| Count how many times patient 10004457 had 'coronary arteriography using two catheters' during their first hospital visit. |

Table 8: The prompt used in Step 3 for SQL regeneration.

# D Word frequencies: 3-gram

| N-gram | 1-gram | 2-gram | 3-gram |
|---|---|---|---|
| **Train set** | ('patient',): 3205<br>('since',): 1572<br>('last',): 1394<br>('hospital',): 1340<br>('first',): 1232<br>('year',): 934<br>('patients',): 861<br>('2100',): 818<br>('visit',): 778<br>('time',): 734 | ('hospital', 'visit'): 608<br>('since', '2100'): 425<br>('first', 'hospital'): 327<br>('last', 'time'): 317<br>('hospital', 'encounter'): 316<br>('last', 'hospital'): 302<br>('since', '1'): 302<br>('1', 'year'): 298<br>('year', 'ago'): 298<br>('first', 'time'): 280 | ('since', '1', 'year'): 298<br>('1', 'year', 'ago'): 298<br>('first', 'hospital', 'visit'): 203<br>('last', 'hospital', 'visit'): 183<br>('within', '2', 'months'): 153<br>('last', 'time', 'patient'): 118<br>('first', 'time', 'patient'): 104<br>('first', 'hospital', 'encounter'): 94<br>('last', 'hospital', 'encounter'): 92<br>('measured', 'last', 'hospital'): 92 |
| **Dev set** | ('patient',): 656<br>('hospital',): 322<br>('since',): 316<br>('patients',): 279<br>('last',): 255<br>('first',): 253<br>('year',): 202<br>('visit',): 199<br>('2100',): 170<br>('time',): 134 | ('hospital', 'visit'): 150<br>('since', '2100'): 82<br>('last', 'hospital'): 72<br>('since', '1'): 70<br>('1', 'year'): 69<br>('year', 'ago'): 69<br>('first', 'time'): 67<br>('hospital', 'encounter'): 65<br>('first', 'hospital'): 59<br>('lab', 'test'): 46 | ('since', '1', 'year'): 69<br>('1', 'year', 'ago'): 69<br>('last', 'hospital', 'visit'): 46<br>('first', 'hospital', 'visit'): 40<br>('within', '2', 'months'): 30<br>('last', 'hospital', 'encounter'): 20<br>('first', 'time', 'patient'): 20<br>('current', 'hospital', 'visit'): 19<br>('arterial', 'blood', 'pressure'): 19<br>('top', 'three', 'frequent'): 18 |
| **Test set** | ('patient',): 620<br>('hospital',): 318<br>('patients',): 306<br>('since',): 304<br>('last',): 265<br>('first',): 234<br>('year',): 217<br>('2100',): 184<br>('visit',): 164<br>('prescribed',): 140 | ('hospital', 'visit'): 128<br>('since', '2100'): 100<br>('hospital', 'encounter'): 81<br>('first', 'hospital'): 72<br>('last', 'hospital'): 72<br>('since', '1'): 57<br>('1', 'year'): 57<br>('year', 'ago'): 57<br>('many', 'patients'): 54<br>('number', 'patients'): 47 | ('since', '1', 'year'): 57<br>('1', 'year', 'ago'): 57<br>('first', 'hospital', 'visit'): 42<br>('last', 'hospital', 'visit'): 42<br>('within', '2', 'months'): 41<br>('last', 'hospital', 'encounter'): 26<br>('measured', 'last', 'hospital'): 23<br>('first', 'time', 'patient'): 20<br>('first', 'hospital', 'encounter'): 20<br>('last', 'time', 'patient'): 18 |
| **Unanswerable Train set** | ('patient',): 252<br>('department',): 49<br>('tell',): 42<br>('procedure',): 41<br>('blood',): 36<br>('dr',): 36<br>('received',): 34<br>('rooms',): 29<br>('test',): 28<br>('last',): 27 | ('received', 'department'): 20<br>('outpatient', 'schedule'): 18<br>('rounds', 'procedures'): 17<br>('another', 'department'): 16<br>('rooms', 'available'): 15<br>('diagnosis', 'patient'): 15<br>('operating', 'rooms'): 14<br>('blood', 'transfusion'): 14<br>('name', 'diagnosis'): 14<br>('genetic', 'test'): 14 | ('name', 'diagnosis', 'patient'): 12<br>('last', 'time', 'patient'): 11<br>('many', 'operating', 'rooms'): 11<br>('appointment', 'another', 'department'): 11<br>('genetic', 'test', 'patient'): 10<br>('subject', 'covid-19', 'testing'): 9<br>('type', 'blood', 'patient'): 9<br>('ward', 'id', 'patient'): 9<br>("today's", 'outpatient', 'schedule'): 9<br>('outpatient', 'schedule', 'dr'): 8 |

Table 9: 3-Gram frequency table with 10 examples sorted in non-increasing order

## E   Examples of Question Templatization

| Input: "When does patient 8016's influenza quarantine end?" | |
|---|---|
| Candidate templetes | "What was the time of <patient number>'s last influenza a/b by dfa microbiology test since 03/2100?" |
| | "Can you tell me when <patient number> had their first rapid respiratory viral screen & culture microbiology test in 08/this year?" |
| | "When did <patient number> depart hospital during this year for the last time?" |
| Reformulated: "When is the end date of patient 8016's influenza quarantine?" | |
| **Input: "Pull up the IDs of patients who were diagnosed with cataract extraction status."** | |
| Candidate templetes | "Number of patients who were diagnosed with <unk>." |
| | "Number of patients who were diagnosed throughout this year with <unk>." |
| | "Tell me the number of patients diagnosed with <unk>." |
| Reformulated: "Retrieve the ids of patients diagnosed with 'cataract extraction status'." | |
| **Input: "How many duloxetine prescription cases were there since 1 year ago?"** | |
| Candidate templetes | "How much duloxetine has been prescribed to <patient number> in 05/2100 in total?" |
| | "How many drugs have been prescribed to <patient number> since 2 months ago?" |
| | "What is the number of drugs <patient number> was prescribed since 1 year ago?" |
| Reformulated: "What is the number of duloxetine prescription cases since 1 year ago?" | |

Table 10: Examples of templatized questions using question masked template.

# ProbGate at EHRSQL 2024: Enhancing SQL Query Generation Accuracy through Probabilistic Threshold Filtering and Error Handling

**Sangryul Kim**[1†]    **Donghee Han**[2]    **Sehyun Kim**[3]

[1]KAIST AI
[2]KAIST Graduate School of Data Science
[3]KAIST Bio and Brain Engineering
{sangryul, handonghee, sehyun}@kaist.ac.kr

## Abstract

Recently, deep learning-based language models have significantly enhanced text-to-SQL tasks, with promising applications in retrieving patient records within the medical domain. One notable challenge in such applications is discerning unanswerable queries. Through fine-tuning model, we demonstrate the feasibility of converting medical record inquiries into SQL queries. Additionally, we introduce an entropy-based method to identify and filter out unanswerable results. We further enhance result quality by filtering low-confidence SQL through log probability-based distribution, while grammatical and schema errors are mitigated by executing queries on the actual database. We experimentally verified that our method can filter unanswerable questions, which can be widely utilized even when the parameters of the model are not accessible, and that it can be effectively utilized in practice[1].

## 1 Introduction

In recent years, the field of natural language processing (NLP) has witnessed remarkable progress driven by transformer-based large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023; Roziere et al., 2023). A prevailing approach involves fine-tuning pre-trained language models with new data across various tasks, facilitating transfer learning (Min et al., 2023). This methodology has proven effective in tasks like document summarization, entity-relationship extraction, document classification, and sentiment analysis. One of the main tasks where these language models are increasingly leveraged is text-to-SQL (Text2SQL), which converts natural language queries into SQL queries (Mellah et al., 2020).

Text2SQL presents unique challenges distinct from conventional NLP tasks. Firstly, it demands



Figure 1: Determines whether a question and the generated SQL are answerable or unanswerable based on the log probability of the tokens generated by the Text2SQL model. If the log probability of a token falls below a certain threshold, we classify the question and SQL as unanswerable.

grammatical correctness, as even minor errors can render SQL queries unexecutable. Unlike document summarization, where semantic correctness compensates for grammatical inaccuracies, SQL queries must adhere strictly to syntax rules (Cao et al., 2023). Secondly, schema awareness is crucial; understanding the database structure is essential for generating accurate SQL queries (Katsogiannis-Meimarakis and Koutrika, 2023). Finally, discerning unanswerable queries is vital, especially in domains like healthcare where incorrect or incomplete information can have severe consequences (Lee et al., 2022). If users do not inspect the SQL queries themselves, but only receive the results of the execution, the results of an incorrect SQL execution can be fatally misleading.

