Leveraging Knowledge Graph-Enhanced LLMs for **Context-Aware Medical Consultation**

Anonymous ACL submission

Abstract

Recent advancements in large language models have significantly influenced the field of online medical consultations. However, critical challenges remain, such as the generation of hallucinated information and the integration of upto-date medical knowledge. To address these issues, we propose Informatics Llama (ILlama), a novel framework that combines retrieval-010 augmented generation with a structured medical knowledge graph. ILlama incorporates relevant medical knowledge by transforming subgraphs from a structured medical knowledge graph into text for retrieval-augmented generation. By generating subgraphs from the medical knowledge graph in advance, specifically focusing on diseases and symptoms, ILlama is able to enhance the accuracy and rel-019 evance of its medical reasoning. This framework enables effective incorporation of causal relationships between symptoms and diseases. Also, it delivers context-aware consultations aligned with user queries. Experimental results on the two medical consultation datasets demonstrate that ILlama outperforms the strong baselines, achieving a semantic similarity F1score of 0.884 when compared with ground truth consultation answers. Furthermore, qualitative analysis of ILlama's responses reveals significant improvements in hallucination reduction and clinical usefulness. These results suggest that ILlama has strong potential as a reliable tool for real-world medical consultation environments.¹

011

012

1 Introduction

Traditional online medical consultation platforms, such as HealthCareMagic² and iCliniq³, rely on medical professionals to answer patient queries and provide expert advice. However, due to their dependence on human labor, these systems face limitations in delivering real-time responses, as experts require significant time to review inquiries and generate appropriate answers (Cao et al., 2022). To address this issue, rule-based medical consultation systems have been introduced (Amato et al., 2017; Mishra et al., 2023; Rosruen and Samanchuen, 2018; Huang et al., 2018). Nevertheless, these systems often struggle to handle complex symptoms and patient-specific queries, as they rely on predefined rules that lack flexibility and adaptability.

041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Recently, large language model (LLM)-based consultation systems, such as ChatDoctor (Li et al., 2023), have emerged as promising alternatives. These systems typically extract keywords from user queries to retrieve relevant medical information from sources like Wikipedia or custom disease databases. However, their reliance on potentially inaccurate keyword extraction may lead to search failures and hallucinations, failing to capture essential disease-symptom relationships. While dense embedding-based retrieval methods (Karpukhin et al., 2020) can alleviate keyword extraction errors, they still have limitations in capturing the complex symptom-disease relationships essential for medical consultations. For example, distinguishing whether shortness of breath and fatigue arise from a serious condition like lung cancer or a more benign cause such as anemia requires an understanding of such causal relationships.

To overcome these limitations, we propose a novel framework for real-time medical consultation, called Informatics Llama (ILlama), which improves the response quality as measured by embedding-based evaluation metrics by incorporating structured medical knowledge. To ensure that ILlama performs reliably not only on known data distributions but also in unfamiliar real-world scenarios, we adopt both in-distribution and outof-distribution evaluation protocols throughout this work. This setup allows us to assess the model's generalizability across different sources of med-

¹The code will be released upon publication.

²https://www.askadoctor24x7.com

³https://www.icliniq.com

ical consultations. ILlama leverages retrievalaugmented generation (RAG) (Lewis et al., 2020b) by incorporating medical knowledge from a structured knowledge graph (KG) built upon the unified medical language system (UMLS)⁴. UMLS is updated regularly, which aligns well with common retraining cycles. Therefore, instead of frequently retraining the language model, it is more efficient to update the KG, enabling practical and timely integration of new medical information.

087

097

100

101

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

123

124

125

127

128

129

Unlike keyword-dependent approaches, ILlama improves both retrieval efficiency and reliability. Keyword-extraction methods often misinterpret the intent of user queries and, in many cases, fail to return relevant results if relevant medical information for the extracted keywords is unavailable. By employing a KG-based retrieval approach (Luo et al., 2025), ILlama effectively represents diseasesymptom causal relationships, enhancing the contextual relevance and accuracy of diagnostic responses.

In addition, ILlama tackles two core limitations prevalent in existing dense embedding-based augmentation systems: (1) the incompleteness of external knowledge representations and (2) the difficulty in aligning user queries with the embedded knowledge space (Varshney et al., 2023). ILlama addresses the incompleteness of the KG by constructing subgraphs that enrich sparse regions with semantically related triples. It is also designed to alleviate the challenge of aligning user queries with the KG structure by embedding each triple and integrating it into the answer generation process. This approach enables more accurate semantic matching and enhances the clinical relevance of the generated responses. Specifically, embedding triples allows the model to retrieve more precise symptomdisease associations, reducing factual errors, while the structured knowledge context provided by the KG improves the alignment of responses with realworld clinical reasoning.

We validate the effectiveness of ILlama using two medical consultation datasets with different characteristics. Specifically, we conducted experiments with publicly available data collected from HealthCareMagic and iCliniq. For in-distribution evaluation, we use the HealthCareMagic dataset, which includes separate training, validation, and test splits. For out-of-distribution evaluation, we use real-world consultation records from the iCliniq platform, serving as the test set. ILlama_{7B}, which is based on the same base model as Chat-Doctor, achieves F1 scores of 0.866 and 0.851 on the in-distribution and out-of-distribution datasets, respectively, and outperforms all baseline models. These scores are computed based on the semantic similarity between the generated responses and the ground truth consultation records. ILlama_{8B}, with a more powerful backbone LLM, further improves these results, achieving scores of 0.884 and 0.871, respectively. The reliability of the generated responses is further supported by a qualitative evaluation that assesses their clinical quality.

