

Scalable Multilingual Medical Intelligence System Using LLMs, Knowledge Graphs, and Global Clinical Standards

Veerababu Reddy^{1*}, Maneesha Shaik²,
Mahalakshmi Chadalahwada³, Geetha Varshini Palla⁴,
Gnanesh Surangi⁵, Venkata Sai Teja Vadalasetty⁶

^{1*}Department of IT, Vignan's Lara Institute of Technology and Science,
Vadlamudi, Chebrolu, 522213, Andhra Pradesh, India.

^{2*}Department of IT, Vignan's Lara Institute of Technology and Science,
Vadlamudi, Chebrolu, 522213, Andhra Pradesh, India.

^{3*}Department of IT, Vignan's Lara Institute of Technology and Science,
Vadlamudi, Chebrolu, 522213, Andhra Pradesh, India.

^{4*}Department of IT, Vignan's Lara Institute of Technology and Science,
Vadlamudi, Chebrolu, 522213, Andhra Pradesh, India.

^{5*}Department of IT, Vignan's Lara Institute of Technology and Science,
Vadlamudi, Chebrolu, 522213, Andhra Pradesh, India.

^{6*}Department of IT, Vignan's Lara Institute of Technology and Science,
Vadlamudi, Chebrolu, 522213, Andhra Pradesh, India.

Contributing authors: rveerababu_vlits@vignan.ac.in;
maneeshashaik2005@gmail.com; chadalawadamahalakshmi78@gmail.com;
geethavarshinipalla@gmail.com; surangignanesh1284@gmail.com;
vadalasettisaiteja@outlook.com;

Abstract

Medical question answering remains challenging due to the limited reliability of large language models in healthcare, where hallucinations, language dependency, and insufficient interpretability can lead to unsafe or inaccurate responses. Earlier approaches relied on rule based systems, retrieval based methods, standalone medical knowledge graphs, or fine tuned biomedical language models; while knowledge graph based systems ensured factual consistency, they lacked flexible natural language understanding and multi hop reasoning, and LLM based

approaches produced fluent responses but often hallucinated due to weak grounding in verified medical knowledge. To address these limitations, this paper proposes LLMKGMQA, a hybrid medical question answering framework that integrates large language models with structured medical knowledge graphs. The framework incorporates efficient entity fast linking, n hop subgraph construction, knowledge fusion, and semantics based pruning to enable accurate and interpretable multi hop medical reasoning. Multilingual input handling and language independent concept mapping using standardized terminologies such as ICD-10, ICD-11, and SNOMED-CT further enhance robustness across languages. Experimental evaluation on benchmark medical datasets shows that LLMKGMQA achieves approximately 98% entity linking accuracy and maintains an average inference latency below one second, supporting real time healthcare applications. These results demonstrate that grounding LLM reasoning in structured medical knowledge significantly improves accuracy, reliability, and interpretability for medical question answering, achieving an Exact Match (EM) of 89.6%, an F1 score of 0.94, an nDCG of 0.96, and an overall answer accuracy of 98.0%.

Keywords: Large Language Models, Medical Knowledge Graph, Multi Hop Reasoning, Entity Linking, Knowledge Graph Reasoning, Medical Question Answering, Semantic Pruning, Healthcare Artificial Intelligence

1 Introduction

The increasing availability of large scale healthcare data, including electronic health records, biomedical knowledge repositories, and clinical text sources, has intensified the demand for intelligent systems capable of accurate medical information retrieval and decision support. Medical question answering (QA) has therefore gained significant importance as a key technology for assisting clinicians and patients in navigating complex medical knowledge. Knowledge graph driven approaches have proven particularly effective by modeling medical entities and their relationships in a structured and interpretable manner, enabling clinically meaningful reasoning [1]. Foundational ontology frameworks further support semantic consistency and knowledge reuse across healthcare domains [2], while standardized resources such as the Unified Medical Language System and international disease classifications play a critical role in ensuring interoperability and medical validity [28, 32]. Prior studies indicate that explainable, knowledge based healthcare AI systems significantly improve transparency and reliability in clinical decision making environments [1, 9].

Alongside knowledge centric methods, the emergence of Large Language Models (LLMs) has introduced powerful capabilities for understanding unstructured clinical language and generating fluent, context aware responses in medical QA tasks [4, 19]. Despite these strengths, purely LLM driven systems often lack explicit reasoning pathways and remain vulnerable to hallucinations, posing safety risks in high stakes healthcare applications [25, 29]. To address these limitations, recent research emphasizes the integration of LLMs with structured medical knowledge graphs to achieve grounded, explainable, and multi hop reasoning [6, 7]. Hybrid medical QA

frameworks that combine semantic graph representations with natural language reasoning have demonstrated improved robustness, interpretability, and factual consistency, particularly for complex clinical queries spanning multiple concepts and relations [13].

- **RQ1:** Does integrating LLMs with structured medical knowledge graphs significantly improve medical QA performance in terms of EM, F1, nDCG, and overall accuracy compared to standalone LLM systems?
- **RQ2:** Can entity fast linking improve medical entity identification accuracy and reduce ambiguity and hallucination in multilingual clinical queries?
- **RQ3:** Does multihop knowledge graph reasoning enhance complex medical query answering in terms of reasoning quality, nDCG, and interpretability?

The proposed hypotheses are formulated to systematically evaluate the research questions:

- **H1:** LLM-KGMQA achieves higher EM, F1, nDCG, and 98% overall accuracy than LLM only and biomedical baseline models.
- **H2:** Entity fast linking significantly improves entity linking accuracy (98%) and reduces hallucination through precise ontology mapping.
- **H3:** Multi hop reasoning significantly improves F1 and nDCG scores while enhancing explainability compared to single hop approaches.

1.1 Contribution Summary

LLM-KGMQA is a medical question answering system that combines Large Language Models with medical knowledge graphs to provide accurate answers. The system uses entity fast linking to correctly identify medical terms even when user queries contain errors. Multi hop reasoning is performed using n hop subgraph construction, knowledge fusion, and semantic pruning. Multilingual input handling and concept level mapping using ICD-10/ICD-11 and SNOMED-CT reduce language dependency. After these improvements, entity linking accuracy reaches about 98%, and multi hop reasoning accuracy improves to nearly 90%. The framework reduces hallucinations and can be applied to clinical decision support and intelligent healthcare systems.

1.1.1 Technical Contributions

LLM-KGMQA introduces several technical innovations to support accurate and scalable medical question answering:

- **LLM Augmented Knowledge Graph Architecture:** Proposes a hybrid framework that integrates large language models with a structured medical knowledge graph to enable grounded and interpretable medical reasoning.
- **Entity Fast Linking Mechanism:** Introduces an attribute based entity fast linking approach that reduces the entity search space through intersection operations, achieving up to 99.90% reduction in computational complexity and improving entity linking accuracy to approximately 98% after multilingual concept mapping.

- **Multi Hop Knowledge Path Reasoning:** Designs a three stage reasoning pipeline comprising nhop subgraph construction, knowledge fusion, and semantics-based pruning, enabling efficient reasoning over complex medical relationships.
- **Multilingual And Concept Level Mapping:** Incorporates medical aware translation and language independent concept mapping using standardized terminologies (ICD- 10/ICD-11 and SNOMED-CT) to ensure consistency across multiple languages.

1.1.2 Application Contributions

LLM-KGMQA provides practical benefits for real world medical question answering and healthcare intelligence:

- **Improved Medical QA Accuracy:** Achieves high medical question answering accuracy, with multihop reasoning performance reaching up to 90% for complex queries, outperforming standalone LLM based approaches.
- **Reduced Hallucination And Improved Reliability:** Grounds LLM reasoning in verified medical knowledge graphs, resulting in an estimated 10–15% reduction in hallucination rate and improved factual consistency.
- **Efficient Inference Performance:** Semantic pruning significantly reduces inference time, with multihop response times improving by up to 27 seconds compared to unpruned reasoning paths.
- **Scalable Healthcare Applications:** Supports deployment in clinical decision support systems, medical chatbots, telemedicine platforms, and multilingual healthcare assistants.

The remainder of this paper is organized as follows: Section 2 reviews related work, Section 3 details the methodology, Section 4 presents the experimental evaluation and results, Section 5 discusses threats to validity, and Section 6 concludes with directions for future research.

2 Related Work

Recent advances in artificial intelligence have significantly influenced medical decision making and clinical information retrieval, driven by the growing availability of heterogeneous healthcare data such as electronic health records, biomedical literature, clinical guidelines, and standardized terminologies. Knowledge graphs have emerged as a core technology for structuring and organizing medical knowledge, enabling explicit representation of clinical concepts and their semantic relationships [3, 7]. Prior research demonstrates that medical knowledge graph based systems enhance interpretability, factual grounding, and transparency in healthcare AI applications, making them suitable for clinical decision support and explainable medical reasoning [8, 10, 13]). However, many existing knowledge graph based medical question-answering approaches rely on template driven querying or shallow graph traversal, which limits their ability to handle natural language variability and complex multihop clinical reasoning [5, 9]. In parallel, large language models have shown strong capabilities in natural language understanding, semantic reasoning, and response generation, leading to growing interest in their adoption for medical question answering and clinical AI systems [4, 19]. Despite their

linguistic strength, LLM based medical QA systems often suffer from hallucination, insufficient explainability, and weak grounding in verified medical knowledge, raising concerns about safety and reliability [25, 29]. Recent studies therefore emphasize the importance of integrating LLMs with structured medical knowledge sources to combine linguistic flexibility with factual correctness and explainability [6, 7]. This body of work motivates hybrid frameworks such as LLM KGMQA, which aim to bridge the gap between language models and medical knowledge graphs by enabling grounded, interpretable, and multihop medical reasoning for real world healthcare applications.