In the domain of Text2SQL, effectively filtering out unanswerable questions presents a significant

---

[†]Corresponding Author
[1]Code and datasets are available at https://github.com/venzino-han/probgate_ehrsql

challenge (Lee et al., 2024a), particularly within the medical field where accuracy is paramount. Existing methodologies for identifying unanswered queries have primarily targeted cases where such queries exhibit discernible patterns (Wang et al., 2023). However, these methods are often tailored to specific model architectures and learning methods, thereby constraining their direct applicability to LLM services accessible via APIs, such as ChatGPT. Given the recent widespread adoption of such managed LLMs across various industries, the need for a more versatile and adaptable approach to filtering unanswered questions becomes increasingly pronounced. This underscores the necessity for innovative solutions that can seamlessly integrate into existing LLM services, ensuring robust performance in diverse application scenarios, including medical contexts.

We address solutions that effectively solve the challenges of Text2SQL tasks through a subset focusing on Electronic Health Records(EHR), utilizing medical questions and corresponding SQL queries relevant to medical systems used in real hospitals (Lee et al., 2022). Specifically, we participate in the EHRSQL Shared Task on Reliable Text-to-SQL Modeling On Electronic Health Records (Lee et al., 2024b). A distinctive feature of this shared task is that under the basic premise of generating appropriate SQL statements for given natural language queries, not all questions are answerable; some are unanswerable. Moreover, beyond merely generating suitable SQL statements for questions, this task is complex as it requires distinguishing between answerable and unanswerable questions and considering the high penalties for incorrectly identifying the questions are answerable or not, thus necessitating both reliability and accuracy in execution.

In this paper, we introduce **Prob**ability **Gate (ProbGate)**, a novel probability-based filtering approach designed for seamless integration with diverse generative language models, without requiring direct access to the model's parameters. Figure 1 illustrates the concept of ProbGate, which leverages the logarithmic probability of individual tokens to assess the uncertainty associated with generated SQL queries. We consider the log probability of specific target tokens as an indicator of how confident the model is and how well it can perform the task without hallucinations. We found that utilizing logarithmic probability-based confidence to identify answerable and unanswerable questions

was very effective, which is a key aspect of this task.

We evaluate the efficacy of ProbGate through experimentation with Electronic Health Record (EHR) SQL dataset (Lee et al., 2022). Specifically, we apply ProbGate to both T5-based (Raffel et al., 2020) Text2SQL models and gpt-3.5-turbo finetuned models, comparing their performance against conventional binary classifiers. Additionally, we train binary classifiers based on T5 and gpt-3.5-turbo model to filter out unaswerable questions. Our experimental findings reveal that ProbGate outperforms binary classifiers in terms of both performance and resilience to shifts in data distribution. These results underscore the potential of ProbGate as a versatile and robust filtering solution for a wide range of applications.

Our contributions and methods can be summarized as follows:

- Through our experiments, we found that the fine-tuned gpt-3.5-turbo performed well at generating SQL queries for questions, but was less able to distinguish and filter out unanswerable questions.

- We present the Probabilistic Threshold Filtering method(ProbGate) to effectively distinguish between answerable and unanswerable questions in datasets containing a mix of both.

- We demonstrate an effective method by creating a single pipeline from training to testing, incorporating SQL execution error handling, showing that it can be applied to similar cases.

## 2 Backgrounds

**Text2SQL** Databases serve as powerful tools for efficiently querying extensive datasets. However, accessing this data often requires users to possess knowledge of query languages like SQL. To democratize this process and render it accessible across proficiency levels, significant research efforts have focused on techniques for interpreting natural language questions and autonomously translating them into SQL queries. Recent strides in deep learning methodologies, particularly transformer-based language models, have spurred the development of text-to-SQL techniques. These approaches aim to bridge the gap between natural language queries and SQL commands, thereby enhancing accessibility and usability in database querying tasks

(Mellah et al., 2020; Katsogiannis-Meimarakis and Koutrika, 2023).

Early Text2SQL research relied on rule-based and template-based methods, but more recently, deep learning-based methodologies have become mainstream (Deng et al., 2022). Deep learning methodologies exhibit robustness on the data they are trained on but often struggle to generalize to unseen database schemas. To mitigate this challenge, researchers have explored approaches to encode database relationships and leverage column relationships using self-attention mechanisms (Wang et al., 2020). In Text2SQL, ensuring the accuracy of generated SQL statements is crucial as even minor errors can lead to failures in query execution. Recent studies have demonstrated the effectiveness of utilizing LLMs like gpt-3.5-turbo to rectify SQL statements derived from natural language queries, addressing the challenge of proofreading SQL output (Pourreza and Rafiei, 2024).

One of the main applications of Text2SQL is its utilization in the healthcare domain, specifically to handle complex tasks within electronic health records (EHRs). Recent research has shown that decomposing these tasks into manageable pieces can improve the performance of multi-table reasoning within EHRs. The authors proposed to iteratively improve SQL queries by incorporating interactive coding and execution feedback mechanisms to learn from the error messages encountered. This iterative improvement process proved to be effective and resulted in noticeable improvements in SQL performance in the healthcare domain (Shi et al., 2024). In a closely related investigation, researchers observed that EHR data is commonly stored in relational databases, which can be represented as directed acyclic graphs. Leveraging this insight, they employed a graph-based methodology to capture the intricate relationships between tables, entities, and values within relational databases (Park et al., 2021).

**Confidence of Generated Tokens** The outputs of LLMs are typically based on a next token prediction method, where the probability of previous tokens is used to predict the next one. During this process, a phenomenon often referred to as 'hallucination' can occur, which results in incorrect inferences about the task(Wang and Sennrich, 2020; Xiao and Wang, 2021; Li et al., 2022). Additionally, previous research has shown that low probability and confidence levels can indicate a lack of knowledge in the model(Kadavath et al., 2022). To overcome this, Jiang et al. (2023) introduced a structure named FLARE, which includes a mechanism where if the probability of a token generated by the model falls below a certain threshold, the token is used as a query to retrieve relevant documents from a retriever. This approach aims to address the lack of knowledge and increase confidence. In our work, we also propose a filtering model using log probability to determine if log probability can effectively distinguish the uncertainty in generated content.

## 3 Methods

From this section, we cover the contents related to the methods. In §3.1, there is a detailed description of the shared task dataset; in §3.2, the main metrics used in the shared task are discussed; and from §3.3 to §3.5, detailed information on the main methods is provided. The entire architecture can be referenced in Figure 2.

### 3.1 Datasets

The dataset employed in this study is sourced from the EHRSQL Shared Task on Reliable Text-to-SQL Modeling On Electronic Health Records(EHRSQL-2024) (Lee et al., 2024b), with the purpose of simplifying access to EHR data by automatically translating natural language questions into corresponding SQL queries. This dataset is referred to as The MIMIC-IV demo version of EHRSQL with additional unanswerable questions. It consists of various questions related to medical records and their corresponding SQL queries, serving as a crucial resource for natural language processing and SQL query generation research. The specific attributes and composition follow the study by EHRSQL (Lee et al., 2022). The EHRSQL dataset is based on questions frequently asked in the medical field, gathered from 222 hospital personnel, including physicians, nurses, insurance assessors, and health records teams. These questions have been reconstructed to reflect various scenarios that can occur in real-world medical contexts and are presented as a dataset annotated with SQL queries aligned with the hierarchical structure of EHR databases.

The primary characteristics of this dataset are as follows: it encapsulates the diverse demands of hospital settings, encompassing tasks from straightforward information retrieval to the more intricate operations such as identifying the top N prescribed

Figure 2: Our method's overall architecture is as follows: During training, we fine-tune the gpt-3.5-turbo model using a dataset from which unanswerable cases have been removed. Subsequently, we identify unanswerable cases using filtering based on log probability and filtering through SQL execution, ultimately deriving the answers.

drugs following a disease diagnosis. Additionally, it incorporates a range of temporal expressions within the questions. Lastly, it includes not only answerable questions but also unanswerable ones that are incompatible with the database schema or require external domain knowledge.

The EHRSQL-2024 task provides a training dataset consisting of questions about medical records, SQL queries corresponding to the MIMIC-IV demo version, and instances annotated as 'null' for unanswerable questions. The test dataset comprises only questions, including types of unanswerable questions that are not included in the training data. The training and test datasets comprise 5124 and 1167 examples, respectively.