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

162

163

164

167

169

170

171

172

173

174

175

176

177

178

In summary, our contributions are three-fold:

- We propose ILlama, a framework that reduces hallucinations by integrating structured medical knowledge and explicit disease-symptom causal relationships into LLMs. This framework facilitates easy knowledge updates by allowing replacement of the underlying UMLSbased KG.
- ILlama utilizes subgraphs from a UMLSbased KG, which are transformed into document form, combined with vector search techniques, enabling precise retrieval and integration of medically relevant knowledge into the answer generation process.
- Our framework achieves state-of-the-art performance across multiple datasets, with the best results observed on the HealthcareMagic dataset, significantly improving the reliability and usefulness of automated medical consultation systems.

2 Related Works

2.1 Early Medical Consultation Systems

Early systems (Amato et al., 2017; Mishra et al., 2023; Rosruen and Samanchuen, 2018; Huang et al., 2018) used rule-based approaches for simple Q&A interactions, easing the burden on healthcare professionals but failing to handle complex symptoms and disease interactions. To overcome this, medical-specialized models using natural language processing technologies (Lee et al., 2020; Yuan et al., 2022; Lu et al., 2022) were developed, yet challenges in incorporating structured medical knowledge and understanding causal relationships between symptoms and diseases remain. LLMs

⁴https://www.nlm.nih.gov/research/umls/ archive/archive_home.html

such as GPT-4 (Achiam et al., 2023) have catalyzed the development of models capable of sophisticated medical consultations (Thirunavukarasu et al., 2023; Li et al., 2024; Toma et al., 2023; Chen et al., 2023; Luo et al., 2022; Yang et al., 2024), although persistent challenges remain, including hallucinations and the incorporation of up-to-date medical information (Vaishya et al., 2023; Hadi et al., 2024). To mitigate these issues, we incorporate a UMLS-based KG that enables accurate identification of disease relationships and contextual information retrieval, thereby supporting more clinically relevant and context-aware consultations.

179

180

181

184

185

187

188

190

192

193

194

195

198

199

204

206

207

210

211

212

214

215

216

217

2.2 Knowledge-Based LLMs in Medical Consulting

To address the limitations of LLMs, such as hallucinations, lack of timely medical knowledge, and insufficient adaptability to patient-specific contexts, recent research has explored the integration of real-world knowledge to enhance performance in medical applications. Among these approaches, the combination of LLMs with KGs has demonstrated effectiveness in incorporating external information. For example, KG-enhanced models have been used for diagnosis prediction (Gao et al., 2025), graph-augmented medical dialogues (Varshney et al., 2023), and factual medical question answering (Guo et al., 2022; Martino et al., 2023). However, many of these methods rely on incomplete KGs, which are often limited in coverage or biased toward certain clinical entities. They also struggle with aligning unstructured user queries to structured graph elements, both of which hinder their clinical precision. In contrast, ILlama introduces subgraph-based retrieval and semantic reranking to improve knowledge relevance and integration, offering more accurate and contextsensitive medical consultations.

3 Method

218The proposed framework consists of three main
components: Retriever, Reranker, and Generator.220Medical knowledge from the KG is first segmented
into subgraphs and transformed into documents in
natural language form, which serve as input across
all stages. Section 3.1 describes how the Retriever
identifies subgraphs semantically relevant to the in-
put query. Section 3.2 presents the Reranker, which
employs a cross-encoder (Reimers and Gurevych,
2019) to rerank the retrieved documents in natural

language form. Section 3.3 explains how the Generator uses the top-ranked documents to generate the final response. The overall process is illustrated in Figure 1.

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

251

252

253

254

255

256

257

258

259

260

261

262

263

264

267

268

269

270

271

272

273

3.1 Retriever: Enhancing Medical Knowledge

In medical consultations, it is essential to provide accurate, context-aware information without hallucinations. Our framework requires comprehensive medical knowledge, particularly regarding causal relationships between symptoms and diseases. To achieve this, we incorporate a KG based on UMLS, which enables the language model to effectively capture these relationships and allows targeted retrieval of relevant medical facts from the KG. This ensures that responses are both precise and contextually appropriate to the user's query.

3.1.1 Triple-Centric Knowledge Structuring for Medical Reasoning

To effectively represent medical knowledge within the model, we adopt the Triple2Seq (Bi et al., 2024) method to segment the UMLS-based KG into meaningful and contextually coherent subgraphs. A segmented KG (i.e., subgraph) is a structured representation where nodes denote medical concepts (e.g., diseases, symptoms, and treatments) and edges define the relationships among them (e.g., "has symptom", "treated by", and "has causes").

Each subgraph \mathcal{T}_g is composed of a center triple \mathcal{T}_c (e.g., Lung Cancer-has symptom-Fatigue), representing a core medical concept, and a set of neighboring context triples \mathcal{T}_N (e.g., Lung Cancer-has causes-Smoking) that provide additional medical facts related to the center triple:

$$\mathcal{T}_g = \mathcal{T}_c \cup \mathcal{T}_N. \tag{1}$$

 \mathcal{T}_N includes all triples connected to the center triple via its neighboring nodes in the KG and is defined as:

$$\mathcal{T}_N = \{ \mathcal{T}_i \mid \mathcal{T}_i \in \mathcal{N} \},\tag{2}$$

where \mathcal{N} denotes the set of nodes that are directly linked to the center concept in the graph. For example, if the center triple corresponds to a disease such as lung cancer, the context triples may include related symptoms (e.g., shortness of breath and fatigue), diagnostic procedures (e.g., chest X-ray), or causes (e.g., smoking or air pollution). By organizing knowledge in this localized and relation-centric



Figure 1: Overall architecture for medical query answering using contextualized subgraphs from the UMLS-based KG. These subgraphs are encoded and stored in a vector database, then combined with the user query to generate the final response using the Llama model.

manner, the model is guided to focus on medically 274 relevant and causally connected concepts, thereby enhancing the contextual consistency of the generated responses. Furthermore, this structure enables more accurate and context-aware diagnosis and consultation based on the patient's reported symptoms, ultimately improving the reliability and practicality of medical dialogue systems.