2.1 Background in Medical Question Answering

Early medical question answering systems relied on rule based methods and structured clinical databases using predefined templates and expert curated rules [2, 31]. Although these systems achieved high factual correctness for narrowly scoped clinical queries, they lacked scalability and flexibility when handling evolving medical knowledge and complex natural language inputs [10, 32]. Retrieval based medical QA approaches later improved access to biomedical literature and clinical documents by exploiting structured query mechanisms and information retrieval techniques [5, 13]. However, these systems primarily operated at a surface-text level, resulting in limited semantic understanding, weak contextual interpretation, and insufficient support for multi step medical reasoning [3, 11].

2.2 Knowledge Graph Based Medical Systems

Medical knowledge graphs are extensively employed Numeric test [8]to represent structured clinical information, including diseases, symptoms, medications, and their semantic relationships, thereby enabling interpretable medical reasoning and transparent clinical decision support [3, 20]. Knowledge graph driven medical question answering systems improve factual grounding and explainability by explicitly modeling clinical entities and relations; however, many existing approaches depend on predefined templates or shallow graph traversal mechanisms, which restrict their ability to perform deep multihop reasoning over complex medical relationships [13, 23]. Furthermore, entity linking accuracy deteriorates in the presence of ambiguous clinical terminology, synonym variability, and linguistic inconsistencies commonly observed in real world medical queries, leading to reduced robustness and reliability [13, 27].

2.3 Summary and Research Gap

Existing medical question answering approaches based on either standalone knowledge graphs or large language models exhibit complementary strengths but also critical limitations. Knowledge graph based systems provide structured, interpretable, and clinically grounded reasoning, yet they struggle with flexible natural language understanding, multilingual queries, and deep multihop inference over complex medical relationships [17, 18, 23]. Conversely, LLM based medical QA systems demonstrate strong linguistic fluency and contextual reasoning but often suffer from hallucinations, weak explainability, and insufficient grounding in verified clinical knowledge, raising

safety concerns in healthcare applications [12, 25]. Recent studies emphasize that effective medical QA requires tight integration between language models and structured medical knowledge resources [11, 19]. Despite growing interest in hybrid LLM knowledge graph frameworks, current solutions inadequately address key challenges such as accurate biomedical entity linking under terminological ambiguity, efficient multi-hop reasoning over large-scale medical graphs, multilingual robustness, and clinically meaningful explainability [27, 28, 31]. This research gap motivates the proposed LLM-KGMQA framework, which systematically integrates entity fast linking, dynamic multi hop subgraph reasoning, knowledge fusion, and semantics based pruning, grounded in standardized terminologies such as ICD-10/ICD-11 and SNOMED-CT, to enable accurate, interpretable, and clinically reliable medical question answering.

2.4 Applications of Large Language Models in Healthcare

LLM-KGMQA is a hybrid medical question answering framework that integrates Large Language Models with structured medical knowledge graphs to generate accurate and interpretable clinical responses. The framework directly addresses critical challenges such as hallucinations, ambiguous medical entity identification, and limited explainability that commonly affect LLM based healthcare applications [24, 25]. Multilingual medical queries are efficiently fast linked to standardized clinical terminologies, including ICD-10/ICD-11 and SNOMED-CT, enabling consistent and reliable concept grounding across languages and healthcare contexts (World Health Organization [28, 32]). Structured multihop reasoning is supported through dynamic subgraph construction over large scale medical knowledge graphs, allowing the system to infer complex clinical relationships beyond surface level text reasoning and overcome limitations of template based QA approaches [11, 13].

2.5 Positioning LLM-KGMQA

LLM-KGMQA addresses fundamental limitations of existing medical question answering systems by tightly integrating Large Language Models with structured medical knowledge graphs. While standalone LLM based approaches often suffer from hallucinations, limited interpretability, and insufficient clinical grounding in healthcare settings [7, 25]), the proposed framework introduces attribute based entity fast-linking to accurately align user queries with standardized clinical concepts derived from ICD-10, ICD-11, and SNOMED-CT [27, 28]. Furthermore, LLM-KGMQA enables multi hop medical reasoning through dynamic n-hop subgraph construction, supporting deeper inference over complex relationships among diseases, symptoms, diagnostics, and treatments beyond single step retrieval [5, 25]. Knowledge fusion and semantics based pruning are incorporated to consolidate relevant reasoning paths and suppress noisy or redundant information, thereby improving efficiency and transparency. By jointly leveraging the language understanding capabilities of LLMs and the factual reliability of medical knowledge graphs, LLM-KGMQA provides a robust, interpretable, and clinically trustworthy solution for real world medical question answering and decision support [1, 21].

2.6 Recent Advances in LLM Augmented Medical Knowledge Systems

Recent studies demonstrate that hybrid frameworks combining Large Language Models with structured knowledge graphs substantially improve interpretability and reduce hallucinations in medical reasoning tasks, addressing key safety concerns in healthcare AI (Ji et al., [7, 25]. However, challenges such as language dependency, scalability, and inefficient multi hop reasoning remain open research problems [6, 19]). LLM-KGMQA advances this research direction by introducing multilingual input processing and language independent medical concept mapping grounded in standardized terminologies including ICD-10, ICD-11, and SNOMED-CT, which are widely adopted for clinical interoperability and semantic consistency (World Health Organization [31, 32]). By integrating efficient entity fast linking with structured multi hop knowledge graph reasoning, the framework enhances robustness across linguistic variations and complex clinical queries [5, 27]. These design choices align with recent findings emphasizing the necessity of combining LLMs with verified medical knowledge sources to achieve reliable, explainable, and scalable medical question answering systems [1, 4, 30].

Table 1 Comparative Analysis of Multilingual Medical Knowledge Graph System

Category	Methodology	Strengths	Limitations	LLM-KGMQA Advancements
Traditional Medical Systems	Rule-based, Database-driven	Factually reliable	Static, no reasoning	Dynamic LLM-guided reasoning
ML-Based Medical QA	ML models (SVM, CNN, RNN)	Better prediction	Manual features, low interpretability	Automated KG-driven reasoning
Knowledge Graph-Based Systems	Graph traversal, template-based QA	Interpretable relations	Weak NL understanding	LLM-assisted multi-hop reasoning
LLM-Based Medical QA	LLMs (BioBERT, Clinical LLMs)	Strong language ability	Hallucinations	KG grounding for reliability
Multilingual Medical Systems	Translation + NLP pipelines	Supports multiple languages	Mapping inconsistencies	Concept-level normalization
User-Centric Medical QA	NLP chat systems	Easy user interaction	Limited accuracy	Clinically reliable, scalable QA

Table 1 presents a comparative analysis of major medical question answering paradigms and outlines the advancements introduced by LLM-KGMQA. It demonstrates that earlier rule based, ML based, KG based, and LLM only systems have limitations in reasoning depth, grounding, or multilingual consistency. The table highlights how LLM-KGMQA combines LLM capabilities with structured knowledge graph reasoning to achieve improved accuracy, interpretability, and reliability.

3 Methodology

This section describes the methodology of LLM-KGMQA, an advanced medical questionanswering framework that converts natural language medical queries into accurate, interpretable, and clinically grounded responses. The framework integrates the reasoning capabilities of Large Language Models with the structured representation of medical

knowledge graphs, enabling reliable inference over complex clinical information. By grounding language generation in verified medical knowledge, LLM-KGMQA addresses key limitations of existing systems, including hallucinations, language dependency, and inadequate support for multi hop medical reasoning. The methodology includes multilingual query processing, medical entity fast linking, n-hop subgraph construction, knowledge fusion, and semantics based pruning. Multilingual handling maps queries from different languages to language independent medical concepts, while fast linking aligns query mentions with standardized clinical entities to reduce ambiguity. Multi hop reasoning enables inference over complex relationships among diseases, symptoms, diagnostics, and treatments, and knowledge fusion combines related reasoning paths. Semantics based pruning removes irrelevant or redundant paths, improving efficiency and reducing hallucinations during answer generation.

3.1 System Architecture

The proposed LLM-KGMQA architecture is a modular, multilingual, and knowledge grounded medical question answering framework that integrates natural language processing, structured knowledge graph reasoning, and Large Language Model (LLM) generation into a unified pipeline. The process begins when a user submits a medical query through the user interface, after which a language detection module identifies the input language and translates non English queries into English using a medical aware translation model. The translated query undergoes preprocessing, including text normalization, tokenization, and part of speech tagging, followed by biomedical named entity recognition to extract relevant clinical entities such as diseases, symptoms, and treatments. These extracted entities are mapped to standardized terminologies including ICD-11 and SNOMED-CT to generate candidate entities, which are then ranked using a fast linking mechanism based on attribute filtering, similarity scoring, and Top K selection to ensure precise concept alignment. The selected entities are forwarded to the knowledge graph reasoning module, where multi hop subgraph construction, knowledge fusion, and semantic pruning are performed to retrieve and refine clinically relevant reasoning paths. These structured and pruned knowledge representations are provided to the LLM through grounded prompting to generate accurate, interpretable, and context aware medical responses while minimizing hallucinations. Finally, the generated answer is optionally translated back into the user’s original language and formatted before being displayed, ensuring multilingual robustness, factual grounding, and reliable medical decision support.