## 3.2 Metric

In the medical and healthcare domains, reliability is particularly emphasized. Therefore, the model's responses must be accurate, and it's better to abstain from answering than to risk errors. From this perspective, we employ the RS (Reliability Score) introduced in TrustSQL(Lee et al., 2024a) to assess the model's performance. The RS assigns scores for accurate predictions, providing an evaluation of the model's performance, while also penalizing incorrect predictions and instances where the model attempts to respond to unanswerable questions.

$$
\phi_c(x) = \begin{cases} 1 & \text{if } x \in Q_{\text{ans}}; g(x) = 1; \text{Acc}(x) = 1, \\ 0 & \text{if } x \in Q_{\text{ans}}; g(x) = 0, \\ -c & \text{if } x \in Q_{\text{ans}}; g(x) = 1; \text{Acc}(x) = 0, \quad (1) \\ -c & \text{if } x \in Q_{\text{una}}; g(x) = 1, \\ 1 & \text{if } x \in Q_{\text{una}}; g(x) = 0. \end{cases}
$$

In EQ(1), $\text{Acc}(x)$ represents the execution accuracy, where for any $x$ belonging to the set of answerable questions ($Q_{ans}$), if $f(x)$ matches the correct answer, it returns 1, and otherwise, it returns 0. The function $g(x)$ indicates whether the model generates an SQL query, where 1 signifies generation and 0 indicates no generation. The parameter $c$ serves as the penalty parameter. A penalty of $-c$ is imposed in two scenarios: when $x$ is in $Q_{ans}$ and the generated query is incorrect, and when $x$ is in the set of unanswerable questions ($Q_{una}$) but a query is generated regardless. The model earns a score of 1 when it correctly answers a question. The final Reward Score (RS) is obtained by calculating the average of $\phi_c(x)$ scores across all samples. The penalty factor $c$ can be adjusted to evaluate the model's reliability, particularly in scenarios requiring high confidence. In our experiments, we consider four options for the penalty, $c = 0, 5, 10, N$, where $N$ represents the total number of samples being evaluated. This metric proves valuable in assessing the model's ability to reliably generate SQL queries and to respond only to questions that are answerable.

690

### 3.3 Fine-Tuning and Prompt Design

**Fine-tuning** As the first step in solving the task, we fine-tune the OpenAI gpt-3.5-turbo-0125 model[2]. This is used for Text2SQL conversion, serving as an easy-to-use baseline and also providing a convenient API for subsequent log probability calculations. To minimize noise in the dataset, we exclude unanswerable data from training, focusing solely on SQL transformation without considering whether the given questions are answerable or not. Out of the 5124 samples in the training set, 450 unanswerable data points were excluded, leaving 4674 question-query pairs that are answerable. These data consist of natural language questions paired with their corresponding correct SQL queries. The example of the input-output format for the training dataset can be found in the Appendix B.

**Prompt** During the training and inference phase, we experiment with various prompt formats to facilitate the model's ability to receive a question and generate the corresponding SQL query accurately. As an illustration, the following structure is utilized for prompts:

> **Optimized Prompt**
>
> "You are 'SQLgpt', an AI designed to convert natural language questions into their corresponding SQL queries. It is imperative that the generated SQL queries conform to the standard SQL format and are not enclosed within quotes (neither single ' nor double "). Your primary objective is to precisely generate the exact SQL query for each presented question."

Such prompts aim to guide the model towards generating the most appropriate SQL query in response to a question while also preventing the occasional generation of SQL queries encased within ' or " symbols, which can potentially lead to errors within the database.

### 3.4 Probabilistic Threshold Filtering (ProbGate)

In the test set of the given task, we can see that answering all questions as unanswerable results in

---

---

**Algorithm 1** ProbGate

```
 1: reserved ← ["SELECT", ...]
 2: procedure CALCLOGBOTTOMK(log, t)
 3:     LogProb ← []
 4:     for x in log do
 5:         if x.token not in reserved then
 6:             LogProb.append(x.logprob)
 7:         end if
 8:     end for
 9:     Keep bottom t values of sorted(LogProb)
10:     return average(LogProb)
11: end procedure
```

a score of 19.97 across all RS metrics. By assuming all questions to be answerable and submitting answers accordingly, we were able to achieve a score of 73.52 on the RS(0) metric, in an effort to understand the performance of the model fine-tuned in the previous step on answerable questions. Interpreting this from a ratio perspective, since we already know that 19.97% of the test set is unanswerable, it implies that 80.03% of it consists of answerable questions. Therefore, we can deduce that the percent accuracy of the model on answerable questions is approximately 91.87%. This implies a percent accuracy of 91.87%, which suggests that to avoid losing points, the threshold for ideally identifying unanswerable questions should be set higher than the scale used to find this threshold, as inferred from the results. Given that the total number of items in the test set for the given task is 1167, we can deduce that to minimize the penalty $-c$ and maximize the score, we find the threshold $k$ in test dataset value should be approximately 425 according to the empirical findings.

To distinguish unanswerable SQL statements, we assume that tokens of each generated SQL with low log probability are likely candidates for unanswerability, considering the log probabilities of the tokens as confidence scores. Since we previously determined the number of unanswerable candidates, or the threshold, we calculate the log probabilities of each SQL token in the test set items, sort them by ascending order of average value of its log probability, and consider all items with indices from the first up to the threshold as unanswerable. We incorporate some additional tricks, taking into account the characteristics of the SQL statement. The given text2SQL task is considered a highly structured sequence-to-sequence task due to the nature of SQL query syntax, which is very struc-

| Model | RS(0) | RS(5) | RS(10) | Rs(N) |
|---|---|---|---|---|
| T5-small FT + Filtering | 47.81 | 45.66 | 43.51 | **-452.19** |
| T5-Large-text2sql-spider FT + Filtering | 74.63 | 59.59 | 44.54 | -3425.37 |
| T5-Large-text2sql-spider FT + Classifier(T5) | 63.80 | 18.23 | -27.34 | -10536.20 |
| T5-Large-text2sql-spider FT + Filtering + Classifier(T5) | 72.74 | 58.56 | 44.37 | -3227.26 |
| gpt-3.5-turbo FT + Classifier(T5) | 90.28 | 51.59 | 12.89 | -8109.02 |
| gpt-3.5-turbo FT + Classifier(gpt-3.5-turbo) | 88.05 | 57.95 | 27.86 | -6911.95 |
| gpt-3.5-turbo FT + ProbGate(t=387) | **85.30** | **80.57** | **75.84** | -1014.70 |

Table 1: Model Selection and Ablation Study in Dev Phase dataset. In the case of the T5-Large model, it is a model that was first fine-tuned using the Spider dataset, which is one of the Text2SQL datasets, and then subsequently trained on the EHRSQL dataset. In abbreviation, "FT" stands for Fine-Tuning. 'Filtering' and 'Classifier' are described in section §4.1.

| Model | RS(0) | RS(5) | RS(10) | Rs(N) |
|---|---|---|---|---|
| gpt-3.5-turbo FT | 73.52 | -58.87 | -191.25 | -30826.47 |
| gpt-3.5-turbo FT + ProbGate(t=450) | 79.43 | 73.01 | 66.58 | -1420.57 |
| gpt-3.5-turbo FT + ProbGate(t=450) + GEF | 79.78 | 75.92 | 72.06 | -820.22 |
| **gpt-3.5-turbo FT + ProbGate(t=425) + GEF** | **81.92** | **78.06** | **74.21** | **-818.08** |

Table 2: The results from applying our methodology during the Test Phase are as follows. The results of ablation at each filtering stage are provided, and it can be observed that there is an improvement in performance at every stage. In abbreviation, "FT" stands for Fine-Tuning, and "GEF" refers to Grammatical Errors Filtering, as introduced in section §3.5.

tured compared to the form of the input. The SQL statement inferred from the model can be broadly divided into two parts: reserved words of SQL syntax such as SELECT, AS, BETWEEN; and entities and attributes. We consider that the model is more likely to hallucinate when generating entities and attributes than when generating reserved words. Hence, when calculating the log probability for each test set item, we exclude reserved words(tokens) and compute it for the remaining tokens. The excluded reserved words can be found in Appendix A. Moreover, to make the distinction between answerable and unanswerable even clearer based on log probability, we also impose a limitation on the value of lowest $t$ tokens($t = 10$ in this case), guiding the calculation towards the average value of these lowest log probability tokens. The algorithm for calculate log probability with one individual data can be found in Algorithm 1.