277

281

283

290

295

296

301

305

3.1.2 Subgraph-to-Text Transformation

UMLS-based KG represents relation-The ships between medical concepts using a subject-predicate-object triple structure, which closely resembles the structure of natural language sentences. Leveraging this property, we convert each triple into a natural sentence. This transformation reconstructs the structural relationships in the graph into a coherent narrative, allowing the model to intuitively understand the meaning and connections between medical entities. As a result, the graph-based knowledge is naturally integrated into the text generation process, enabling the model to learn richer contextual information.

We further define subgraphs consisting of a center triple and its related neighboring triples. All triples within each subgraph are converted into natural sentences and aggregated into a single document, forming a semantically coherent and logically structured unit of knowledge. A detailed example of this subgraph-to-document transformation, including the rule-based sentence structure and the resulting document format, is provided in Appendix **B**.

Pseudo Query Generation for 3.1.3 **Fine-Tuning Medical Search System**

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

322

323

324

325

326

327

328

329

330

332

333

334

335

336

340

In our framework, document-form subgraph encoder and reranker models pre-trained on general domain data are not sufficient to accurately retrieve medical information grounded in a UMLSbased KG. To improve their ability to understand and retrieve UMLS-specific representations, these models should be fine-tuned on domain-specific data. However, manually constructing high-quality query-document pairs is impractical and costly. To address this, we propose an automated pipeline based on frozen Llama3.1_{8B} (Dubey et al., 2024) models that generates and filters training data without human supervision. The pipeline consists of two core components: a pseudo query generator, which produces queries reflecting key contents of each document, and an evaluator, which filters these queries based on two criteria, patient centeredness ([Patient/notPatient]) and document relevance ([Relevant/Irrelevant]).

As illustrated in Figure 2, the system generates multiple candidate queries per document, evaluates them, and filters those that meet the training standards. Although the evaluator operates in a zero-shot setting without parameter updates, it consistently selects high-quality query-document pairs and generalizes well across unseen pairs. These filtered pairs are subsequently used to fine-tune the document-form subgraph encoder and reranker models, contributing to improved retrieval accuracy and consistency. Details on the objective functions used for each model are provided in Appendix C. Furthermore, while our pipeline focuses on medical consultation documents in this study, it can be



Figure 2: Overview of training data generation process for the document-form subgraph encoder and reranker. A fixed query generator creates questions that a patient is likely to ask, and an evaluator checks if they match a patient-like style and are relevant to the document. If they don't meet the criteria, they are regenerated. The final data trains *bge-large-en v1.5* and *bge-reranker-large*, enhancing the model's ability to understand and process patient-oriented queries.

easily adapted to other domains by adjusting the evaluation criteria.

341

342

343

345

347

351

353

357

361

363

370

3.1.4 Document Embedding and Vector Retrieval

We fine-tuned the *bge-large-en-v1.5* (Xiao et al., 2024) model to generate embeddings for documents derived from the subgraph, optimizing its ability to capture semantic nuances. These embeddings are stored in a vector database using FAISS (Johnson et al., 2019), which is optimized for large-scale similarity searches. By integrating maximum inner product search (Shrivastava and Li, 2014), we efficiently retrieve relevant documents, ensuring low-latency and high-precision results, crucial for real-time applications like conversational agents.

3.2 Reranker: Filtering for Exact Knowledge

To enhance the accuracy of retrieved documents, we implemented a reranking process using the fine-tuned *bge-reranker-large*⁵ model. The crossencoder simultaneously encodes both the user's query and the documents, effectively capturing complex interactions and evaluating the relevance and specificity of each document in relation to the user's needs. Based on the reranking scores, documents are reordered to prioritize those most relevant to the query and tailored to the user's context. This ensures the model can reflect more precise and relevant information in the final response, ultimately providing answers that are accurate and aligned with the user's reported medical concerns.

bge-reranker-large

3.3 Generator: Generating Patient-Centered Medical Consultation

In the final stage, we generate medically accurate and context-aware responses using reranked documents. We fine-tune Llama2_{7B} (Touvron et al., 2023) and Llama3.1_{8B} on real medical consultation data, allowing the model to learn associations between patient queries and retrieved knowledge. Unlike methods that rely solely on synthetic prompts, our framework uses actual consultation records with retrieved documents integrated during finetuning. This helps the model better understand semantic relationships between queries and supporting knowledge, grounding its generation in clinically relevant context. As a result, ILlama can deliver more accurate and tailored responses while reducing hallucinations and speculative content.

4 Experiments

	HealthcareMagic	iCliniq
# dialogues	112,165	1,380
# tokens	27,475,545	313,735
Avg. # tokens per dialogue	245.01	227.34
Max # tokens per dialogue	2,544	1,001
Min # tokens per dialogue	78	60

Table 1: Statistics of the datasets used for training, validation, and testing, showing the distribution of dialogues and token counts.