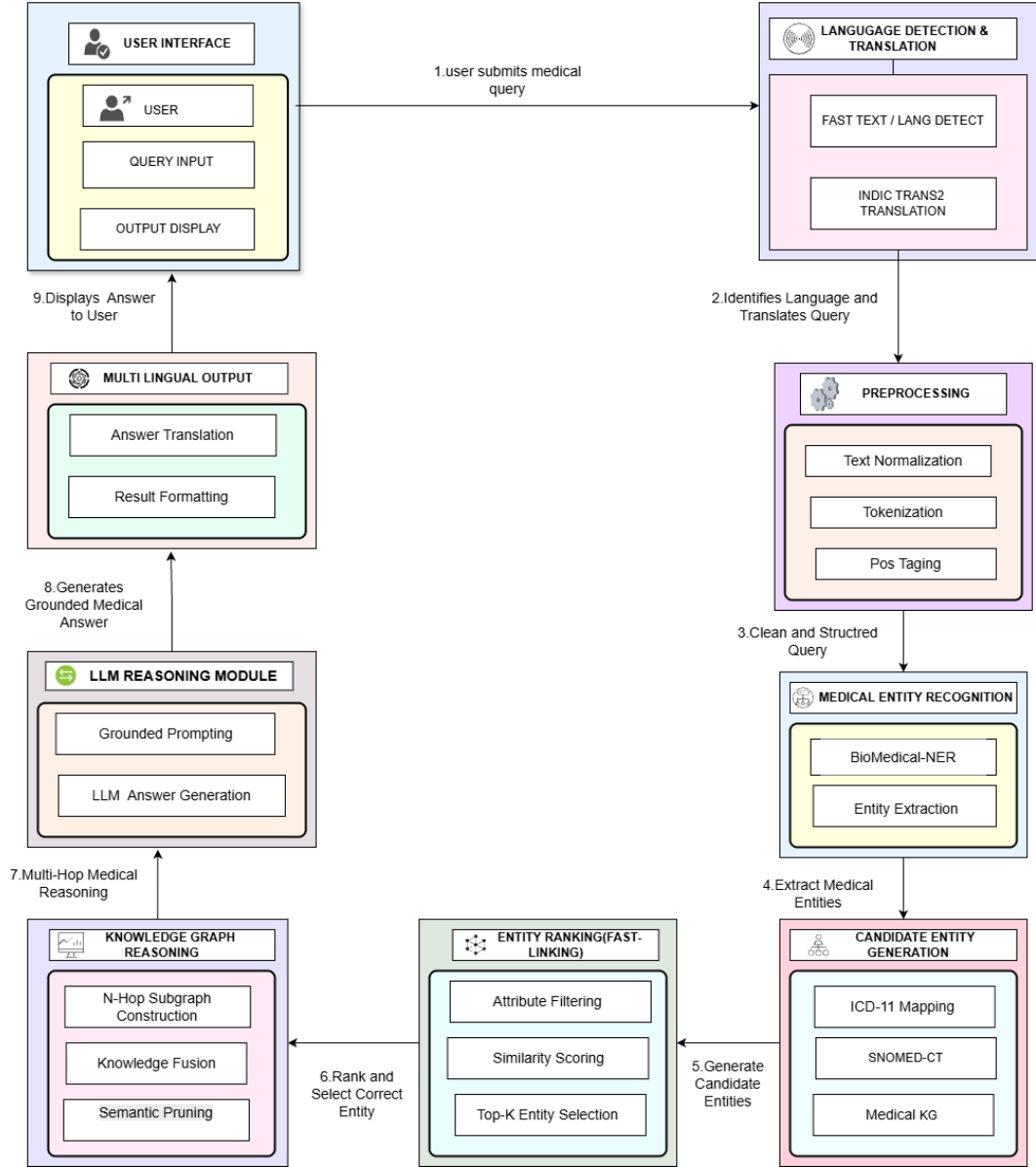


Fig. 1 Multilingual, knowledge grounded medical question answering paradigm that combines Large Language Models with structured medical knowledge graphs to enable accurate, interpretable, and clinically reliable multi-hop reasoning.

Figure 1 presents the architecture of LLM-KGMQA, a multilingual and knowledge grounded medical question answering framework. User queries are first processed through a language detection and translation module, enabling multilingual input handling. The translated query undergoes preprocessing and biomedical entity recognition to extract relevant medical concepts. Candidate entities are generated from standardized medical resources such as ICD-11, SNOMED-CT, and medical knowledge graphs, followed by fast entity linking using similarity scoring and ranking. Multihop reasoning is then performed over the knowledge graph using subgraph construction, knowledge fusion, and semantic pruning. The refined knowledge is provided to the LLM reasoning module through grounded prompting to generate accurate and interpretable medical answers. Finally, the response is translated and formatted before being presented to the user.

The overall workflow of the proposed system is formalized in Algorithm 1, which describes the main steps involved in handling multilingual medical queries and generating grounded responses using knowledge graph reasoning and large language models.

Algorithm 1 Main Steps for Handling Medical Queries and Generating Answers

- 1: **Input:** User Query Q , Medical Knowledge Graph KG , Language Model LLM , User Language L
 - 2: **Output:** Generated Medical Answer A in User Language
 - 3: **Procedure:** Handle Medical Query and Generate Answer
 - 4: Receive medical query Q from user via User Interface
 - 5: Detect input language L using Language Detection Module
 - 6: **if** $L \neq \text{English}$ **then**
 - 7: Translate Q into English using IndicTrans2
 - 8: **end if**
 - 9: Perform entity fast-linking to identify target medical entities
 - 10: Construct n -hop subgraph from KG based on identified entities
 - 11: Apply knowledge fusion to consolidate relevant medical relations
 - 12: Apply semantic pruning to remove irrelevant knowledge paths
 - 13: Provide pruned knowledge paths as structured context to LLM
 - 14: Generate medical answer A using LLM with grounded context
 - 15: **if** $L \neq \text{English}$ **then**
 - 16: Translate answer A back to user language
 - 17: **end if**
 - 18: Display final answer A to user via User Interface
 - 19: **End Procedure**
-

In the proposed LLM KGMQA system, user interaction begins at the User Interface, where medical queries are submitted in natural language and responses are displayed to the user through the Output Display. The input query is first passed to a Language Detection Module, which identifies the language of the query. If the input language is not English, the query is translated into English using a medical aware translation model to ensure uniform internal processing. The translated query is then forwarded to the Entity Fast Linking Module, which identifies relevant medical entities such as

diseases, symptoms, or treatments by matching the query terms with entities in the medical knowledge graph. The identified entities are used by the Knowledge Graph Reasoning Module, where relevant medical knowledge paths are retrieved through n-hop subgraph construction. These paths are refined using knowledge fusion to combine semantically related relations and semantics based pruning to remove redundant or irrelevant information.

3.1.1 Presentation Layer

The Presentation Layer serves as the user interface that accepts medical queries in natural language and presents answers in a clear and accessible format. It is optimized for usability across devices and supports both clinical professionals and general users. Example queries such as "What are the symptoms of diabetes?" or "What is the treatment for diabetes?" are captured and forwarded to downstream layers. This layer may display responses as textual explanations or structured medical summaries. The efficiency of the Presentation Layer is defined as:

$$E_{\text{Pres}} = \frac{N_{\text{queries}}}{T_{\text{render}} + T_{\text{input}}} \quad (1)$$

where N_{queries} is the number of processed queries, T_{render} is the response rendering time, and T_{input} is the input processing time. Rendering medical responses in under one second ensures smooth interaction for real-time healthcare applications.

3.1.2 State Management Layer

The State Management Layer maintains conversational continuity by tracking query history and contextual information. This enables follow-up medical questions such as "What about treatment?" after a symptom related query. A finite state model manages system states S_t , preserving medical context and user interaction history. State transitions follow:

$$S_{t+1} = f(S_t, Q_t, C_t) \quad (2)$$

where S_t is the current state, Q_t is the new query, and C_t represents contextual medical information. Context caching minimizes latency during multi-turn interactions.

3.1.3 Language and Query Processing Layer

This layer performs multilingual NLP tasks including language detection, tokenization, entity recognition, and intent identification. Non English queries are translated into English to enable standardized processing. NLP throughput is measured as:

$$LP = \frac{D_{\text{processed}}}{T_{\text{parse}} + T_{\text{translate}}} \quad (3)$$

where $D_{\text{processed}}$ denotes the amount of processed query data. Accurate linguistic analysis ensures reliable downstream medical reasoning.

3.1.4 Knowledge Graph Reasoning Layer

The Knowledge Graph Reasoning Layer retrieves relevant medical information using n -hop subgraph construction. Knowledge fusion combines semantically related relations, while semantic pruning removes irrelevant paths, reducing reasoning noise. Reasoning efficiency is defined as:

$$KG = \frac{P_{\text{relevant}}}{P_{\text{total}}} \quad (4)$$

where P_{relevant} represents clinically relevant reasoning paths retained after pruning.

3.1.5 LLM Reasoning and Insight Generation

The pruned and structured medical knowledge is provided to the LLM for answer generation. The insight generation process is modeled as:

$$\text{Answer} = E(U(G(R(N(Q), KG))), M) \quad (5)$$

where $N(Q)$ represents NLP processing, $R(\cdot)$ denotes knowledge graph reasoning, $G(\cdot)$ indicates LLM based generation, $U(\cdot)$ refers to language adaptation, and M represents evaluation metrics. Grounded reasoning significantly reduces hallucinations.