### 3.5 Grammatical Errors Filtering

In the last stage, we execute generated answerable SQL queries filtered by ProbGate through given database, if there is an error when executing SQL queries, we consider them unanswerable. The necessity of this stage arises because grammatical

errors that are not fully caught by the previous ProbGate stage can only be detected by actual execution Although the query might actually have an answer and could be an answerable example, we consider it unanswerable to avoid penalties. This is because we can convert the penalty for incorrect answers, the $-c$ score, into 0. Reflecting on real-world scenarios, generating a response from the model indicating it does not know the answer could be more beneficial for the model's robustness and safety than returning incorrect results.

## 4 Results and Analysis

### 4.1 Model Selection and Ablation Study

As our final methodology, the base model gpt-3.5-turbo is relatively difficult to access the weights or perform additional analysis compared to other open-source models, so we use one of the Seq2Seq models, the T5 model, as a comparison model. Additionally, we compare using filtering based on maximum entropy, as utilized in (Lee et al., 2022), as our filtering model. Lastly, we also train a binary classifier with both T5 and gpt-3.5-turbo to distinguish between answerable and unanswerable questions to see its impact on performance. The re-

Figure 3: **Left** - Log Probability Distribution of the Fine-Tuned Model, **Right** - Log Probability Distribution of the Unfine-Tuned Model

sults are shown in Table 1, and conclusively, none of the methodologies surpasses the performance of the methodology applying gpt-3.5-turbo FT + ProbGate. The reason for this is observed in the accuracy of Text2SQL, where the gpt-3.5-turbo model, with its larger parameters and more advanced tuning methods, outperforms models from the T5 series. Additionally, it is interpreted that the Classifier does not show significant effectiveness due to the too different distribution between the training and the remaining dataset, and the task's high penalty for errors.

## 4.2 ProbGate and Grammatical Errors Filtering

The best results for the test set are achieved using our pipeline architecture, as shown in Table 2. The process involves fine-tuning the gpt-3.5-turbo model with data excluding unanswerable data, then prioritizing the filtering of unanswerable data with ProbGate set to a threshold of 425, and finally applying Grammatical Errors Filtering. This sequence shows progressively better metric values. Additionally, we can interpret that the smaller the gap between the scores of RS{0, 5, 10, N}, the fewer penalties our model receives. Our final architecture can be seen as achieving the narrowest gap among these scores.

## 4.3 Log Probability Distribution between Answerable and Unanswerable.

In this section, we analyze the log probability distribution of SQL queries generated by the gpt-3.5-turbo model and compare the differences in distribution based on whether the model is finetuned or not. For the experiments with the finetuned model, we first divide the training dataset into a 7:3 ra-

tio, using 70% of dataset to finetune gpt-3.5-turbo with only answerable data. The remaining 30% includes both answerable and unanswerable data, enabling the extraction of log probabilities during the model's SQL inference process. In left graph of Figure 3, red represents null data, while blue indicates answerable data. The X-axis represents the log probability, and the Y-axis represents the number of data points with that log probability. As a result, it is observed that answerable data exhibited higher log probabilities, whereas null data show relatively lower probabilities, revealing the uncertainty in the generated SQL. The right graph of Figure 3 displays the log probability distribution of SQL generated by an unfine-tuned gpt-3.5-turbo model under the same conditions. The difference in log probability distributions based on answerability is not significant, making it difficult to distinguish labels in the distribution. These results underscore the effectiveness of fine-tuning on answerable data, indicating that fine-tuning significantly increases the log probability of the model for answerable data while also creating a discernible distribution difference with unanswerable data. By leveraging this distributional difference, ProbGate suggests that by setting an optimal threshold to treat all data that is either unanswerable or has uncertain generation outcomes as unanswerable, it can enhance response stability and reliability.

## 5 Conclusion

We participate in the EHRSQL Shared Task on Reliable Text-to-SQL Modeling On Electronic Health Records, as detailed in (Lee et al., 2024b), aiming to develop a reliable and high-performance Text2SQL method. This encompasses the chal-

lenge of generating appropriate SQL for answerable questions while also distinguishing unanswerable questions within datasets that include them. To solve this, we fine-tune LLMs on the training dataset and then employ a filtering pipeline called ProbGate, which consists of a combination of probabilistic threshold filtering and grammatical errors filtering, effectively executing the task. Additionally, through an ablation study and detailed analysis, we demonstrate that our method can be effectively used for tasks with a high sensitivity to errors. Ultimately, using this method, we conclude the shared task with a team ranking of 3rd place.

## Limitations

The methodology discussed here is central to solving competitive, contest-style shared tasks, with discussions taking place at a time when labels for the development and test sets, excluding training data, have not been disclosed. Therefore, our methodology greedily constructs the architecture to maximize the score on the main evaluation metric of the shared task, RS(10). Consequently, the primary parameters used in the model (e.g., threshold value, t value of ProbGate, etc.) can be specifically adjusted for the data and are sensitive to new datasets, meaning parameter values have a significant impact on the overall performance of the architecture. The performance of the basic model, which depends on the performance of the Fine-tuning model, is tied to a specific model (gpt-3.5-turbo) that is not open-sourced. Therefore, additional experiments with Text2SQL specialized open-source LLMs(Li et al., 2023) are needed. These limitations increase in severity when the distribution of unanswered questions differs between training and test datasets. Therefore, further research on unanswered question filtering approaches from an out-of-distribution detection perspective is warranted.

## Ethics Statement

Throughout this research, we are using the gpt-3.5-turbo model as a baseline. It's acknowledged that depending on the inputs provided by users, the model's outputs may include harmful content or exhibit unintended biases. Recognizing and addressing these potential issues is essential for deploying this technology in real-world production environments. This entails a necessity for additional engineering tuning aimed at minimizing such side effects, highlighting a commitment to responsible AI use and the importance of continual improvement to ensure ethical deployment. Furthermore, the gpt-3.5-turbo model, which is used as our primary method, has not publicly disclosed its weights or training processes. There is also a risk that private data may be exposed during fine-tuning. Therefore, when handling sensitive data, it is advisable to switch to an open-source model or exercise caution.

## Acknowledgments

## References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Ruisheng Cao, Lu Chen, Jieyu Li, Hanchong Zhang, Hongshen Xu, Wangyou Zhang, and Kai Yu. 2023. A heterogeneous graph to abstract syntax tree framework for text-to-sql. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(11):13796–13813.

Naihao Deng, Yulong Chen, and Yue Zhang. 2022. Recent advances in text-to-SQL: A survey of what we have and what we expect. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2166–2187, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992, Singapore. Association for Computational Linguistics.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.

George Katsogiannis-Meimarakis and Georgia Koutrika. 2023. A survey on deep learning approaches for text-to-sql. *The VLDB Journal*, 32(4):905–936.

Gyubok Lee, Woosog Chay, Seonhee Cho, and Edward Choi. 2024a. Trustsql: A reliability benchmark for

text-to-sql models with diverse unanswerable questions. *arXiv preprint arXiv:2403.15879*.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024b. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Raymond Li, Loubna Ben allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia LI, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Joel Lamy-Poirier, Joao Monteiro, Nicolas Gontier, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Ben Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason T Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Urvashi Bhattacharyya, Wenhao Yu, Sasha Luccioni, Paulo Villegas, Fedor Zhdanov, Tony Lee, Nadav Timor, Jennifer Ding, Claire S Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro Von Werra, and Harm de Vries. 2023. Starcoder: may the source be with you! *Transactions on Machine Learning Research*. Reproducibility Certification.

Shaobo Li, Xiaoguang Li, Lifeng Shang, Zhenhua Dong, Chengjie Sun, Bingquan Liu, Zhenzhou Ji, Xin Jiang, and Qun Liu. 2022. How pre-trained language models capture factual knowledge? a causal-inspired analysis. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1720–1732, Dublin, Ireland. Association for Computational Linguistics.

Youssef Mellah, Hassane El Ettifouri, Toumi Bouchentouf, and Mohammed Ghaouth Belkasmi. 2020. Artificial neural networks for text-to-sql task: State of the art. In *Advances in Smart Technologies Applications and Case Studies*, pages 557–565, Cham. Springer International Publishing.

Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Comput. Surv.*, 56(2).

Junwoo Park, Youngwoo Cho, Haneol Lee, Jaegul Choo, and Edward Choi. 2021. Knowledge graph-based question answering with electronic health records. In *Machine Learning for Healthcare Conference*, pages 36–53. PMLR.

