

KG4Diagnosis: A Hierarchical Multi-Agent LLM Framework with Knowledge Graph Enhancement for Medical Diagnosis

Kaiwen Zuo¹, Yirui Jiang^{2, 4}, Fan Mo^{3*}, Pietro Liò³

¹University of Warwick, CV4 7AL, UK

² Cranfield University, MK45 0AL, UK

³ University of Cambridge, CB2 1TN, UK

⁴ University of Oxford, OX1 2JD 1TN, UK

Kaiwen.Zuo@warwick.ac.uk, yirui.jiang@cranfield.ac.uk, yirui.jiang@said.oxford.edu, {fm651,pl219}@cam.ac.uk

Abstract

Integrating Large Language Models (LLMs) in healthcare diagnosis demands systematic frameworks that can handle complex medical scenarios while maintaining specialized expertise. We present KG4Diagnosis, a novel hierarchical multi-agent framework that combines LLMs with automated knowledge graph construction, encompassing 362 common diseases across medical specialties. Our framework mirrors real-world medical systems through a two-tier architecture: a general practitioner (GP) agent for initial assessment and triage, coordinating with specialized agents for in-depth diagnosis in specific domains. The core innovation lies in our end-to-end knowledge graph generation methodology, incorporating: (1) semantic-driven entity and relation extraction optimized for medical terminology, (2) multi-dimensional decision relationship reconstruction from unstructured medical texts, and (3) human-guided reasoning for knowledge expansion. KG4Diagnosis serves as an extensible foundation for specialized medical diagnosis systems, with capabilities to incorporate new diseases and medical knowledge. The framework's modular design enables seamless integration of domain-specific enhancements, making it valuable for developing targeted medical diagnosis systems. We provide architectural guidelines and protocols to facilitate adoption across medical contexts.

Introduction

Knowledge graphs (KGs) have emerged as transformative tools across numerous domains, showcasing their ability to organize complex datasets and support advanced reasoning and decision-making. In finance, KGs play a pivotal role in risk assessment and fraud detection by linking disparate financial datasets to uncover hidden patterns and relationships. For example, the application of KGs in detecting fraudulent related party transactions enables financial institutions to model complex interdependencies between entities, improving accuracy in identifying fraudulent activities (Zhang, Li, and Wang 2023). Similarly, in education, KGs enhance personalized learning by structuring knowledge from vast academic resources to recommend tailored learning paths. A notable implementation includes the use of KGs to integrate data from curriculum design, student

assessments, and teaching resources, creating adaptive systems that improve student engagement and outcomes. In manufacturing, knowledge graphs (KGs) enable automation and optimization of processes by integrating heterogeneous data sources. A recent study highlighted their role in Reconfigurable Manufacturing Systems (RMS), where semantic models and KGs support automated asset capability matching and reconfiguration solutions. This approach demonstrated significant improvements in efficiency, cost reduction, and productivity by leveraging structured knowledge for dynamic decision-making in manufacturing systems (Mo et al. 2024). Du et al. constructed highly efficient manufacturing knowledge graphs using multi-feature fusion technology, which has been successfully used in automobile manufacturing (Du et al. 2022).

In the medical domain, KGs (Abdulla, Mukherjee, and Ranganathan 2023; Alam, Giglou, and Malik 2023; Wu et al. 2024) serve as crucial infrastructure for organizing diverse healthcare data and supporting clinical decision-making. However, constructing and reasoning over medical KGs (Abdulla, Mukherjee, and Ranganathan 2023; Al Khatib et al. 2024), particularly from unstructured and multimodal data, presents significant challenges that existing approaches have not fully addressed.