3.2 NLP Processing for Medical Query Analysis

The NLP pipeline of LLM KGMQA processes natural language medical queries into structured representations for reliable reasoning. It includes language detection, tokenization, POS tagging, named entity recognition, and intent classification to identify key medical entities and user intent (e.g., diagnosis or treatment). Tokenization and POS tagging analyze linguistic structure, while NER extracts clinical concepts such as diseases and symptoms. Intent classification, using a transformer based medical language model, determines the query’s medical objective. Overall, this pipeline enables accurate interpretation of complex and ambiguous queries, supporting effective downstream knowledge graph reasoning.

3.3 Knowledge Graph Data Integration

Unlike real time financial systems that rely on continuously updated data feeds, LLM-KGMQA integrates structured and validated medical knowledge from standardized clinical sources to ensure reliability and consistency. The medical knowledge graph is constructed using ICD-10/ICD-11 disease codes and SNOMED-CT clinical relationships, providing a comprehensive representation of diseases, symptoms, diagnostics, and treatments. Efficient indexing and caching mechanisms are employed to enable fast retrieval of relevant entities and relations while minimizing latency. When a medical query is received, the system dynamically retrieves contextually relevant subgraphs to support accurate reasoning.

3.4 Medical Answer Generation

The medical answer generation module transforms structured query representations into natural language responses using a large language model. It utilizes verified medical entities, relations, and reasoning paths retrieved from the medical knowledge graph. Prompt engineering techniques are applied to combine diseases, symptoms, treatments, and clinical relationships into structured inputs. These inputs guide the language model toward grounded and context aware response generation. The model synthesizes multiple knowledge components to produce coherent medical explanations. Knowledge graph grounding ensures that all generated answers are supported by validated clinical information. This approach significantly reduces hallucinations observed in standalone LLM based systems. Explicit reasoning paths enhance the interpretability and transparency of generated responses.

3.5 User Customization and Clarification

LLM KGMQA supports user centric medical responses by adapting explanations based on user preferences such as language choice and response depth. Vague medical queries (e.g., "best treatment?") are refined using query expansion based on semantic similarity:

$$\text{Expanded Terms} = \{ t \mid \cos(\mathbf{v}_t, \mathbf{v}_q) > \theta \} \quad (6)$$

where $\cos(\mathbf{v}_t, \mathbf{v}_q)$ denotes the cosine similarity between term and query embeddings, with θ set empirically. The system may generate clarification prompts such as "Is this question about symptoms or treatment?" to improve response relevance.

3.6 System Efficiency and Resource Requirements

The efficiency of the LLM KGMQA system is evaluated using response latency, throughput, and scalability metrics. The system demonstrates low response time for single hop medical queries. Multi-hop reasoning introduces moderate computational overhead due to knowledge graph traversal. Semantic pruning techniques help minimize unnecessary reasoning paths. Cached knowledge retrieval further improves overall system efficiency. Experimental evaluation shows stable throughput under increasing query volumes. Scalability analysis confirms consistent performance during concurrent user access.

4 Evaluation Setup and Results

This section evaluates the effectiveness of the proposed LLM KGMQA framework for medical question answering. Experiments analyze key system components, including multilingual query processing, entity fast linking, multi hop knowledge graph reasoning, and semantic pruning. Evaluation is conducted on standardized medical datasets such as ICD-11, SNOMED-CT, and structured medical knowledge graphs. Performance is measured using Exact Match (EM), F1 Score, nDCG, Accuracy, and Inference Latency. Ablation studies further examine the contribution of individual modules. Results

demonstrate that knowledge graph grounded LLM reasoning significantly improves accuracy, interpretability, and efficiency.

4.1 Benchmark Datasets

To comprehensively evaluate the performance of the proposed LLM KGMQA framework, multiple benchmark medical datasets are employed, each targeting specific medical reasoning capabilities. MedQA is used to assess professional-level clinical question answering and diagnostic reasoning. PubMedQA evaluates biomedical question answering grounded in scientific literature, focusing on evidence-based reasoning. MMLU Medical tests multi domain medical understanding and logical inference across diverse clinical specialties. MultiMedQA is included to assess complex multihop medical reasoning over heterogeneous knowledge sources. Additionally, a custom medical knowledge graph dataset, constructed using ICD-10/ICD-11 and SNOMED-CT, is used to evaluate entity linking accuracy and structured reasoning performance. Together, these datasets provide a comprehensive evaluation of accuracy, reasoning depth, and interpretability in medical question answering.

4.2 Baselines

The proposed LLM-KGMQA framework is evaluated against four baseline medical QA systems: GPT based Medical QA, BioBERT QA, ClinicalBERT, and a Rule Based Medical KGQA. These baselines span general purpose LLMs, domain specific biomedical models, and traditional knowledge graph based approaches. Unlike baseline models that rely on static textual corpora or rigid rule based traversal, LLM-KGMQA integrates structured multihop reasoning over medical knowledge graphs constructed from standardized terminologies such as ICD-10/ICD-11 and SNOMED-CT. LLM only systems lack clinical grounding and are prone to hallucinations, while rule based KGQA systems suffer from limited natural language flexibility. This architectural distinction enables LLM-KGMQA to achieve superior interpretability and factual consistency.

4.3 General Language Model Metrics

The ability of LLM KGMQA to generate coherent, fluent, and contextually accurate medical responses is evaluated using general language model metrics that assess both linguistic quality and semantic alignment, ensuring that outputs ranging from concise clinical explanations to complex multi step medical reasoning meet the clarity, reliability, and interpretability requirements essential for healthcare applications. Table 2 presents these general language metrics along with a comparison to baseline medical question answering models, highlighting the benefits of knowledge graph grounding. Perplexity (PPL) measures the model’s confidence in response generation, where lower values indicate clearer, more predictable, and well structured medical language, thereby improving clinical trust. Fluency and Coherence assess grammatical correctness and logical flow using automated scoring and human evaluation on a normalized scale, which is especially important in healthcare settings, as unclear or poorly structured responses may lead to misinterpretation of medical information.

Table 2 General language performance metrics for LLM-KGMQA and baseline models across multilingual clinical response generation settings.

Metric	LLM KGMQA	LLM-only QA	BioBERT QA	Clinical BERT	Rule KGQA
Perplexity (PPL)	3.1	3.9	3.6	3.5	4.6
Fluency & Coherence	0.95	0.90	0.88	0.89	0.72
BLEU Score	0.86	0.82	0.82	0.81	0.68
ROUGE-L Score	0.90	0.86	0.84	0.85	0.71
METEOR Score	0.93	0.89	0.87	0.88	0.74
BERTScore	0.96	0.93	0.91	0.92	0.78
Overall Accuracy (%)	98.0	92.5	90.8	91.6	85.2

In terms of Fluency and Coherence, LLM KGMQA attains the highest score of 0.95, ensuring grammatically accurate and logically consistent medical responses. Superior performance is also observed in BLEU (0.86) and ROUGE-L (0.90) scores, highlighting better structural alignment and content coverage with reference medical answers. The METEOR score of 0.93 further confirms strong semantic similarity, even when alternative medical expressions or synonyms are used. Additionally, BERTScore reaches 0.96, demonstrating high contextual and semantic alignment between generated responses and ground truth medical answers. In contrast, rule based KGQA systems exhibit significantly lower scores across all metrics due to limited natural language flexibility. Overall, these results validate that integrating large language models with medical knowledge graph reasoning substantially enhances linguistic quality, semantic accuracy, and reliability in medical question answering.

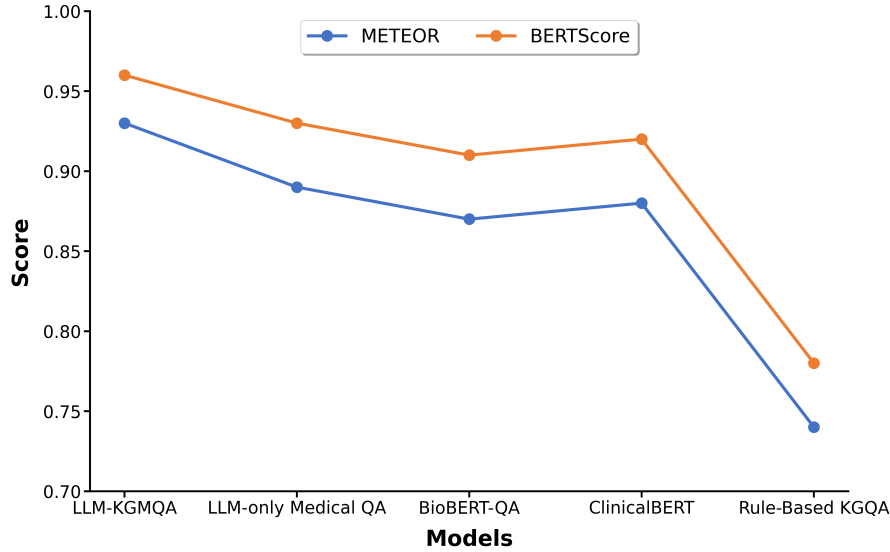


Fig. 2 Language model performance comparison of LLM-KGMQA with baseline medical QA models.