Mohammadreza Pourreza and Davood Rafiei. 2024. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Advances in Neural Information Processing Systems*, 36.

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.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.

Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce Ho, Carl Yang, and May D Wang. 2024. Ehragent: Code empowers large language models for complex tabular reasoning on electronic health records. *arXiv preprint arXiv:2401.07128*.

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*.

Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2020. RAT-SQL: Relation-aware schema encoding and linking for text-to-SQL parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7567–7578, Online. Association for Computational Linguistics.

Bing Wang, Yan Gao, Zhoujun Li, and Jian-Guang Lou. 2023. Know what I don't know: Handling ambiguous and unknown questions for text-to-SQL. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5701–5714, Toronto, Canada. Association for Computational Linguistics.

Chaojun Wang and Rico Sennrich. 2020. On exposure bias, hallucination and domain shift in neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3544–3552, Online. Association for Computational Linguistics.

Yijun Xiao and William Yang Wang. 2021. On hallucination and predictive uncertainty in conditional language generation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2734–2744, Online. Association for Computational Linguistics.

## A  Reserved Words List

This refers to a list of reserved words in SQL that we used in our experiment.

["SELECT", "AS", "IN", "COUNT", "FROM", "WHERE", "AND", "OR", "INSERT", "UPDATE", "DELETE", "CRE-ATE", "DROP", "ALTER", "JOIN", "ON", "GROUP BY", "ORDER BY", "HAVING", "LIMIT", "UNION", "DIS-TINCT", "INDEX", "TABLE", "VIEW", "TRIGGER", "PRIMARY KEY", "FOREIGN KEY", "NULL", "NOT NULL", "UNIQUE", "CHECK", "DEFAULT", "INDEX", "SEQUENCE", "EXEC", "LIKE", "BETWEEN", "EXISTS", "CASE", "WHEN", "THEN", "ELSE", "END", "CAST", "CHAR", "VARCHAR", "BOOLEAN", "INTEGER", "DATE", "IN-TERVAL", "TIME", "TIMESTAMP", "YEAR", "MONTH", "DAY", "HOUR", "MINUTE", "SECOND", "ZONE", "CURRENT_DATE", "CURRENT_TIME", "CURRENT_TIMESTAMP", "TRUE", "FALSE"]

## B  Input and Output Format

This is the input and output format according to the training specifications of gpt-3.5 turbo.

```
{
    'messages': [
        {'role': 'system', 'content': 'You are 'SQLgpt', an AI
            designed to convert natural language questions into their
            corresponding SQL queries. Your primary goal is to
            accurately generate the exact SQL query for each question
            presented to you.'},

        {'role': 'user', 'content': <Answerable Question>},

        {'role': 'assistant', 'content': <Correct SQL>}
    ]
}
```

# LTRC-IIITH at EHRSQL 2024: Enhancing Reliability of Text-to-SQL Systems through Abstention and Confidence Thresholding

**Jerrin John Thomas, Pruthwik Mishra, Dipti Sharma, Parameswari Krishnamurthy**

LTRC, International Institute of Information Technology, Hyderabad, India

{jerrin.thomas, pruthwik.mishra}@research.iiit.ac.in

{dipti, param.krishna}@iiit.ac.in

## Abstract

In this paper, we present our work in the EHRSQL 2024 shared task which tackles reliable text-to-SQL modeling on Electronic Health Records. Our proposed system tackles the task with three modules - abstention module, text-to-SQL generation module, and reliability module. The abstention module identifies whether the question is answerable given the database schema. If the question is answerable, the text-to-SQL generation module generates the SQL query and associated confidence score. The reliability module has two key components - confidence score thresholding, which rejects generations with confidence below a predefined level, and error filtering, which identifies and excludes SQL queries that result in execution errors. In the official leaderboard for the task, our system ranks 6th. We have also made the source code public[1].

## 1 Introduction

Electronic Health Records (EHRs) have revolutionized healthcare by serving as comprehensive digital repositories of medical histories of patients. They capture every step, from initial admission and diagnosis to treatment plans and discharge summaries. While EHRs are invaluable for clinical data storage and retrieval, unlocking their full potential goes beyond basic searches. Traditional methods often necessitate proficiency in Structured Query Language (SQL), a complex hurdle for many healthcare providers, especially those pressed for time. To bridge this gap and make EHR data more accessible, researchers are exploring the development of question-answering systems that leverage text-to-SQL models. These systems empower users to ask questions in plain natural language and receive answers directly retrieved from the EHR data, streamlining the process of extracting valuable insights from patient data.

In this task (Lee et al., 2024), we tackle the problem of developing a reliable text-to-SQL model tailored for an EHR database, ensuring accurate responses while abstaining from providing incorrect answers. This model must handle a diverse range of topics relevant to clinical settings, such as patient demographics, vital signs, and disease survival rates. The model should accurately generate SQL queries for answerable questions, abstain from providing erroneous answers, and recognize and abstain from addressing unanswerable questions, whether they extend beyond the database schema or are impossible to solve using SQL alone. The spectrum of unanswerable questions also encompasses adversarially crafted queries designed to mislead text-to-SQL models. Successfully tackling this task will yield a robust and scalable question-answering system for EHRs, significantly enhancing how clinicians leverage clinical knowledge.

## 2 Related Work

With the advent of deep learning models, there has been a renewed interest in generating text-to-SQL models in the medical domain. Wang et al. (2020) tackle the problem by developing a deep learning-based approach that adapts a sequence-to-sequence architecture to directly generate SQL queries for a given question. The model further performs the necessary edits using an attentive copying mechanism and task-specific lookup tables. Additionally, they release a large-scale dataset called MIMIC-SQL that generates SQL queries from questions in this domain.

Several papers utilize pre-trained BERT (Devlin et al., 2019) models, as their foundation blocks for their text-to-SQL systems. These models leverage large amounts of text data to learn effective representations of language in their pre-training stage, which are then fine-tuned for the specific task of EHR-based question answering. Pan et al. (2021)

---

[1] https://github.com/jr-john/ehrsql-2024

Figure 1: System Workflow

develop a BERT-based model to convert medical text into an intermediate representation that can then be translated into SQL. Gao et al. (2023) explore using open-source Large Language Models (LLMs) and supervised fine-tuning methods for text-to-SQL tasks. Tarbell et al. (2023) investigate the generalizability of the text-to-SQL models across different EHR systems and data formats. They also introduce a data augmentation approach to improve generalizability.

There has been limited prior work in improving the reliability of text-to-SQL systems, and this task aims to address this gap.

## 3  Dataset

The dataset provided for the shared task is the Medical Information Mart for Intensive Care IV (MIMIC-IV) demo version of EHRSQL with additional unanswerable questions (Lee et al., 2022). MIMIC-IV (Johnson et al., 2023, 2021; Goldberger et al., 2000) is a large database containing de-identified patient information from Beth Israel Deaconess Medical Center. It is a relational database consisting of twenty-six tables, which stores the data collected during routine clinical care, including patient demographics, vital signs, diagnoses, medications, and procedures.

EHRSQL (Lee et al., 2022) is a dataset designed to evaluate and enhance text-to-SQL systems specifically tailored for EHRs. It contains real-world questions gathered from various hospital staff, including doctors, nurses, insurance specialists, and record-keeping personnel. This ensures that the dataset incorporates practical queries healthcare professionals ask daily. The questions

range from simple data lookups to complex calculations, such as determining patient survival rates. Recognizing that not all questions have answers within the EHR data, EHRSQL empowers text-to-SQL systems to identify cases when low confidence or missing information necessitates abstaining from a response.

The given dataset consists of 5,124 samples in the training set, 1,163 samples for validation, and 1,167 samples for testing. The database encompasses eighteen different tables adapted from MIMIC-IV. The training data is exclusively used for model training, while the validation data is not utilized during training. No external data sources are incorporated and no data augmentation techniques are employed.

## 4  System Description

Our system contains three modules - abstention module, text-to-SQL generation module, and reliability module.

### 4.1  Pre-Trained Model Description

We choose SQLCoder-7b-2, a fine-tuned implementation of CodeLlama-7b (Rozière et al., 2024), as the pre-trained text-to-SQL model.[2] This model outperforms GPT-4 (as of Feb 5, 2024) and GPT-4-Turbo (as of Feb 5, 2024) in text-to-SQL benchmarks.