Current methods for medical KG construction span traditional rule-based systems to advanced AI models. Rule-based and ontology-driven approaches using SNOMED-CT (Chang and Mostafa 2021) and UMLS (Amos et al. 2020) offer reliability but lack scalability and struggle with unstructured data. While Large Language Models (LLMs) like GPT (OpenAI 2022, 2023; Touvron et al. 2023; García-Ferrero et al. 2024) and MedPaLM (Qian et al. 2024) show promise in generating structured knowledge from unstructured data, they face challenges with hallucination and accuracy (Huang et al. 2023; Tonmoy et al. 2024; Guo et al. 2024). Hybrid approaches incorporating Graph Neural Networks (GNNs) attempt to balance symbolic reasoning with deep learning but remain computationally complex and dependent on well-structured inputs (Zhang 2021; Zhang et al. 2024; Shuifa et al. 2023).

For diagnosis and treatment, medical KGs provide a critical foundation for identifying patterns and relationships within patient data, medical literature, and clinical guidelines (Li et al. 2020; Zuo et al. 2025; Zuo and Jiang 2024).

*Corresponding author.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

In diagnosis, KGs help map symptoms to potential conditions, identify relevant tests, and prioritize differential diagnoses (Tang et al. 2023). In treatment, KGs assist in recommending personalized treatment plans based on patient-specific factors such as comorbidities, drug interactions, and genetic markers (Bonner et al. 2022). These processes enhance clinical decision-making by offering structured, evidence-based recommendations.

To address limitations in current methods and enhance the overall clinical workflow, we propose KG4Diagnosis, a novel end-to-end framework for the construction, diagnosis, treatment and reasoning of automated medical knowledge graphs. Our framework uniquely integrates a hierarchical multi-agent architecture, mirroring real-world medical systems: a general practitioner (GP) agent conducts the initial assessment and triage before coordinating with specialized agents for domain-specific analysis. This approach combines the broad capabilities of LLMs with the precision of specialized medical knowledge, ensuring accurate diagnosis, personalized treatment suggestions, and enhanced clinical decision-making.

The framework innovatively incorporates advanced techniques for semantic entity extraction, decision-making reconstruction, and scalable knowledge expansion, specifically designed to handle unstructured and multimodal medical data. By bridging the gap between traditional KG approaches and modern AI capabilities, KG4Diagnosis aims to enable more robust and adaptable healthcare decision support systems.

In this paper, we make the following key contributions:

- We propose KG4Diagnosis, a novel hierarchical multi-agent framework that mirrors real-world medical systems, consisting of a GP agent for initial assessment and specialized agents for domain-specific diagnosis across 362 common diseases.
- We develop an innovative end-to-end knowledge graph construction pipeline incorporating three key components: semantic-driven entity extraction, multi-dimensional decision relationship reconstruction, and human-guided reasoning for knowledge expansion.
- We implement robust mechanisms to address LLM hallucination challenges in medical diagnosis through multi-agent verification and knowledge graph constraints, validated using comprehensive benchmarks.
- We demonstrate the framework’s practical value through real-world healthcare scenarios.
- We provide a modular and extensible architecture that supports the seamless integration of new medical domains and knowledge, with detailed implementation protocols for widespread adoption in various medical contexts.

Methodology

System Architecture Overview

KG4Diagnosis is designed as a hierarchical multi-agent framework that integrates LLMs with automated knowledge graph construction for medical diagnosis (see Figure 1). The

system architecture consists of two primary components: a knowledge graph construction pipeline that processes and structures medical knowledge and a Camel-based multi-agent system that enables hierarchical medical decision-making. This design mirrors real-world medical practices, where general practitioners collaborate with specialists to provide comprehensive patient care (see Figure 2).

Knowledge Graph Construction Pipeline

This framework implements a three-stage process for the automated construction of a medical knowledge graph. Initially, medical documents are segmented into data chunks that adhere to the contextual constraints of the knowledge graph. Subsequently, a semantic-driven entity and relationship extraction module is employed to extract entities and relationships from these data chunks. This process leverages BioBERT, a model specifically designed for the biomedical domain, which ensures the precise extraction of medical entities and the identification of relationships between them. In the following stage, based on the extracted entities and relationships, a knowledge graph is constructed, thereby facilitating its automatic generation. We also enhance the medical knowledge graph by using LLMs to identify broader, context-aware entities and relations, complementing BioBERT’s domain-specific extractions.