4.1 Datasets

We use two types of data in ILlama, namely a UMLS-based KG and real-world medical consulta-

391

371

372

373

375

376

377

378

380

381

382

384

386

387

⁵https://huggingface.co/BAAI/

Category		Ι	n-Distributi	on		Οι	ıt-of-Distrib	ution	
Model	F1	METEOR	BLEU-4	ROUGE-2	PPL↓ F1	METEOR	BLEU-4	ROUGE-2	PPL↓
			Ba	selines witho	ut Retrieval				
BART _{Large}	0.837	0.059	0.0	0.038	1.398 0.83	8 0.063	0.0	0.023	1.404
T5 _{Large}	0.840	0.061	0.0	0.031	1.475 0.84	3 0.069	0.0	0.020	1.417
Llama2 _{7B} w/ LoRA	0.838	0.192	0.031	0.052	7.940 0.83	8 0.201	0.029	0.050	10.926
Llama3.1 _{8B} w/ LoRA	0.844	0.230	0.061	0.074	10.659 0.84	1 0.222	0.029	0.048	14.259
			Base	elines without	Fine-Tuning				
Gemma29B	0.836	0.180	0.016	0.036	19.211 0.84	1 0.201	0.022	0.040	16.648
Yi1.59B	0.832	0.168	0.015	0.034	15.552 0.83	5 0.188	0.021	0.038	14.371
Falcon37B	0.839	0.135	0.008	0.025	39.389 0.84	4 0.156	0.012	0.028	34.445
DeepSeek-R1 _{8B}	0.832	0.175	0.012	0.028	43.653 0.83	7 0.191	0.014	0.030	40.259
Baselines with Fine-Tuning & Retrieval									
Llama27B w/ LoRA	0.837	0.191	0.029	0.050	7.939 0.83	9 0.203	0.029	0.050	10.926
Llama3.18B w/ LoRA	0.786	0.222	0.010	0.024	10.657 0.78	9 0.199	0.006	0.019	14.259
ChatDoctor	0.846	0.218	0.008	0.022	10.009 0.84	5 0.211	0.035	0.045	12.676
Ours									
ILlama _{7B}	0.866	0.203	0.037	0.058	7.939 0.85	1 0.213	0.041	0.048	10.924
ILlama _{8B}	0.884	0.231	0.063	0.075	7.659 0.87	1 0.222	0.030	0.049	10.259

Table 2: Performance comparison across baselines categorized into three groups: without retrieval, without finetuning, and with fine-tuning & retrieval. Metrics such as F1, METEOR, BLEU, ROUGE, and PPL are evaluated for both in-distribution and out-of-distribution datasets. The highlighted row represents our proposed method, demonstrating superior performance across most metrics.

tion records. The KG provides structured clinical relationships that support precise retrieval, while the consultation data enables response generation grounded in authentic patient–doctor interactions. Detailed descriptions of each dataset are provided in the following subsections.

4.1.1 Datasets for Medical Retrieval

395

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

We construct our KG using the 2024 release of the UMLS Metathesaurus⁶, which comprises approximately 20K entities, 22 relation types, and 250K triples. This structured resource provides a semantic backbone for our RAG framework, enabling precise retrieval and integration of clinically relevant knowledge. Grounding generation in this KG enhances factual accuracy, reduces hallucinations, and supports context-aware medical responses.

4.1.2 Datasets for Medical Consultation

To evaluate the performance of ILlama, we used medical consultation records from two real-world platforms: HealthcareMagic and iCliniq. The HealthcareMagic dataset, specifically collected for medical question answering tasks, consists of single-turn interactions where each patient query is paired with a response from a licensed medical professional. We split this dataset into training, validation, and in-distribution test sets using an 8:1:1

> ⁶https://www.nlm.nih.gov/research/umls/ licensedcontent/umlsknowledgesources.html

ratio. In contrast, the iCliniq dataset, which follows a similar single-turn format, was used exclusively as an out-of-distribution test set. This separation allows us to evaluate the model's generalization performance on unseen queries from a different source, minimizing the risk of data leakage and ensuring a fair comparison. Both datasets are publicly available for academic research and have been deidentified to protect user privacy. As these records often include patient-reported details such as age and symptoms, the model implicitly learns to adapt responses to clinical contexts during fine-tuning. Detailed dataset statistics are provided in Table 1. 418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

4.2 Metrics

We evaluated our model using semantic and quantitative metrics to assess the accuracy and contextual appropriateness of generated responses. For semantic evaluation, we adopted BERTScore (Zhang* et al., 2020) with RoBERTa_{Large} (Liu et al., 2019), which measure contextual similarity using deep embeddings. This approach is particularly suitable for handling the nuances of medical language where lexical overlap is often limited (Hanna and Bojar, 2021). We report the F1 score from BERTScore as our primary semantic similarity metric.

In addition to semantic evaluation, we also employed lexical metrics such as ROUGE-2 (Lin, 2004), BLEU-4 (Papineni et al., 2002), Perplexity



Figure 3: Comparison of responses from ChatDoctor, ILlama_{7B}, and ILlama_{8B} regarding COVID-19 symptoms. Highlighted sections indicate usefulness, ambiguity, hallucinations, and grammatical errors.

(PPL) (Lavie and Agarwal, 2007), and METEOR to evaluate lexical accuracy, fluency, and coherence. In particular, METEOR considers synonymy, stemming, and paraphrasing, offering a flexible assessment of sentence-level similarity. This evaluation framework enables a comprehensive assessment of the model's ability to deliver precise, fluent, and contextually relevant medical responses.

4.3 Baselines

Baselines without Retrieval These models rely solely on fine-tuned capabilities on domain-specific data, without any retrieval mechanism. We examine BART_{Large} (Lewis et al., 2020a), $T5_{Large}$ (Raffel et al., 2020), and Llama2_{7B} w/ LoRA (Hu et al., 2022) and Llama3.1_{8B} w/ LoRA.