The model is fine-tuned on a comprehensive SQL curriculum, ranging from basic clauses to underrepresented categories like date functions and advanced operations like window functions. It

---

[2]https://huggingface.co/defog/sqlcoder-7b-2

Figure 2: Distribution of confidence scores for test predictions, with scores below -1 excluded

contains hand-crafted SQL queries as well as augmented data from WikiSQL (Zhong et al., 2017). It has thirteen different schemas and questions of varying levels of difficulty. The schemas are quite complex with four to twenty tables. Each question in the dataset has been classified into "easy", "medium", "hard", and "extra-hard" categories. The model is fine-tuned in two stages. First, the base CodeLlama model is fine-tuned on easy and medium questions. Then, the resulting model is fine-tuned on hard and extra-hard questions.

## 4.2 Text-to-SQL Generation Module

The text-to-SQL generation module has a pretrained text-to-SQL model that is fine-tuned on the training data for EHR domain adaptation. We use the prompt template of the pre-trained model for fine-tuning (the prompt template is provided). Following the template, we provide the table metadata as DDL (Data Definition Language) commands for creating the tables. For each table in the database, a "CREATE" statement is generated that includes all the fields and information about the primary key. Each field in the database has a descriptive comment explaining its purpose. The comments at the end of the metadata provide information about the foreign key relationships between the tables. These foreign key constraints define the dependencies between the tables.

The training data is processed and prepared for

fine-tuning the model. We train this model for two epochs till convergence. If the model cannot answer the question with the available database schema, the system returns "null".

---

**Prompt for Text-to-SQL Generation**

### Task
Generate a SQL query to answer [QUESTION]{user_question}[/QUESTION]
### Instructions
If you cannot answer the question with the available database schema, return 'I do not know'
### Database Schema
The query will run on a database with the following schema:
{table_metadata_string}
### Answer
Given the database schema, here is the SQL query that answers [QUESTION]{user_question}[/QUESTION]
[SQL]{answer}[/SQL]

---

### 4.2.1 Prompt Engineering

We experiment with the prompt by providing the table metadata without any comments. The foreign keys in each table are declared using the "FOREIGN KEY" constraint.

699

| Experiment | RS(0) | RS(5) | RS(10) | RS(N) |
|---|---|---|---|---|
| Text-to-SQL Generation (Initial Prompt) | 33.59 | -295.46 | -624.51 | -76766.41 |
| Text-to-SQL Generation | 78.58 | -25.96 | -130.51 | -24321.42 |
| Text-to-SQL Generation + Error Filtering | 82.60 | 13.20 | -56.21 | -16117.40 |
| *Error Filtering* | | | | |
| + Abstention + Text-to-SQL Generation | **83.55** | 25.28 | -32.99 | -13516.45 |
| + Abstention (Multi-Task) + Text-to-SQL Generation (Multi-Task) | 71.47 | -9.94 | -91.35 | -18928.53 |
| *Abstention (Multi-Task) + Text-to-SQL Generation + Error Filtering* | | | | |
| + Confidence Thresholding (Threshold = -1) | 78.66 | 54.24 | 29.82 | -5621.34 |
| + Confidence Thresholding (Threshold = -0.5) | 69.92 | **55.78** | 41.65 | -3230.07 |
| + Confidence Thresholding (Threshold = -0.4) | 66.84 | 55.27 | **43.70** | -2633.16 |
| + Confidence Thresholding (Threshold = -0.35) | 65.38 | 54.24 | 43.10 | -2534.61 |
| + Confidence Thresholding (Threshold = -0.3) | 63.92 | 53.21 | 42.50 | -2436.08 |
| + Confidence Thresholding (Threshold = -0.2) | 58.44 | 51.17 | 43.87 | **-1641.55** |

Table 1: Evaluation Results for Different Experiments (Best Results in Bold)

## 4.3 Abstention Module

**Prompt for Abstention**

```
### Task
Classify whether the question is answerable or unanswerable - [QUESTION]{user_question}[/QUESTION]
### Instructions
- Remember that answerable question is one that can be answered with the given database
- Remember that unanswerable question is one that cannot be answered with the given database
### Database Schema
The query will run on a database with the following schema: {table_metadata_string}
### Answer
Given the database schema, here is the class of [QUESTION]{user_question}[/QUESTION]
[CLASS]{answer}[/CLASS]
```

The abstention module has a pre-trained text-to-SQL model (SQLCoder-7b-2) that is fine-tuned to classify whether a question is answerable given the database schema (the prompt template is provided). We generate the data for fine-tuning the abstention model by taking all the unanswerable questions and randomly sampling the same number of answerable questions, thus preventing class imbalance. The model is fine-tuned for six epochs till convergence. If this model classifies the question as unanswerable, it returns "null".

### 4.3.1 Multi-task Training

We further experiment on the abstention module by training a multi-task model on both text-to-SQL and abstention tasks. Multi-task models are those which are trained to perform multiple related tasks. The training data of both of the tasks are combined and the model is fine-tuned for one epoch till convergence.

## 4.4 Reliability Module

The reliability module has two checks - confidence score thresholding and error filtering. The SQL query is returned if both conditions are met.

### 4.4.1 Confidence Score Thresholding

The confidence score is calculated by summing up the log probabilities of the generated tokens. It checks whether the confidence score of the generated SQL query is above a certain threshold, returning "null" if it does not satisfy this criterion. The confidence score distribution is plotted for the test dataset and the threshold is chosen heuristically.

### 4.4.2 Error Filtering

The generated SQL query is executed on the database and returns "null" if there is an error in execution.

## 5 Experiments

The system is developed incrementally, allowing us to evaluate each module after its introduction. We

begin with a baseline system consisting solely of the text-to-generation module and an experiment is conducted with different representations of the table metadata. Error filtering is introduced through the execution of the query. Next, we integrate the abstention module and compare its performance to the multi-task model trained on both text-to-SQL generation and abstention tasks. Finally, we incorporate the reliability check of confidence score thresholding, experimenting with different threshold values to optimize performance.

## 5.1 Experimental Setup

We perform each fine-tuning using 4-bit Quantized Low-Rank Adaptation (QLoRA) (Dettmers et al., 2023). QLoRA is applied to all the linear layers of the model and the LoRA rank and alpha are chosen as 32 and 64 respectively. A paged 8-bit Adam optimizer with weight decay (Loshchilov and Hutter, 2019) is used with a learning rate of 2.5e-5 on a linear scheduler. We fine-tune the model with a batch size of 8 and 2 gradient accumulation steps. The inference is optimized using vLLM (library for LLM inference and serving) (Kwon et al., 2023). We use greedy decoding for inference and the tokens generated are limited to 4096.

## 5.2 Evaluation Metrics

The system is evaluated using reliability score (RS). RS metric rewards accurate SQL generation for answerable questions and abstaining from answering unanswerable questions. It penalizes incorrect generation or attempts to answer unanswerable questions. The aggregate RS is the mean of individual scores represented as a percentage.

The severity of the penalty can be adjusted by a parameter $c$. A higher value of $c$ leads to stricter evaluation. RS(0) does not penalize any mistakes ($c = 0$). In RS(5), every accurate prediction earns a +1 reward, while each mistake results in a -5 penalty. This means every 5 correct predictions weigh the same as one incorrect prediction.

## 6 Results

The evaluation results of the different experiments can be seen in Table-1. We achieve the best performance with the system of abstention module + text-to-SQL generation module + reliability module, with a confidence score threshold of -0.4 (See Fig 2 for confidence score distribution). Our submission ranks 6th on the leaderboard.

## 6.1 Limitations

The system suffers from the effects of cascading errors. Each module has its own intricacies and potential points of failure. If any of the modules makes an incorrect prediction, the subsequent modules will likely propagate and amplify the error.

The reliability module's performance may heavily depend on the chosen confidence score threshold. Setting the threshold requires careful consideration. A high threshold might reject good queries, while a low threshold might allow unreliable ones.

Despite the reliability checks, there is still a possibility of false positives (accepting unreliable queries) or false negatives (rejecting reliable queries). Balancing between these two extremes is crucial for the system's overall reliability and performance.

## 7 Conclusion

In this work, we present a system consisting of several layers that contribute to reliable text-to-SQL modeling. By filtering out unanswerable questions based on the provided database schema, the system avoids generating incorrect SQL queries and focuses on its strengths. This makes the system robust to unanswerable questions. The reliability module ensures a certain level of confidence in the generated query before using it. Error filtering avoids errors during execution.