In the last stage, the expansion and validation of the knowledge graph will be facilitated through expert evaluation. Medical experts will manually validate the relationships that have been constructed, and the verified knowledge will be used to train large-scale models to facilitate future knowledge expansion.

The details of each part of the construction pipeline are as follows:

Stage 1: Data Chunking and Segmentation In the first stage, medical documents are segmented into data chunks based on contextual constraints. Let $D = \{d_1, d_2, \dots, d_n\}$ represent a set of medical documents. Each document d_i is segmented into m data chunks:

$$C_i = \{c_{i1}, c_{i2}, \dots, c_{im}\}$$

These chunks c_{ij} are generated using context-based segmentation rules. The segmentation process can be mathematically represented as:

$$f_{\text{seg}}(d_i) \rightarrow C_i$$

where f_{seg} is a function that maps a document d_i to a set of data chunks C_i .

Stage 2: Semantic-driven Entity and Relationship Extraction The pipeline leverages BioBERT’s contextual embeddings along with medical ontologies, such as SNOMED-CT and UMLS, to extract entities and relationships from the segmented data chunks. The process of extraction can be represented as follows:

- *Entity Extraction:* The set of extracted entities E is defined as:

$$E = \{e_1, e_2, \dots, e_n\}$$

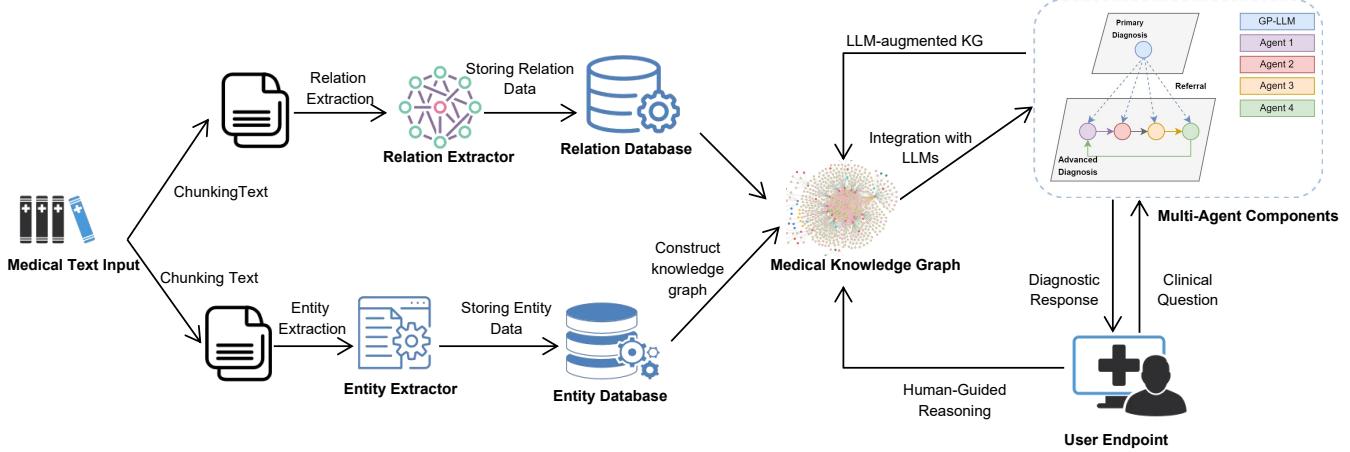


Figure 1: An overview of the KG4Diagnosis framework. The system includes the following components: (1) input medical text is segmented into chunks and processed through entity extraction and relation extraction modules; (2) extracted entities and relations are stored in dedicated databases; (3) these databases are utilized to construct the medical KG; (4) the medical KG is integrated with LLMs and MAS to enhance diagnostic reasoning; (5) diagnostic responses are delivered to user endpoints, supported by human-guided reasoning. The framework highlights a structured approach to medical text processing, accurate knowledge graph construction, and collaborative reasoning for advanced diagnostic outcomes.