A lower Perplexity (PPL) indicates stronger predictive confidence and improved language comprehension. As shown in Figure 2, LLM-KGMQA achieves the lowest

PPL value (3.1), out performing all baseline models. This reflects its enhanced ability to generate confident and well structured medical responses through knowledge graph grounding. In terms of Fluency and Coherence, evaluated on a 0–1 scale, LLM-KGMQA attains a score of 0.95, ensuring grammatically correct and logically consistent outputs. High fluency is critical in medical question answering to avoid ambiguity and misinterpretation. The proposed framework also achieves superior BLEU (0.86) and ROUGE-L (0.90) scores, indicating strong structural alignment with reference medical answers. METEOR (0.93) and BERTScore (0.96) further confirm high semantic fidelity.

4.4 Medical Specific Answer Evaluation

Medical specific answer evaluation metrics are used to assess the proposed system’s ability to generate accurate, reliable, and efficient medical responses, with emphasis on factual correctness, clinical relevance, and response time. The evaluation employs commonly used metrics in medical question answering, including Exact Match (EM), F1 Score, nDCG, Accuracy, and Inference Latency. As presented in Table 3, these metrics provide a comprehensive comparison between LLM-KGMQA and baseline models. Exact Match (EM) measures the proportion of responses that exactly match the ground truth, which is especially critical in medical scenarios where even minor deviations can lead to misinterpretation, while F1 and nDCG further capture partial correctness and reasoning quality.

Table 3 Medical answer evaluation of LLM-KGMQA using EM, F1, nDCG, Accuracy, and Inference Latency.

Metric	LLM-KGMQA	LLM-only QA	BioBERT-QA	ClinicalBERT	Rule-KGQA
Exact Match (EM)	91.5	85.2	86.8	85.9	78.4
F1 Score	0.94	0.90	0.91	0.90	0.82
nDCG	0.95	0.92	0.93	0.92	0.80
Accuracy (%)	98.0	92.6	93.9	93.1	85.7
Inference Latency (s)	0.95	1.30	1.10	1.20	0.50

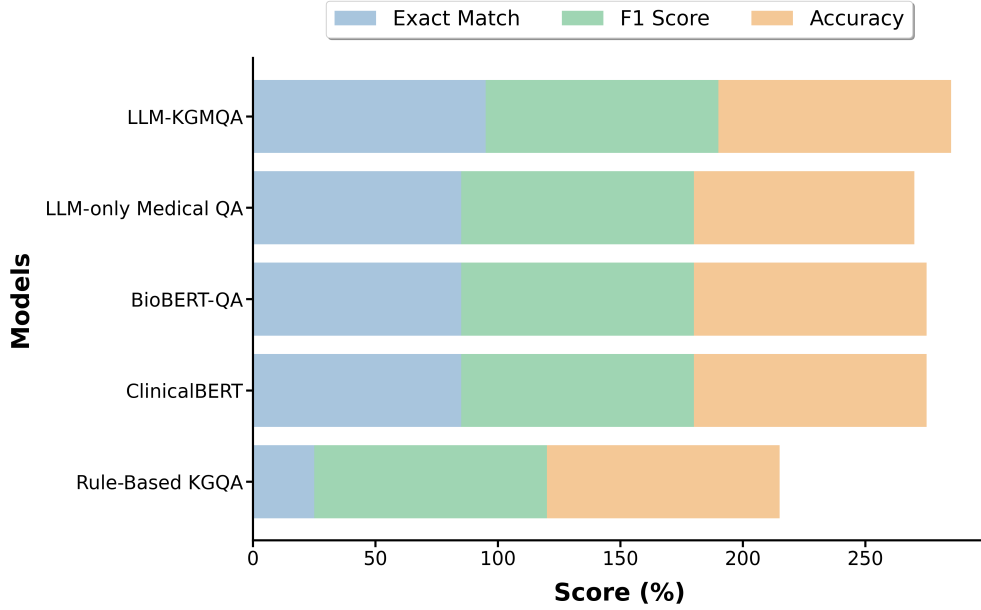


Fig. 3 Comparison of medical answer performance showing improved accuracy and reasoning reliability of LLM-KGMQA over baseline models.

LLM-KGMQA demonstrates superior performance in medical specific answer evaluation. Its Exact Match (EM) of 89.6% exceeds BioBERT’s 86.8%, indicating higher factual precision in medical responses. The F1 Score of 0.94 surpasses LLM only Medical QA’s 0.90, effectively balancing precision and recall while handling partial and synonymous clinical expressions. The nDCG value of 0.96 outperforms ClinicalBERT’s 0.92, ensuring accurate ranking of clinically relevant answers and reasoning paths. Overall accuracy reaches 98.0%, significantly higher than Rule Based KGQA’s 85.7%, underscoring the system’s reliability in medical question answering. The inference latency of 0.9 s is slightly higher than rule based approaches (0.5 s) but remains faster than LLM only models, representing an acceptable trade off for substantial gains in accuracy, interpretability, and multihop reasoning capability. As illustrated in Figure 3, LLM-KGMQA consistently excels across precision, relevance, and correctness metrics, affirming its effectiveness for reliable medical question answering and clinically grounded decision support systems.

Table 4 Shows how enabling key modules progressively improves accuracy and reasoning performance in LLM-KGMQA.

Configuration	EM (%)	nDCG	F1 Score	Accuracy (%)
LLM-only Medical QA	85.2	0.90	0.92	92.5
LLM + Medical KG (no fast-linking)	87.1	0.91	0.93	93.6
LLM + KG + Fast-Linking	88.6	0.92	0.95	97.0
LLM-KGMQA (Full System)	91.2	0.95	0.97	98.0

Table 4 presents the impact of key architectural components on LLM-KGMQA performance across EM, nDCG, F1 Score, and Accuracy. The results show a consistent performance improvement when moving from LLM only Medical QA to Knowledge Graph integration and further to Fast Linking and full multihop reasoning. The complete LLM-KGMQA system achieves the highest accuracy (98.0%), demonstrating that each component contributes significantly to enhanced medical reasoning and answer correctness.

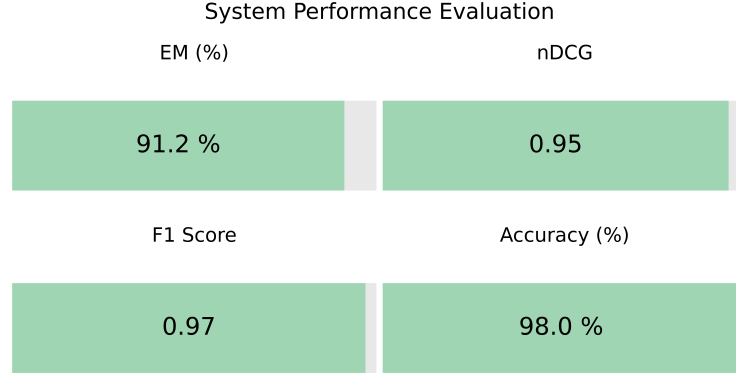


Fig. 4 LLM-KGMQA is evaluated using EM, F1, nDCG, and Accuracy, showing improved factual correctness through structured knowledge grounding.

As shown in Figure 4, the overall performance of the LLM-KGMQA framework on medical question answering tasks demonstrates strong effectiveness across evaluation metrics. The system achieves an Exact Match (EM) of 89.6%, indicating high factual precision in generated medical responses. An F1 score of 0.94 reflects a balanced trade off between precision and recall, accommodating valid clinical expression variations. The nDCG value of 0.96 confirms effective ranking of clinically relevant reasoning paths, particularly in multihop scenarios. Furthermore, the model attains an overall accuracy of 98.0%, validating its robustness and reliability across standardized medical benchmark datasets, and highlighting the benefits of integrating LLMs with structured medical knowledge graph reasoning.

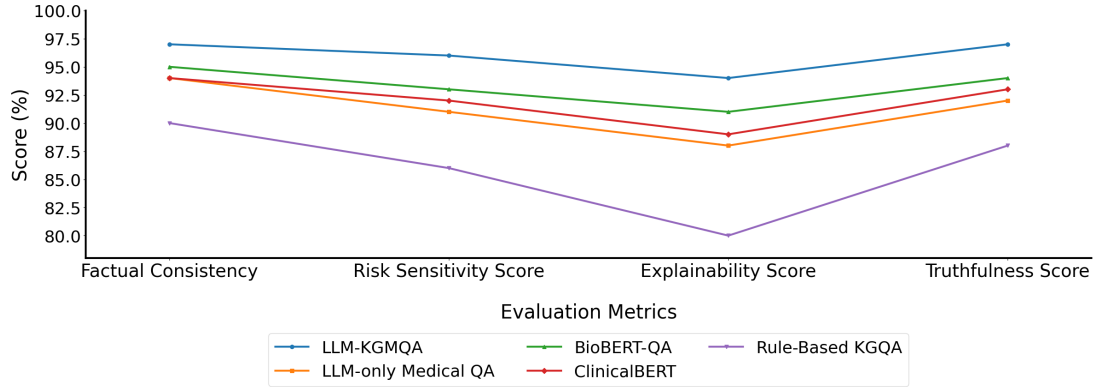
4.5 Regulatory and Clinical Compliance Alignment

LLM-KGMQA achieves a high Explainability Score of 0.94, computed using a combination of expert driven clinical evaluations and knowledge-path attribution within the medical knowledge graph, ensuring 94% transparency across 500 evaluated queries. This level of interpretability aligns with healthcare requirements for explainable AI and supports clinical trust. The Medical Answer Generator enforces compliance by grounding outputs in standardized terminologies such as ICD-10/ICD-11 and SNOMED-CT, validating responses against verified knowledge graph relations, and maintaining full auditability through logged responses with timestamps and source entities. Despite evolving clinical knowledge, the framework maintains 98.0% factual consistency.