Fine-tuning the pre-trained models on EHR data specifically helps the system understand the medical language and schema, leading to more accurate SQL generation for EHR-related queries in comparison to models trained on generic text-to-SQL tasks. Using pre-trained text-to-SQL models as a starting point helps the system leverage existing knowledge and reduces the amount of training data required for fine-tuning. This has led to resource efficiency as no external training data is used.

The modular design allows for easier development, maintenance, and potential future improvements to each specific module. This facilitates adaptation to evolving requirements or changes in the dataset or task. Overall, this system presents a promising approach for reliable text-to-SQL generation in the EHR domain. However, the potential limitations need to be managed to ensure the system's robust and reliable performance in real-world applications.

# References

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

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. *CoRR*, abs/2308.15363.

Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. 2000. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *Circulation*, 101(23):e215–e220.

Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2021. Mimic-iv (version 1.0).

Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2023. Mimic-iv clinical database demo (version 2.2).

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*.

Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. 2022. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601.

Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. 2024. Overview of the ehrsql 2024 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Youcheng Pan, Chenghao Wang, Baotian Hu, Yang Xiang, Xiaolong Wang, Qingcai Chen, Junjie Chen, Jingcheng Du, et al. 2021. A bert-based generation model to transform medical texts to sql queries for electronic medical records: Model development and validation. *JMIR Medical Informatics*, 9(12):e32698.

Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code. *Preprint*, arXiv:2308.12950.

Richard Tarbell, Kim-Kwang Raymond Choo, Glenn Dietrich, and Anthony Rios. 2023. Towards understanding the generalization of medical text-to-SQL models and datasets. *AMIA Annual Symposium Proceedings*, 2023:669–678.

Ping Wang, Tian Shi, and Chandan K. Reddy. 2020. Text-to-sql generation for question answering on electronic medical records. In *Proceedings of The Web Conference 2020*, WWW '20, page 350–361, New York, NY, USA. Association for Computing Machinery.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *CoRR*, abs/1709.00103.

# LTRC-IIITH at MEDIQA-M3G 2024: Medical Visual Question Answering with Vision-Language Models

**Jerrin John Thomas, Sushvin Marimuthu, Parameswari Krishnamurthy**
LTRC, International Institute of Information Technology, Hyderabad, India
{jerrin.thomas, sushvin.marimuthu}@research.iiit.ac.in
{param.krishna}@iiit.ac.in

## Abstract

In this paper, we present our work to the MEDIQA-M3G 2024 shared task, which tackles multilingual and multimodal medical answer generation. Our system consists of a lightweight Vision-and-Language Transformer (ViLT) model which is fine-tuned for the clinical dermatology visual question-answering task. In the official leaderboard for the task, our system ranks 6th. After the challenge, we experiment with training the ViLT model on more data. We also explore the capabilities of large Vision-Language Models (VLMs) such as Gemini and LLaVA.

## 1 Introduction

The rapid evolution of telecommunication technologies, coupled with increased healthcare demands and the recent challenges posed by the pandemic, has accelerated the adoption of remote clinical diagnosis and treatment. Alongside conventional live consultations conducted via telephone or video, asynchronous methods such as e-visits, emails, and messaging chats have emerged as practical and cost-effective alternatives.

This task (wai Yim et al., 2024a) focuses on addressing the challenge of generating suitable textual responses to queries in clinical dermatology, taking into account multimodal inputs such as clinical history, queries, and accompanying images.

This paper describes our proposed solution. We fine-tune ViLT (Kim et al., 2021) for the visual question-answering task with the training data provided for the challenge. We choose ViLT due to its lightweight nature and ability to handle both visual and textual inputs efficiently. After the challenge, we also explore how large Vision-Language Models (VLMs) such as Gemini (Team, 2024) and LLaVA (Liu et al., 2023) perform in this task. These models stand at the top of the multi-modal benchmarks such as Massive Multi-discipline Mul-

timodal Understanding and Reasoning benchmark (MMMU) (Yue et al., 2023).

## 2 Related Work

Previous research has predominantly focused on consumer health question-answering but has been limited to textual inputs (Ben Abacha et al., 2019). Similarly, existing work on visual question-answering has primarily concentrated on radiology images, lacking integration with additional clinical text inputs (Abacha et al., 2019). Moreover, while significant research has been conducted on dermatology image classification, the emphasis has largely been on lesion malignancy classification for dermatoscopy images (Li et al., 2022).

Recently, there has been a surge in the development of multimodal models, particularly large vision-language models (VLMs). These models integrate both textual and visual information, allowing them to understand and generate content that combines both modalities. VLMs typically employ techniques such as joint embedding to unify the representations of text and images in the same embedding space. During training, they utilize datasets that contain interleaved text and images, enabling the model to associate textual descriptions with visual content effectively. This process enables VLMs to grasp nuanced relationships between words and visual elements, facilitating tasks like visual question-answering (VQA). In VQA, these models can accurately respond to questions about images by understanding the content of both the image and accompanying text, showcasing their ability to comprehend and synthesize information across modalities.

For the medical domain, VLMs have been trained on medical corpora and developed for various clinical tasks. MedBLIP (Chen et al., 2023) aids computer-aided diagnosis (CAD) in the medical field. It tackles the challenge of combin-

703

Figure 1: ViLT Model Architecture (from Kim et al. (2021))

ing image and text data from electronic health records for medical diagnosis. The model shows promising results in classifying healthy, mildly impaired, and Alzheimer's patients and also demonstrates the ability to answer medical questions based on visual information. PMC-LLaMA (Wu et al., 2024), designed specifically for medical applications, demonstrates superior performance on medical question-answering tasks. Med-Flamingo (Moor et al., 2023) is a model that can learn from small datasets by embracing in-context learning for the multi-modal medical domain. BiomedGPT (Zhang et al., 2024) is a unified model designed to handle diverse medical data and perform various tasks such as diagnosis and summarization. LLaVA-Med (Li et al., 2023) utilizes a massive dataset of biomedical images and captions from PubMed Central and employs the powerful language model GPT-4 to create diverse training examples.

## 3 Dataset

The dataset provided for the shared task is DermaVQA: A Multilingual Visual Question Answering Dataset for Dermatology (wai Yim et al., 2024b). It is translated and adapted from Chinese telemedicine datasets. The given dataset consists of 842 samples in the training set, 56 for validation and 100 for testing. Each sample has a query, clinical history, and one or more associated images. The textual content is provided in three languages - Chinese, English, and Spanish. The English and Spanish versions of the training set are generated through machine translation from the Chinese original. The validation and test datasets are manually translated by human translators. This approach allows for comprehensive testing and validation of

models across multiple languages while ensuring the integrity and quality of the data through both machine and human translations.

| Dataset | Size |
|---------|------|
| **Train** | 842 |
| **Valid** | 56 |
| **Test** | 100 |

Table 1: Dataset Splits

## 4 System Description

Our system is made of a fine-tuned Vision-and-Language Transformer (ViLT) (Kim et al., 2021) model. ViLT is lightweight and can handle data of both textual and visual modalities.

### 4.1 Model Description

ViLT is a pre-trained multimodal model that simplifies the processing of visual inputs by treating them in the same convolution-free manner as text inputs. This approach reduces the computational complexity of the model-specific components compared to the transformer component for multimodal interactions. ViLT is pre-trained on three objectives: image-text matching, masked language modeling, and word-patch alignment. For the visual question-answering task, we employ a ViLT model with a classifier head on top, which consists of a linear layer applied to the final hidden state of the [CLS] token. This architecture allows the model to leverage its pre-trained multimodal representations to effectively answer questions about visual inputs. The key advantages of ViLT are its simplicity and computational efficiency

Figure 2: After-Challenge System Workflow

## 4.2 Data Pre-Processing

In the pre-processing step, we focus solely on the English content, ignoring the Chinese and Spanish content. We select specific fields from the data, namely `image_ids`, `query_title_en` for questions and `query_content_en` for labels. We proceed to structure the dataset by flattening it and organizing it into tuples containing the image IDs, questions, and labels.

Following the dataset flattening, we encode both the images and texts using the ViltProcessor, a processor tailored for our model. This encoding step is crucial for transforming the raw textual and visual inputs into formats suitable for ingestion by the model. By leveraging the capabilities of the Vilt-Processor, we ensure that the data is prepared as required for the subsequent training process. With the pre-processing complete, we obtain a refined and standardized dataset ready for training our model on visual question-answering tasks.

Initially, we build a dataset with only 200 samples and after the challenge, we use all 842 samples for training. We process the data in batches of 200 samples each and merge all the processed data at the end. This approach allows us to effectively manage memory usage without sacrificing the richness of our dataset, ensuring robust model training and analysis.