Baselines without Fine-Tuning Models in this category use a retrieval mechanism but are not fine-tuned on domain-specific data. These baselines enhance their performance by leveraging the PubMed dataset (Xiong et al., 2024) for retrieval of pertinent biomedical literature, which provides a rich source of domain-specific information without the need for additional fine-tuning. These include Gemma2_{9B} (Team et al., 2024), Yi1.5_{9B} (Young et al., 2024), Falcon3_{7B} (Team, 2024), and DeepSeek-R1_{8B} (DeepSeek-AI et al., 2025).

472 Baselines with Fine-Tuning & Retrieval This
473 category includes models that undergo fine-tuning
474 on domain-specific data and use retrieval. Mod475 els include Llama2_{7B} w/ LoRA (Hu et al., 2022),
476 Llama3.1_{8B} w/ LoRA, and ChatDoctor, although
477 ChatDoctor does not use PubMed for retrieval.

4.4 Result

479 As shown in Table 2, ILlama consistently outper-480 formed the baselines across both in-distribution and out-of-distribution evaluations. In the indistribution setting, ILlama_{8B} achieved the best F1 score of 0.884 and METEOR of 0.231, surpassing all baseline models. Notably, ILlama_{7B} also showed strong performance (F1: 0.866, METEOR: 0.203), outperforming ChatDoctor, which shares the same base model, across all major metrics including F1, METEOR, and PPL. These results highlight the effectiveness of ILlama's RAG framework in producing accurate and coherent responses.

In the out-of-distribution setting, ILlama maintained robust generalization performance. ILlama_{8B} achieved an F1 of 0.871 and METEOR of 0.222, with minimal performance drop compared to its in-distribution results. This demonstrates ILlama's ability to adapt to unseen queries and linguistic variations from different data sources. The integration of embedding-based vector retrieval and a structured medical KG played a key role in improving factual consistency while minimizing hallucinations. Overall, ILlama achieved state-of-theart performance across F1, METEOR, and PPL, validating its reliability and generalization in real-world medical consultation scenarios.

5 Analysis

5.1 Qualitative Analysis of Outputs

As shown in Figure 3, we present a qualitative comparison of ChatDoctor, ILlama_{7B}, and ILlama_{8B} in response to a COVID-19 related query, alongside the underlying reasoning represented through contextualized subgraphs extracted from the UMLSbased KG. While ChatDoctor exhibited frequent hallucinations, such as recommending antibiotics for viral infections, ILlama_{7B} demonstrated improved clinical reasoning but still included unnecessary suggestions. ILlama_{8B} provided the most balanced response, delivering accurate medical guid-



Figure 4: ILlama and ChatDoctor performance comparison in search latency, answer latency, end-to-end latency and throughput.



Figure 5: F1 score comparison between subgraph and full graph reasoning

ance and appropriate follow-up steps. These evaluations were conducted using the OpenAI o1 model⁷ (Liu et al., 2023). The prompts used for this assessment are provided in Appendix D. The reasoning process is grounded in causal and diagnostic relationships (e.g., Viral Infection–URI–Cough/Fever or Chest X-Ray to rule out Lung Infection), captured within the subgraph structure.

518

519

520

521

524

525

526

534

535

538

539

540

541

542

5.2 Latency Analysis in Real-Time Medical Consultation Systems

In Figure 4, we present a comparison of latency and throughput between ILlama and ChatDoctor. ILlama consistently demonstrates lower latency across search, answer, and end-to-end processing. For example, ILlama's end-to-end latency ranges from approximately 5,538 to 6,507 ms, whereas ChatDoctor's ranges from around 6,921 to over 26,491 ms. This gap stems from ChatDoctor's reliance on LLM-based keyword extraction followed by live API calls to external sources such as Wikipedia, whose articles average more than 900 tokens and thus markedly inflate the answer-latency portion of the overall response time. In contrast, ILlama uses pre-indexed graph-based retrieval with documents averaging around 130 tokens, enabling higher throughput with reduced delay. These results highlight ILlama's efficiency and its suitability for real-time medical consultation.

543

544

545

546

547

549

550

551

552

553

554

555

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

5.3 Full-Graph vs. Subgraph Search

We compared the effectiveness of subgraph-based retrieval using Triple2Seq with that of a full-graph approach. As shown in Figure 5, the subgraph method produces more accurate and reliable responses. Full-graph retrieval, while comprehensive, often includes loosely connected or clinically irrelevant nodes, which can overwhelm the generation process with extraneous information and increase the risk of hallucinations. In contrast, subgraph retrieval narrows the focus to a central medical concept and its semantically related neighbors, allowing the model to process only the most relevant context. This targeted representation helps align retrieved knowledge more precisely with the user's query, resulting in improved factual accuracy and reduced noise. Such precision is particularly important in the medical domain, where irrelevant or overly broad context can compromise the safety and trustworthiness of responses.

6 Conclusion

In this study, we proposed ILlama, a retrievalaugmented medical consultation framework that integrates structured KGs and cross-encoder reranking. ILlama captures causal relationships between symptoms and diseases, and leverages patient demographics during training to support advice while reducing hallucinations. Experiments show strong performance across accuracy, latency, and contextual relevance in both in-distribution and out-ofdistribution settings, highlighting the promise of structured medical knowledge in LLMs for scalable healthcare applications.