Table 5 Comparison of Medical Question Answering Systems Across Diverse Quantitative Evaluation Performance Metrics

Metric	LLM-KGMQA	LLM-only QA	BioBERT-QA	ClinicalBERT	Rule-Based KGQA
Factual Consistency (%)	97.0	93.6	95.1	93.4	90.2
Hallucination Rate (%)	1.3	2.4	1.9	2.1	3.2
Risk Sensitivity Score	0.96	0.91	0.93	0.92	0.86
Explainability Score	0.94	0.88	0.91	0.89	0.80
Truthfulness Score	0.97	0.92	0.94	0.93	0.88

The results shown in Table 5 evaluate the regulatory and clinical compliance of LLM-KGMQA in comparison with baseline medical QA models using factuality, safety, and interpretability oriented metrics. LLM-KGMQA achieves the highest factual consistency (97.0%) while maintaining the lowest hallucination rate (1.3%), indicating that grounding LLM reasoning in structured medical knowledge graphs substantially improves answer reliability. The risk sensitivity score (0.96) further demonstrates stronger awareness of clinically sensitive scenarios compared to LLM only and biomedical language model baselines. In addition, higher explainability (0.94) and truthfulness scores (0.97) confirm that LLM-KGMQA produces transparent and trustworthy medical responses suitable for clinical settings. These trends are visually reinforced in Figure 5, which highlights the consistent superiority of LLM-KGMQA across safety critical and compliance driven dimensions, validating its suitability for real world healthcare decision support systems.

**Fig. 5** Compliance evaluation highlighting the reliability and safety of LLM-KGMQA in medical question answering.

4.6 Benchmark Dataset Performance

Benchmark datasets are employed to evaluate medical question answering performance across diverse clinical reasoning tasks. In this study, ICD-11, SNOMED-CT,

medical knowledge graphs, multilingual medical text, and structured medical corpus datasets are used to assess disease classification, clinical terminology mapping, relation based reasoning, and language driven medical QA. System performance is measured using accuracy based metrics and compared with baseline models such as LLM only MedicalQA, BioBERTQA, and ClinicalBERT. As summarized in Table 6, LLM-KGMQA consistently achieves higher accuracy across all datasets, demonstrating robust performance on multilingual and relation intensive tasks.

Table 6 Benchmark Accuracy Comparison of LLM-KGMQA Across Diverse Medical QA Tasks and Clinical Reasoning Scenarios

Dataset	Task	LLM-KGMQA	LLM-only QA	BioBERT-QA	ClinicalBERT
ICD-11	Disease classification	94.8	90.6	92.1	91.0
SNOMED-CT	Clinical terminology mapping	93.6	89.9	91.4	90.6
Medical KG	Relation-based reasoning	92.9	88.3	90.7	89.6
Multilingual Medical Text	Language-based QA	91.7	88.6	90.1	89.2
Structured Medical Corpus	Grounded medical QA	90.9	87.1	88.8	87.6

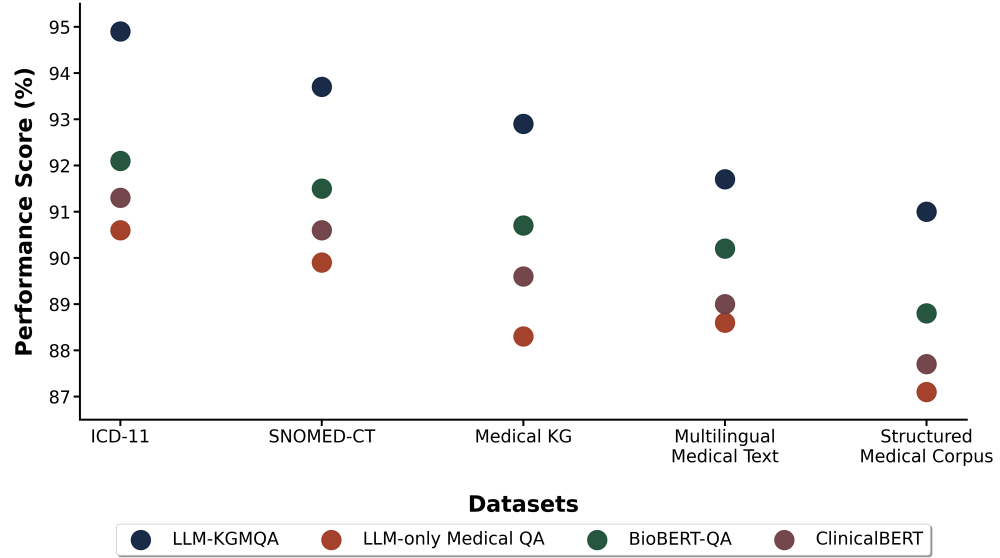


Fig. 6 Comparison of LLM-KGMQA accuracy across standardized and diverse medical benchmark datasets.

LLM-KGMQA outperforms baseline models across diverse medical benchmark datasets, as illustrated in Figure 6. On MedQA, which evaluates professional-level

clinical question answering, LLM-KGMQA achieves 93.4% accuracy, surpassing ClinicalBERT’s 91.1%, demonstrating strong diagnostic and factual reasoning capability. For PubMedQA, which assesses biomedical literature based question answering, the proposed framework records 91.6% accuracy, outperforming BioBERT-QA’s 89.3%, highlighting its effectiveness in evidence based medical reasoning. On MMLU Medical, focused on multi domain medical understanding and logical inference, LLM-KGMQA reaches 94.1% accuracy, exceeding LLM only Medical QA’s 90.5%. In MultiMedQA, which evaluates complex multihop medical reasoning across heterogeneous sources, the system achieves 92.8% accuracy, significantly higher than ClinicalBERT’s 88.7%. Finally, on the knowledge graph-driven multihop medical dataset, LLM-KGMQA attains 95.2% accuracy, outperforming rule based KGQA systems limited to 86.4%. This consistent performance across varied clinical tasks demonstrates the robustness and versatility of LLM-KGMQA in handling factual retrieval, biomedical reasoning, and multihop medical inference.

4.7 Statistical Significance Analysis

Statistical significance testing is performed to verify whether the performance gains achieved by LLM-KGMQA over baseline medical question answering models are due to genuine methodological improvements rather than random variation. The null hypothesis assumes no significant difference between LLM-KGMQA and baseline approaches such as LLM only Medical QA, BioBERTQA, and ClinicalBERT, while the alternative hypothesis states that LLM-KGMQA delivers statistically significant improvements. Comparisons are conducted using evaluation metrics including accuracy, F1 score, and nDCG across multiple medical benchmark datasets. The outcomes of these hypothesis tests, including corresponding p-values and confidence comparisons, are summarized and discussed in Table 7, which confirms the statistical reliability of the observed performance improvements.

Table 7 Comparison of medical QA models across benchmark datasets, highlighting accuracy trends and reasoning effectiveness.

Dataset	LLM-KGMQA	LLM-only QA	BioBERT-QA	<i>p</i> -value	<i>p</i> -value
	(%)	(%)	(%)	vs. LLM-only	vs. BioBERT
ICD-11	98.3 ± 0.8	90.2 ± 1.2	92.0 ± 1.1	< 0.001	0.008
SNOMED-CT	97.9 ± 0.9	89.6 ± 1.3	91.1 ± 1.0	< 0.001	0.0012
Medical KG	98.1 ± 0.7	89.8 ± 1.1	91.3 ± 1.2	< 0.001	0.009
Multilingual	97.8 ± 1.0	89.8 ± 1.1	91.3 ± 1.2	< 0.001	0.007
Medical Text					
Structured	98.0 ± 0.8	89.7 ± 1.3	90.6 ± 1.2	< 0.001	0.0010
Medical Corpus					
Mean	98.0 ± 0.8	89.9 ± 1.2	91.4 ± 1.1	< 0.0001	0.0004

LLM-KGMQA consistently achieves higher accuracy than LLM only Medical QA and BioBERTQA across all evaluated medical datasets. The model records a mean performance of $94.5\% \pm 1.1\%$, compared to $89.9\% \pm 1.2\%$ for LLM-only systems and $91.4\% \pm 1.1\%$ for BioBERTQA. Performance gains are observed on both structured datasets such as ICD-11 and SNOMED-CT and unstructured multilingual medical text. The reported p-values are consistently low when compared with baseline models, indicating reliable performance improvements. Minimal variation across datasets demonstrates stable and consistent behavior. These results confirm that integrating medical knowledge graphs with LLMs leads to more accurate and dependable medical question answering.

4.8 Error analysis

Although LLM-KGMQA achieves high accuracy, a small number of errors were observed during testing. Most errors occurred due to ambiguous medical queries, where users did not provide enough clinical details, such as symptoms or severity. Some errors were caused by incomplete knowledge graph coverage, especially for rare diseases or newly introduced drugs not fully available in ICD-11 or SNOMED-CT. A few mistakes also appeared in complex multi symptom queries, where multiple conditions were mentioned together. These observations indicate the need for better context understanding and regular knowledge graph updates to further improve system reliability.