## 4.3 Fine-Tuning

Leveraging the ViLT Processor, we seamlessly load our data into the model and the ViLT model is fine-tuned using the processed dataset. During fine-tuning, we tune the hyperparameters to suit our objectives effectively. The batch size is set to 4 and the learning rate at $5e-5$. We train the model for 10 epochs. This fine-tuning process allows the ViLT model to adapt and specialize to the nuances of the specific visual question-answering task, ensuring that it can effectively comprehend and respond to questions of the clinical dermatology domain. With these hyperparameters in place, we aim to achieve optimal performance and robustness in our model's ability to answer questions accurately and comprehensively.

## 4.4 Inference

Following the completion of the fine-tuning, we perform inference on the model. We follow the same pre-processing steps. We utilize `encounter_id` as the question identifier, `query_title_en` as the question itself, and `image_ids` as additional contextual information. Leveraging the model's learned representations and understanding of visual and textual inputs, we generate predictions for each sample.

Once the predictions are obtained, we format the results into a JSON file, organized as an array of JSON objects. Each JSON object contain `encounter_id` as a unique identifier and a responses array. Within this array, we include predicted responses in English (`content_en`), leaving the corresponding fields for Chinese (`content_zh`) and Spanish (`content_es`) empty, as our training and prediction efforts were focused solely on the English language.

## 5 After-Challenge Experiments

After the challenge, we develop a system containing two modules - visual question-answering module and translation module. We experiment with two models - Gemini 1.0 Pro Vision (Team, 2024)

| Models | Chinese | | English | | Spanish | |
|---|---|---|---|---|---|---|
| | DeltaBleu | BERTScore | DeltaBleu | BERTScore | DeltaBleu | BERTScore |
| ViLT (200 Samples) | - | - | 0.46 | 0.83 | - | - |
| ViLT (842 Samples) | - | - | 0.52 | 0.82 | - | - |
| LLaVA - 1.6 34B | 1.60 | 0.64 | 0.53 | 0.82 | 0.88 | 0.76 |
| Gemini 1.0 Pro Vision | 2.70 | 0.59 | 0.86 | 0.70 | 1.39 | 0.66 |

Table 2: Evaluation Results for Different Experiments

and Llava-1.6 34B (Liu et al., 2024).

## 5.1 Model Description

Gemini 1.0 Pro Vision is capable of comprehending inputs from both textual and visual sources, which can be both images or videos, yielding contextually relevant textual outputs. Serving as a foundational model, It excels across a spectrum of multimodal tasks, including visual comprehension, classification, summarization, and content generation from diverse visual inputs such as photographs, documents, infographics, and screenshots.

LLaVA-1.6 34B is a large vision-language model that stands out for its ability to understand and process both text and visual data, making it highly capable in general-purpose visual and language tasks. It is an auto-regressive language model, based on the transformer architecture, and has 34 billion parameters. It is fine-tuned on multi-modal instruction following data. LLaVA-1.6 34B has even surpassed the performance of models like Gemini Pro on some benchmarks.

## 5.2 Visual Question-Answering Module

> **VQA Prompt**
>
> You are a clinical dermatology assistant who can generate clinical responses, given the clinical history and a query, along with one or more associated images. Be concise and do not give additional information other than answering the query.
> {Associated Images}
> {Clinical History}
>
> {Query}

For the visual question-answering (VQA) module, we process the dataset and extract only the encounter_id, English content, and image_ids from the responses. Subsequently, we employ the model to predict results using a prompt created from the clinical history, query, and associated im-

ages (the prompt template is provided). These predictions are then stored in a JSON format, associating each encounter_id with a responses array containing the predicted English data. As for Chinese and Spanish content, we leave those fields empty, reflecting our current focus on English language prediction.

## 5.3 Translation Module

For the translation module, we use Gemini 1.0 Pro. The prompt template only has a simple instruction - "Translate from English to {language}". The English responses are translated into both Spanish and Chinese languages.

## 6 Evaluation Metrics

The evaluation process utilizes deltaBLEU (Galley et al., 2015), a metric that accounts for multiple correct responses. These responses are weighted based on various factors, including completeness, consistency with the most commonly provided answer as determined by human assessment, as well as author rank level and author validation level. The completeness metric assigns a score on a scale of {0.0, 0.5, 1.0}, indicating the extent to which the original query's question was addressed. A score of 1.0 signifies a fully answered query, 0.5 indicates a partial response, and 0.0 reflects a lack of response. If the query doesn't explicitly specify, it's assumed to seek information on both the disease and its treatment. Contains Most Frequent Answer rating is given on a scale of {0.0, 1.0}. A score of 1.0 is assigned if the response aligns with the most frequently provided answer.

## 7 Results

The evaluation results of the different experiments can be seen in Table-2. We submitted the system of ViLT trained on 200 samples for the challenge. Our submission ranks 8th in the leaderboard.

# 8 Conclusion

This work investigates vision-language models like ViLT, Gemini, and LLaVA for the challenging multilingual and multi-modal medical answer generation task in dermatology. A modular system with separate visual QA and translation components shows improved performance over the initial ViLT approach. A key strength is leveraging powerful multi-modal models that can effectively integrate visual and textual clinical data. Future efforts should focus on utilizing larger, more diverse datasets, incorporating stronger multi-modal reasoning, and rigorous evaluations by medical experts to ensure clinical utility and safety before real-world deployment. Overall, this mult-imodal approach holds promise but requires further advancements to be reliable for remote diagnosis and treatment.

# References

Asma Ben Abacha, Sadid A. Hasan, Vivek Datla, Joey Liu, Dina Demner-Fushman, and Henning Müller. 2019. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. In *CEUR Workshop Proceedings*, pages 9–12.

Asma Ben Abacha, Chaitanya Shivade, and Dina Demner-Fushman. 2019. Overview of the MEDIQA 2019 shared task on textual inference, question entailment and question answering. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 370–379, Florence, Italy. Association for Computational Linguistics.

Qiuhui Chen, Xinyue Hu, Zirui Wang, and Yi Hong. 2023. Medblip: Bootstrapping language-image pre-training from 3d medical images and texts. *Preprint*, arXiv:2305.10799.

Michel Galley, Chris Brockett, Alessandro Sordoni, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 445–450, Beijing, China. Association for Computational Linguistics.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 5583–5594. PMLR.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Preprint*, arXiv:2306.00890.

Zhouxiao Li, Konstantin Christoph Koban, Thilo Ludwig Schenck, Riccardo Enzo Giunta, Qingfeng Li, and Yangbai Sun. 2022. Artificial intelligence in dermatology image analysis: Current developments and future trends. *Journal of Clinical Medicine*, 11(22).

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. Llava-next: Improved reasoning, ocr, and world knowledge.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning.

Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Cyril Zakka, Yash Dalmia, Eduardo Pontes Reis, Pranav Rajpurkar, and Jure Leskovec. 2023. Med-flamingo: A multimodal medical few-shot learner. ArXiv:2307.15189.

Gemini Team. 2024. Gemini: A family of highly capable multimodal models. *Preprint*, arXiv:2312.11805.

Wen wai Yim, Asma Ben Abacha, Velvin Fu, Zhaoyi Sun, Fei Xia, Meliha Yetisgen, and Martin Krallinger. 2024a. Overview of the mediqa-m3g 2024 shared task on multilingual and multimodal medical answer generation. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, Mexico City, Mexico. Association for Computational Linguistics.

Wen wai Yim, Velvin Fu, Zhaoyi Sun, Asma Ben Abacha, Meliha Yetisgen, and Fei Xia. 2024b. Dermavqa: A multilingual visual question answering dataset for dermatology. *CoRR*.

Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2024. Pmc-llama: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association*, page ocae045.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *Preprint*, arXiv:2311.16502.

Kai Zhang, Jun Yu, Eashan Adhikarla, Rong Zhou, Zhiling Yan, Yixin Liu, Zhengliang Liu, Lifang He, Brian Davison, Xiang Li, Hui Ren, Sunyang Fu, James Zou, Wei Liu, Jing Huang, Chen Chen, Yuyin Zhou, Tianming Liu, Xun Chen, Yong Chen, Quanzheng Li, Hongfang Liu, and Lichao Sun. 2024. Biomedgpt: A unified and generalist biomedical generative pre-trained transformer for vision, language, and multimodal tasks. *Preprint*, arXiv:2305.17100.

# Author Index