⁷https://openai.com/o1/

579 Limitations

580 Although ILlama improves the accuracy of medical consultations and reduces hallucinations by leveraging a UMLS-based KG and embedding-based 582 retrieval, several limitations remain. First, the cov-583 erage of the KG and datasets is relatively narrow, 584 primarily reflecting specific diseases and linguistic patterns. This limits the model's generalizability to broader clinical domains and multilingual, multicultural contexts. Expanding the training data to include more diverse and representative medical 589 590 cases is necessary to improve robustness. Second, the KG may not capture the latest clinical updates such as emerging diseases, new treatments, or revised guidelines. Without real-time synchroniza-593 tion, the model may generate outdated or clinically 594 irrelevant responses. Third, while ILlama achieves 595 strong performance on standard metrics such as 596 F1, METEOR, and PPL, these metrics do not fully capture clinical safety or decision-making validity. Future work should explore automated evaluation methods or deployment in simulated clinical environments to further validate the model's clinical reliability without increasing the burden on human experts. Lastly, the current system relies on relatively large-scale models, which may limit deployment in resource-constrained settings. Future directions include developing lightweight variants 606 and adapting the framework for multilingual and 607 cross-cultural applications to enable broader adoption in global healthcare environments.

Ethical Considerations

610

611

612

613

614

615

616

617

618

624

625

628

While our model, ILlama, aims to enhance medical consultations by reducing hallucinations and incorporating up-to-date medical knowledge, it is not 100% accurate and may not always provide correct diagnoses. In the medical field, inaccuracies can lead to severe consequences, including misdiagnosis and inappropriate treatment recommendations, which could potentially be life-threatening.

Therefore, it is crucial to acknowledge the limitations of AI-based medical consultation systems. Further research is necessary to improve the model's accuracy and reliability. Additionally, such systems should be used to support, not replace, professional medical advice. We recommend that users consult qualified healthcare professionals for personalized medical guidance and that our model be used as a supplementary tool to enhance access to medical information.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*. 629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

- Flora Amato, Stefano Marrone, Vincenzo Moscato, Gabriele Piantadosi, Antonio Picariello, Carlo Sansone, and 1 others. 2017. Chatbots meet ehealth: Automatizing healthcare. In *WAIAH@ AI* IA*, pages 40–49.
- Zhen Bi, Siyuan Cheng, Jing Chen, Xiaozhuan Liang, Feiyu Xiong, and Ningyu Zhang. 2024. Relphormer: Relational graph transformer for knowledge graph representations. *Neurocomputing*, 566:127044.
- Bolin Cao, Wensen Huang, Naipeng Chao, Guang Yang, and Ningzheng Luo. 2022. Patient activeness during online medical consultation in china: multilevel analysis. *Journal of Medical Internet Research*, 24(5):e35557.
- Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, and 1 others. 2023. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yanjun Gao, Ruizhe Li, Emma Croxford, John Caskey, Brian W Patterson, Matthew Churpek, Timothy Miller, Dmitriy Dligach, and Majid Afshar. 2025. Leveraging medical knowledge graphs into large language models for diagnosis prediction: Design and application study. *JMIR AI*, 4:e58670.
- Quan Guo, Shuai Cao, and Zhang Yi. 2022. A medical question answering system using large language models and knowledge graphs. *International Journal of Intelligent Systems*, 37(11):8548–8564.
- Muhammad Usman Hadi, Qasem Al Tashi, Abbas Shah, Rizwan Qureshi, Amgad Muneer, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, and 1 others. 2024. Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. *Authorea Preprints*.

- Michael Hanna and Ondřej Bojar. 2021. A fine-grained analysis of bertscore. In *Proceedings of the Sixth Conference on Machine Translation*, pages 507–517.
- 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.
- Chin-Yuan Huang, Ming-Chin Yang, Chin-Yu Huang, Yu-Jui Chen, Meng-Lin Wu, and Kai-Wen Chen.
 2018. A chatbot-supported smart wireless interactive healthcare system for weight control and health promotion. In 2018 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM), pages 1791–1795.

695

700

701

702

704

705 706

707

710

711

712

713

714 715

716

717

718

719

721

722

723

724

726

727

730

731

732

733

734 735

736

737

740

- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for opendomain question answering. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6769–6781, Online. Association for Computational Linguistics.
- Alon Lavie and Abhaya Agarwal. 2007. Meteor: an automatic metric for mt evaluation with high levels of correlation with human judgments. In *Proceedings* of the Second Workshop on Statistical Machine Translation, StatMT '07, page 228–231, USA. Association for Computational Linguistics.
- 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. 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, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020b. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Jianning Li, Amin Dada, Behrus Puladi, Jens Kleesiek, and Jan Egger. 2024. Chatgpt in healthcare: a taxonomy and systematic review. *Computer Methods and Programs in Biomedicine*, page 108013.

Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6). 741

742

743

744

745

746

747

748

749

750

751

752

754

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

778

779

781

782

783

784

785

786

787

788

789

790

791

792

793

- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 2511–2522, Singapore. Association for Computational Linguistics.
- 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.
- 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.
- Hang Luo, Jian Zhang, and Chujun Li. 2025. Causal graphs meet thoughts: Enhancing complex reasoning in graph-augmented llms. *arXiv preprint arXiv:2501.14892.*
- 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.
- Ariana Martino, Michael Iannelli, and Coleen Truong. 2023. Knowledge injection to counter large language model (llm) hallucination. In *European Semantic Web Conference*, pages 182–185. Springer.
- Ritwik Mishra, Simranjeet Singh, Jasmeet Kaur, Pushpendra Singh, and Rajiv Shah. 2023. Hindi chatbot for supporting maternal and child health related queries in rural India. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 69–77, Toronto, Canada. Association for Computational Linguistics.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting on Association for Computational Linguistics, ACL '02, page 311–318, USA. Association for Computational Linguistics.

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.