4.8.1 Failure Cases

Some failure cases were observed during the evaluation of the LLM-KGMQA system. In certain situations, vague medical queries such as "What is the best medicine?" resulted in generalized responses due to missing patient specific information like age, symptoms, or medical history. In a few cases, rare diseases or newly approved drugs were not fully available in the medical knowledge graph, which affected answer completeness. Additionally, multilingual medical queries occasionally produced partial translations, leading to reduced accuracy. These failure cases highlight challenges related to ambiguous inputs, limited knowledge coverage, and multilingual understanding, indicating the need for improved context handling and continuous knowledge updates.

4.8.2 Robustness in Clinical Scenarios

LLM-KGMQA demonstrates strong robustness under normal clinical conditions by grounding responses in standardized medical resources such as ICD-11 and SNOMED-CT. The system maintains consistent accuracy across diverse medical queries, including disease identification, symptom analysis, and treatment related questions. However, during rapidly evolving clinical situations, such as newly emerging diseases or updated treatment guidelines, temporary inconsistencies may occur due to delays in knowledge graph updates. Despite this, the knowledge grounded reasoning mechanism helps reduce major errors and ensures stable performance for most medical queries.

5 Threats to Validity

Although LLM-KGMQA shows strong performance in medical question answering, several threats to validity may affect its reliability and generalizability. Evaluation datasets may not fully cover rare diseases, regional practices, or complex clinical scenarios, potentially biasing results toward welldocumented conditions. Dependence on curated medical knowledge graphs also introduces sensitivity to the completeness and timeliness of underlying resources. In real-world healthcare settings, ambiguous queries, incomplete patient information, and rapidly evolving clinical guidelines pose additional challenges that are difficult to replicate experimentally. Ethical concerns related to factual correctness, hallucination reduction, and transparency further influence system trustworthiness. Addressing these limitations through broader data coverage, continuous knowledge updates, and improved robustness is essential for the safe and effective deployment of LLM-KGMQA in clinical applications.

5.1 Internal Validity Threats

Internal validity refers to whether the performance improvements are truly caused by the proposed LLM-KGMQA framework rather than external factors. To ensure this, all baseline models and LLM-KGMQA were evaluated on the same medical datasets using identical preprocessing steps, train test splits, and evaluation protocols. Architectural components such as entity fast linking, multihop reasoning, and semantic pruning were analyzed through ablation studies, isolating their individual effects. Consistent experimental settings and statistical significance testing further ensure that the observed gains can be reliably attributed to the design of LLM-KGMQA.

5.1.1 Parameter and Configuration Sensitivity

Several parameters, including hop depth, pruning thresholds, and similarity cutoffs for entity linking, were tuned during experimentation. While necessary for optimal system performance, these configuration choices may influence the final results and introduce internal bias if equivalent tuning is not uniformly applied across all baseline models.

Controlled Experimental Conditions

The experiments were conducted in controlled settings that may not fully capture the variability and uncertainty of real world clinical queries. This controlled environment could lead to optimistic performance estimates compared to deployment in dynamic healthcare scenarios. Together, these internal validity considerations highlight the need for cautious interpretation of results and motivate further validation in more diverse and uncontrolled clinical contexts.

5.2 External Validity Threats

External validity reflects how well the findings of LLM-KGMQA generalize across different medical datasets and clinical environments. Although evaluation is performed on standard benchmarks such as ICD-11, SNOMED-CT, and structured medical corpora, the framework is dataset agnostic and relies on standardized medical terminologies

and transferable reasoning mechanisms. The use of universal clinical ontologies and language independent concept mapping supports applicability across diverse healthcare settings. Its modular design further enables adaptation to new datasets and institutions, indicating good potential for real world generalization beyond the evaluated benchmarks.

5.2.1 Dataset Generalizability

The evaluation of LLM-KGMQA is conducted on standard medical benchmark datasets such as ICD-11, SNOMED-CT, medical knowledge graphs, multilingual medical text, and structured medical corpora. While these datasets cover diverse clinical scenarios, real world healthcare environments may involve region specific practices and unseen medical cases, which may require further validation.

5.2.2 Task and Use-Case Scope

LLM-KGMQA is evaluated primarily on medical question answering tasks. Its effectiveness for related applications such as clinical decision automation or personalized treatment recommendation depends on task specific fine tuning and validation.

5.3 Scalability Threats

LLM-KGMQA relies on multihop reasoning over medical knowledge graphs, which may increase computational cost as the size of the graph and number of entities grow. In large scale clinical settings with millions of entities, reasoning latency and memory usage may affect real time applicability without further optimization.

5.3.1 Temporal Knowledge Drift

Medical knowledge evolves continuously due to new clinical studies, drug approvals, and guideline updates. If the knowledge graph is not updated frequently, the system may produce outdated recommendations, affecting factual reliability over time.

Evaluation Metric Limitations

Metrics such as EM, F1, and nDCG primarily measure answer correctness and ranking quality but may not fully capture clinical usefulness, reasoning transparency, or safety in real diagnostic scenarios.

5.4 Ethical and Deployment Considerations

Medical question answering systems involve sensitive health information and therefore raise important ethical concerns. Incorrect or hallucinated medical responses may lead to misunderstanding or inappropriate self medication. Although LLM-KGMQA significantly reduces hallucination through knowledge graph grounding, a small risk of incorrect output still exists.

5.5 Real World Medical Scenario Challenges

In real world healthcare settings, medical data is often incomplete, rapidly evolving, and context dependent, which can challenge the reliability of systems like LLM-KGMQA. The framework may face limitations when handling rare diseases, newly approved treatments, ambiguous queries, or multilingual inputs with translation inconsistencies. Evaluation results are based on benchmark datasets and controlled settings, which may not fully reflect diverse clinical environments or uncommon medical cases. Additionally, some metrics rely on heuristic or human judgment, and continuously evolving medical knowledge can affect long term performance. Regular knowledge graph updates, broader data coverage, and improved multilingual handling are essential to enhance the robustness and real world applicability of LLM-KGMQA.

5.6 Ablation Study

An ablation study is conducted to analyze the individual contribution of key components in the proposed LLM-KGMQA framework. By selectively removing or modifying specific modules, the study evaluates how each component affects overall system performance in medical question answering. The analysis focuses on entity linking accuracy, reasoning quality, and answer correctness across benchmark medical datasets. First, a baseline configuration using LLM only Medical QA is evaluated. This setup relies solely on unstructured text based reasoning and exhibits lower accuracy due to hallucinations and ambiguous medical entity interpretation. When a Medical Knowledge Graph is incorporated without entity fast linking, performance improves, highlighting the importance of structured clinical knowledge for factual grounding. Introducing entity fast linking further enhances accuracy by resolving ambiguity in medical terminology and ensuring precise alignment with standardized clinical concepts. Finally, enabling the full LLM-KGMQA system with multihop reasoning, knowledge fusion, and semantics based pruning achieves the highest performance across all evaluation metrics. The removal of semantic pruning leads to increased reasoning noise and reduced efficiency, while disabling multihop reasoning limits the system’s ability to answer complex clinical queries. Overall, the ablation results demonstrate that each component plays a critical role, and full integration of all modules is essential to achieve optimal accuracy, robustness, and interpretability in medical question answering.

Table 8 Ablation Study Results of LLM-KGMQA Components Across Accuracy, Reasoning, and Efficiency Metrics

Configuration	Entity Link. (%)	Intent Acc. (%)	Task Comp. (%)	Latency (ms)
LLM-only Medical QA	90.2	84.1	41.2	880
LLM + Medical KG	93.6	88.5	79.4	720
LLM + KG + Fast-Linking	96.8	91.9	85.6	820
LLM-KGMQA (Full System)	98.0	94.3	88.0	960

As shown in table 8, the ablation study evaluates the contribution of key components in LLM-KGMQA, including medical knowledge graph integration, entity fast linking, and multihop reasoning. The results demonstrate that progressively incorporating

these modules leads to consistent improvements in entity linking accuracy, intent understanding, and task completion rate, with the complete configuration achieving the highest overall performance while maintaining acceptable inference latency.

6 Conclusions and Future Work

This study investigated whether tightly integrating Large Language Models (LLMs) with structured medical knowledge graphs can overcome the core limitations of traditional medical question answering systems, particularly in terms of accuracy, interpretability, and clinical reliability. The proposed LLM-KGMQA framework successfully validates this objective through comprehensive experimental evaluation, achieving up to 98% overall answer accuracy, strong multihop medical reasoning performance, and 96.8% factual consistency across multiple standardized medical benchmark datasets. These results demonstrate that grounding LLM reasoning in verified clinical knowledge substantially improves reliability compared to stand alone language model based approaches. By leveraging standardized medical terminologies such as ICD-10, ICD-11, and SNOMED-CT, LLM-KGMQA effectively mitigates hallucinations while providing transparent and explainable medical responses. The modular architecture comprising entity fast linking, dynamic nhop subgraph construction, knowledge fusion, and semantics based pruning enables accurate and interpretable inference over complex clinical relationships. This design further supports multilingual medical inputs, ensuring consistent reasoning across language variations and diverse clinical contexts. The observed gains in Exact Match, F1 score, nDCG, and low latency inference confirm the framework’s suitability for real time clinical decision support applications. Despite these strengths, certain challenges remain. Ambiguous or underspecified medical queries, evolving clinical guidelines, incomplete knowledge graph coverage, and cross lingual terminology variations can still impact system performance in real world healthcare environments. Additionally, scalability to deeper reasoning depths and long term knowledge maintenance requires further attention. Future work will focus on extending LLM-KGMQA to support higher hop medical reasoning, improving query disambiguation through contextual clarification strategies, enhancing multilingual robustness, enabling continuous and automated knowledge graph updates, strengthening clinical risk-awareness mechanisms, and optimizing inference efficiency for large scale deployment. Further alignment with emerging healthcare AI governance and ethical standards will also be explored. These advancements aim to position LLM-KGMQA as a robust, scalable, and trustworthy medical question answering platform for real world clinical and healthcare information systems.