795

796

798

812

814

818 819

821

822

823

824

825

829

831

832

833 834

835

838

841

842

844 845

- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. 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.
 - Nudtaporn Rosruen and Taweesak Samanchuen. 2018. Chatbot utilization for medical consultant system. In 2018 3rd Technology Innovation Management and Engineering Science International Conference (TIMES-iCON), pages 1–5.
 - Anshumali Shrivastava and Ping Li. 2014. Asymmetric lsh (alsh) for sublinear time maximum inner product search (mips). *Advances in neural information processing systems*, 27.
- Falcon-LLM Team. 2024. The falcon 3 family of open models.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- 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.
- 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. *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, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Raju Vaishya, Anoop Misra, and Abhishek Vaish. 2023. Chatgpt: Is this version good for healthcare and research? *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, 17(4):102744.
- Deeksha Varshney, Aizan Zafar, Niranshu Kumar Behera, and Asif Ekbal. 2023. Knowledge grounded medical dialogue generation using augmented graphs. *Scientific Reports*, 13(1):3310.

Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. C-pack: Packed resources for general chinese embeddings. In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24, page 641–649, New York, NY, USA. Association for Computing Machinery. 850

851

852

853

854

855

856

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

- Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6233–6251, Bangkok, Thailand. Association for Computational Linguistics.
- Kailai Yang, Tianlin Zhang, Ziyan Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. Mentallama: interpretable mental health analysis on social media with large language models. In *Proceedings* of the ACM on Web Conference 2024, pages 4489– 4500.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, and 1 others. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*.
- 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 *International Conference on Learning Representations*.

A Implementation Details

896

900

901

902

903

904

905

906

907

908

909

910

911

912

913 914

915

916

917

918

919

920

921

922

926

928

932

This study fine-tuned a medical domain-specific model using LoRA with a configuration of r = 16, $lora_alpha = 16$, and $lora_dropout = 0$. The learning rate started at 2×10^{-5} , with 10% of the total steps dedicated to warmup. A linear scheduler was used for adjusting the learning rate during training. The model was trained for 3 epochs, and the maximum sequence length was set to 4,096 for handling complex queries, and the training was conducted on two NVIDIA RTX A6000 48GB GPUs.

For retrieval, 50 documents were fetched using maximum inner product search from the FAISS vector store. The top 10 documents from this set were selected for final use after reranking. This approach improved the model's ability to address medical queries by leveraging dense retrieval methods, enhancing both retrieval accuracy and response quality.

B Example Document from UMLS Subgraph

Table 3 presents an example of subgraph-to-text conversion used in our system. The subgraph is constructed around the central triple (*Lung cancer – has symptom – fatigue*) from the UMLS-based KG. All triples connected to the central node are included and expressed as simple natural language sentences using a rule-based template. Each relation type (e.g., *has symptom, diagnosed by, treated by*) is mapped to a consistent sentence pattern, such as "*X has symptom Y*" or "*X can be diagnosed by Y*." This consistency facilitates automatic transformation and retrieval in downstream components. The resulting document serves as a structured and semantically coherent unit of medical knowledge for training and inference.

C Objective Functions for Fine-Tuning the Medical Search System

C.1 Document-form Subgraph Encoder

The encoder is trained with the InfoNCE loss (Oord et al., 2018), which is a contrastive learning objective widely used in self-supervised learning. Given a set of N random samples $X = \{x_1, \ldots, x_N\}$ containing one positive sample x_{t+k} from the true conditional distribution $p(x_{t+k} | c_t)$ and N-1 negative samples drawn from a proposal distribution $p(x_{t+k})$, the loss is formulated as:

$$\mathcal{L}_N = -\mathbb{E}_X \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum\limits_{x_j \in X} f_k(x_j, c_t)} \right],$$

where $f_k(x, c_t)$ denotes a scoring function (e.g., a dot product or similarity function) that estimates the compatibility between context c_t and future sample x. Optimizing this loss leads $f_k(x_{t+k}, c_t)$ to approximate the density ratio:

933

934

935

936

937

938

939

940

941

942

943

944

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

$$f_k(x_{t+k}, c_t) \propto \frac{p(x_{t+k} \mid c_t)}{p(x_{t+k})}.$$

C.2 Reranker

The reranker model adopts a cross-encoder architecture and is fine-tuned with a binary crossentropy loss. Given a query-document pair (q, d)and a binary label $y \in \{0, 1\}$ indicating relevance, the model predicts a scalar relevance score $\hat{y} = \text{sigmoid}(s(q, d))$, and the loss is computed as:

$$\mathcal{L}_{BCE} = -\left[y \log \hat{y} + (1 - y) \log(1 - \hat{y})\right]$$
945

This objective encourages the model to produce high scores for relevant documents and low scores for irrelevant ones, improving the quality of the final ranking.

D Evaluation Prompt Design

To support the qualitative evaluation of model outputs, we designed three structured prompts targeting hallucination detection, grammatical correctness, and patient helpfulness. These prompts were used with the OpenAI o1 model to evaluate responses generated by ILlama. Table 5 presents the full text of each prompt. Each includes clear task instructions, placeholders for the model-generated response, and, in the helpfulness case, the original patient question. The prompts instruct the model to make a binary decision and identify specific parts of the response when relevant.

The hallucination prompt assesses whether a response contains fabricated or unsupported information. The grammatical prompt checks for language correctness. The helpfulness prompt determines whether the response includes content that would be useful to a patient, based on the given question. Evaluations were conducted in a zero-shot setting, and the prompt design aimed to guide the model toward accurate and consistent judgments without fine-tuning. This allowed for scalable and focusedassessment of clinical response quality.

E Algorithm

974

975

976

977

978

979

980

981

982

E.1 ILlama Algorithm

This algorithm, as shown in Algorithm 1, retrieves and reranks relevant documents for context-aware medical consultations. It combines FAISS search results, reranks them with a cross encoder, and generates a contextually accurate response using Llama, maintaining optimal performance and accuracy throughout the process.