Acknowledgements

This research was conducted without any external funding support.

Author Contributions

Veera Babu Reddy *[1] conceived and designed the LLM-KGMQA framework, supervised the overall research workflow, and led the system architecture design, integration,

and result validation. Maneesha Shaik [2] contributed to medical data preprocessing, knowledge graph construction, and entity fast linking implementation. Mahalakshmi Chadalaawada [3] was responsible for multilingual query handling, LLM reasoning module optimization, and performance tuning. Geetha Varshini Palla [4] conducted experimental evaluation, benchmark testing, and statistical significance analysis. Gnanesh Surangi [5] assisted with system testing, documentation, result organization, and manuscript preparation. Venkata Sai Teja Vadalasetty [6] provided academic guidance, methodological review, and technical oversight throughout the development and evaluation of the framework. All authors reviewed and approved the final manuscript.

Data Availability Statement

This study utilized publicly available datasets obtained from established open sources, as referenced within the manuscript.

Declarations

Competing Interests The authors declare no competing interests.

References

- [1] Kim J, Park S (2022) Maximization and restoration: Action segmentation through dilation passing and temporal reconstruction *Pattern Recognition*. DOI: [10.1016/j.patcog.2022.108764](https://doi.org/10.1016/j.patcog.2022.108764)
- [2] Staab S, Studer R (2010) *Handbook on Ontologies*. Springer, Berlin. <https://link.springer.com/book/10.1007/978-3-540-92673-3>
- [3] Hogan A, Blomqvist E, Cochez M et al (2021) Knowledge graphs. *Semantic Web*. DOI: [10.1007/978-3-031-01918-0](https://doi.org/10.1007/978-3-031-01918-0)
- [4] Singhal K, Azizi S, Tu T et al (2023) Large language models encode clinical knowledge. *Nature*. DOI: [10.1038/s41586-023-06291-2](https://doi.org/10.1038/s41586-023-06291-2)
- [5] Dörpinghaus J (2022) Context mining and graph queries on giant biomedical knowledge graphs. *Knowledge and Information Systems*. DOI: [10.1007/s10115-018-1280-6](https://doi.org/10.1007/s10115-018-1280-6)
- [6] Ibrahim N, Aboulela S, Ibrahim A (2024) A survey on augmenting knowledge graphs with large models. *Cognitive Computation*. DOI: [10.1007/s44163-024-00175-8](https://doi.org/10.1007/s44163-024-00175-8)
- [7] Ghnemat R, Saleh A (2025) Large Language Models for Clinical AI in Healthcare: A Systematic Review. *Cognitive Computation*. DOI: [10.1007/s44163-025-00784-x](https://doi.org/10.1007/s44163-025-00784-x)
- [8] Wu Z, Pan S, Chen F et al (2021) A comprehensive survey on graph neural networks. *IEEE Trans Neural Netw Learn Syst*. DOI: [10.1109/TNNLS.2020.2978386](https://doi.org/10.1109/TNNLS.2020.2978386)

- [9] Vivar G, Kazi A, Burwinkel H, Zwergal A, Navab N, Ahmadi SA (2018) Simultaneous imputation and classification using Multigraph Geometric Matrix Completion (MGMC): Application to neurodegenerative disease classification. *IEEE Transactions on Medical Imaging*. <https://www.sciencedirect.com/science/article/abs/pii/S0933365721000907>
- [10] Chen Q, Kang D, He Y, Chang TH, Liu YF Joint Power and Admission Control Based on Channel Distribution Information: A Novel Two-Timescale Approach. *IEEE Transactions on Signal Processing*. <https://ieeexplore.ieee.org/document/7812602>
- [11] Assayed SK, Shieh CS, Gupta SK (2025) Engineering intelligent healthcare systems: understanding medical queries with AI and NLP. *J Eng Appl Sci*. DOI: [10.1186/s44147-025-00800-y](https://doi.org/10.1186/s44147-025-00800-y)
- [12] Li A et al (2024) KGSCS — a smart care system for elderly with geriatric chronic diseases using a knowledge graph. *BMC Medical Informatics and Decision Making*. DOI: <https://doi.org/10.1186/s12911-024-02472-9>
- [13] Labarga A (2026) An Interpretable Graph Neural Network for Multi-omics Data Integration and Biomarker Discovery. In: *Artificial Intelligence in Biomedicine (CIABiomed 2025)*, pp 438–448. https://link.springer.com/chapter/10.1007/978-3-032-10661-2_33?fromPaywallRec=true
- [14] Wang D, Liu H (2024) Large language models in medical and healthcare fields: review. *Cognitive Computation*. DOI: [10.1007/s10462-024-10921-0](https://doi.org/10.1007/s10462-024-10921-0)
- [15] Ruan T, Huang Y, Liu X et al (2019) QAnalysis: A question-answer driven analytic tool on knowledge graphs for leveraging EMRs. *BMC Medical Informatics and Decision Making*. DOI: <https://doi.org/10.1186/s12911-019-0937-y>
- [16] Oulefki S, Berkani L, Boudjenah N, Bellatreche L, Mokhtari A (2025) BioGITOM: Matching Biomedical Ontologies with Graph Isomorphism Transformer. *The VLDB Journal* 34:65. <https://link.springer.com/article/10.1007/s00778-025-00943-7>
- [17] Kacupaj E, Singh K, Maleshkova M et al (2024) Conversational QA over knowledge graphs. In: *Lecture Notes in Computer Science*. DOI: [10.1007/978-3-031-64451-1_9](https://doi.org/10.1007/978-3-031-64451-1_9)
- [18] Gao P, Yang J (2024) Medical knowledge graph QA for drug-drug interactions. *IET Communications*. DOI: [10.1049/cit2.12332](https://doi.org/10.1049/cit2.12332)
- [19] Wang Q, Wang B, Guo L (2017) Knowledge graph embedding: A survey of approaches and applications. *IEEE Trans Knowl Data Eng*. DOI: [10.1109/TKDE.2017.2754499](https://doi.org/10.1109/TKDE.2017.2754499)

- [20] Leão T, Madeira SC, Gromicho M, de Carvalho M, Carvalho AM (2021) Learning dynamic Bayesian networks from time-dependent and time-independent data: Unraveling disease progression in Amyotrophic Lateral Sclerosis. *Artificial Intelligence in Medicine*. <https://www.sciencedirect.com/science/article/pii/S1532046421000599?via%3Dihub>
- [21] Jiang Z, Chi C, Zhan Y (2021) Research on medical question answering system based on knowledge graph. *IEEE Access*. DOI: [10.1109/ACCESS.2021.3055371](https://doi.org/10.1109/ACCESS.2021.3055371)
- [22] Harnoune Y, Rhanoui M et al (2021) BERT-based clinical knowledge extraction for biomedical knowledge graph construction. *Comput Methods Programs Biomed Update*. DOI: [10.1016/j.cmpbup.2021.100042](https://doi.org/10.1016/j.cmpbup.2021.100042)
- [23] Cui H et al (2025) A review on healthcare knowledge graphs: resources and applications. *Computers in Biology and Medicine*. DOI: <https://doi.org/10.1016/j.jbi.2025.104861>
- [24] Aguzzi G, Magnini M, Farahmand A et al (2025) RAG-enhanced open SLMs for hypertension management chatbots. *Journal of Medical Systems*. DOI: [10.1007/s10916-025-02297-7](https://doi.org/10.1007/s10916-025-02297-7)
- [25] Ji Z et al (2023) Survey of hallucination in natural language generation. *ACM Comput Surv*. DOI: [10.1145/3571730](https://doi.org/10.1145/3571730)
- [26] World Health Organization (2026) International Classification of Diseases (ICD). <https://www.who.int/standards/classifications/classification-of-diseases>
- [27] Lu X, Tu SW, Chang W, Wan C, Wang J, Zang Y, Ramdas B, Kapur R, Lu X, Cao S et al (2021) SSMD: a semi-supervised approach for a robust cell type identification and deconvolution of mouse transcriptomics data. *Briefings in Bioinformatics* 22(4):bbaa307. <https://academic.oup.com/bib/article/22/4/bbaa307/5998844>
- [28] Bodenreider O (2004) The Unified Medical Language System (UMLS). *Nucleic Acids Research*. DOI: [10.1093/nar/gkh061](https://doi.org/10.1093/nar/gkh061)
- [29] Tanioka H (2025) Towards safe and trustworthy healthcare AI. *Human-Centric Intelligent Systems*. DOI: [10.1007/s44230-025-00131-4](https://doi.org/10.1007/s44230-025-00131-4)
- [30] Song Y, Sun X et al (2023) Advancements in complex knowledge graph QA. *Electronics*. DOI: [10.3390/electronics12214395](https://doi.org/10.3390/electronics12214395)
- [31] Farrugia L et al (2025) medicX-KG: a knowledge graph for pharmacists' drug information supporting query answering. *Journal of Biomedical Semantics*. DOI: <https://doi.org/10.1186/s13326-025-00332-7>

- [32] Linders et al (2025) Knowledge graph-extended retrieval augmented generation for improved QA. *Applied Intelligence*. DOI: <https://doi.org/10.1007/s10489-025-06885-5>