E.2 Pseudo Query Generation Algorithm

This algorithm, as shown in Algorithm 2, generates patient-style queries and evaluates them to obtain 985 (q, d) pairs for training the encoder and reranker if the conditions are met. Based on the input prompt 987 and documents derived from the KG, the pseudo 988 query generator (Llama3.18B) creates a query. The 989 evaluator (Llama3.18B) then checks if the generated query meets the "patient-style" and "relevant" conditions. If the conditions are satisfied, the (q, d)992 pairs are stored for document-form subgraph en-993 994 coder and reranker training; otherwise, the query is regenerated, and the evaluation process is repeated. 995

Subject	Relation	Object	Document-form subgraph
Lung Cancer	has symptom	Fatigue	Lung cancer has symptom fatigue.
Lung Cancer	has symptom	Shortness of Breath	Lung cancer has symptom shortness of breath.
Shortness of Breath	is symptom of	Anemia	Shortness of breath is symptom of anemia.
Fatigue	is symptom of	Anemia	Fatigue is symptom of anemia.
Lung Cancer	has symptom	Chronic Cough	Lung cancer has symptom chronic cough.
Lung Cancer	diagnosed by	Chest X-Ray	Lung cancer can be diagnosed by chest X-ray.
Lung Cancer	has cause	Smoking	Lung cancer has cause smoking.
Lung Cancer	has cause	Air Pollution	Lung cancer has cause air pollution.
Lung Cancer	treated by	Surgery	Lung cancer is treated by surgery.
Surgery	isa	lobectomy	Surgery is a lobectomy.

Table 3: Example of subgraph-to-text conversion for a document centered on the triple (*Lung cancer – has symptom – fatigue*).

Prompting Category	Input Prompt
ILlama's prompt	You are a medical assistant specializing in providing expert consultations
	for medical inquiries. Your role is to deliver accurate, user-friendly
	medical information, clarify symptoms, explain potential
	medical conditions, and recommend next steps with empathy
	and professionalism. When formulating your response,
	to ensure clarity and accuracy, user-friendly answer in your response.
	<pre>### Context {context}</pre>
	### Input {query}
	### Response

Table 4: Prompt used for ILlama inference

Algorithm 1 ILlama Algorithm for Medical Query Answering

1: Input:

- q: User query
- KG: UMLS-based KG
- *DB*: FAISS vector database (Encoded Subgraph Documents)
- T2S: Triple2Seq
- QE: Query Encoder
- CE: Cross Encoder
- Llama: Llama Model
- 2: **Output:** Final response r
- 3: Step 1: UMLS-based KG Processing
- 4: $KG_{sub} \leftarrow T2S.split(KG)$
- 5: $D_{sub} \leftarrow \text{Convert } G_{sub}$ to text-based documents
- 6: Store D_{sub} in FAISS Vector Database
- 7: Step 2: Query Encoding
- 8: $q_{emb} \leftarrow QE.encode(q)$
- 9: Step 3: Retrieval & Reranking from Vector DB
- 10: $D_{top50} \leftarrow DB.retrieve(q_{emb}, k = 50)$
- 11: for each document d in D_{top50} do
- 12: $s_d \leftarrow CE.score(q, d)$
- 13: end for
- 14: $D_{top10} \leftarrow$ Select top-10 documents based on s_d
- 15: Step 4: Response Generation
- 16: $input \leftarrow q + D_{top10}$
- 17: $r \leftarrow Llama.generate(input)$
- 18: **Return** r

Algorithm 2 Patient-Style Pseudo Query Generation and Evaluation

1: **Input:**

- *p*: Prompt for query generation
- *d*: Graph Document (Derived from KG)
- QG: Query Generator (Llama3.18B)
- *Eval*: Evaluator (Llama3.1_{8B})
- 2: **Output:** (q, d) pairs for training Encoder and Reranker
- 3: Step 1: Generate Query

```
4: q \leftarrow QG.generate(p, d)
```

- 5: Step 2: Evaluate Query
- 6: $(s_1, s_2) \leftarrow Eval.check(q, d)$
- 7: if $s_1 ==$ Patient-Style and $s_2 ==$ Relevant then
- 8: **Store** (q, d) for training Encoder and Reranker

9: **else**

- 10: **Regenerate** q using QG
- 11: **Repeat from Step 2**
- 12: **end if**

13: **Return** (q, d) pairs

Prompting Category	Input Prompt
Hallucination Evaluation	The following is a response generated by a model. Carefully read the response and evaluate whether it contains hallucinations based on logical consistency and factual accuracy. A hallucination refers to information that is fabricated or unsupported by evidence.
	### InstructionsIf a hallucination is found, pinpoint the exact part.If no hallucination is found, respond with "No hallucination."
	<pre>### Model Response: {model response}</pre>
Grammatical Error Evaluation	### Evaluation: The following is a response generated by a model. Carefully read the response and identify any grammatical errors.
	### Instructions- If grammatical errors are found, pinpoint the exact part.- If no grammatical errors are found, respond with"No grammatical errors."
	### Model Response: {model response}
	### Evaluation:
Helpful Information for Patients Evaluation	The following is a patient's question and a response generated by a model. Carefully read the response and identify any words or phrases that could be helpful to the patient.
	 ### Instructions Pinpoint the exact words or phrases in the model's response that are relevant to the patient's question. If no helpful information is found, respond with "No helpful information."
	<pre>### Patient's Question: {question}</pre>
	<pre>### Model Response: {model response}</pre>
	### Evaluation:

Table 5: Prompt for evaluation ILlama