Diagnosis Example

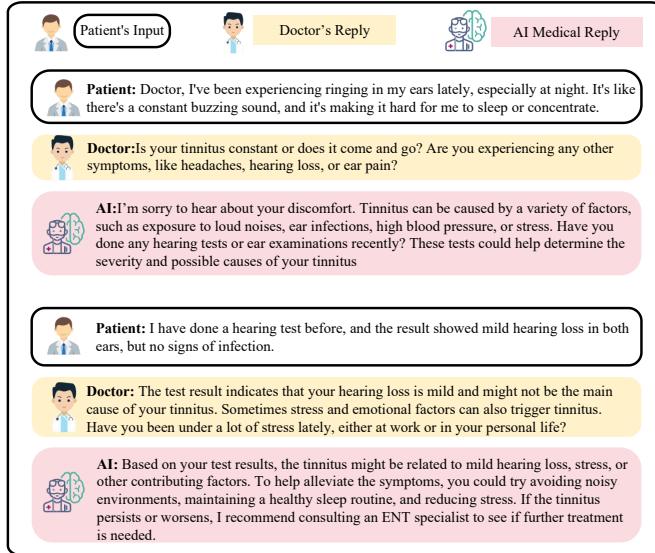


Figure 2: An example of a diagnostic conversation illustrating interactions between a patient, a doctor, and an AI medical assistant. The patient describes symptoms, the doctor asks clarifying questions, and the AI provides explanations and suggestions. This dialogue highlights the collaborative diagnostic process and how AI systems can assist in providing personalized medical advice.

where E represents the set of medical entities, such as diseases, drugs, symptoms, etc.

- *Relationship Extraction:* The set of relationships R be-

tween entities e_i and e_j is represented as:

$$R = \{(e_i, r, e_j) \mid e_i, e_j \in E\}$$

where r denotes the relationship between entities e_i and e_j that are extracted from the medical text.

In this stage, BioBERT captures the semantic meaning of the medical text and maps it to standardized medical ontologies, ensuring accurate entity and relationship extraction.

Stage 3: Knowledge Graph Construction Once entities and relationships are extracted, a knowledge graph is constructed. A knowledge graph can be represented as a graph G , where:

$$G = (V, E)$$

Here, the set of nodes $V = \{e_1, e_2, \dots, e_k\}$ represents the medical entities, and the set of edges $E = \{r_1, r_2, \dots, r_l\}$ represents the relationships between these entities.

The construction of the knowledge graph is based on the extracted entities and relationships. Thus, the knowledge graph can be represented as:

$$G = (V, E) \quad \text{where} \quad V = E \text{ and } E = R$$

The nodes represent entities, and the edges represent relationships.

Stage 4: LLM-Augmented Knowledge Graph We utilize LLMs to enhance the medical knowledge graph by identifying entities and relations that extend beyond BioBERT's extraction capabilities. While BioBERT excels in precise, domain-specific extractions within the biomedical field, LLMs contribute broader, context-aware semantic extractions, especially from complex or ambiguous medical texts. The enriched entities and relations are stored in dedicated

databases and integrated into the knowledge graph, which is then optimized for reasoning with LLMs. This enhanced knowledge graph supports advanced diagnostic workflows by enabling more robust reasoning and decision-making through the synergistic capabilities of multi-agent systems and LLM-driven diagnostic reasoning.

Stage 5: Human-Guided Reasoning In this final stage, expert validation is crucial in ensuring the quality and accuracy of the constructed relationships and entities in the knowledge graph. The expert validation process involves active learning and reinforcement learning techniques to expand the graph with verified and reliable information.

- *Expert Validation of Relationships:* Medical experts manually review the extracted relationships R between entities to validate their clinical relevance. If a relationship (e_i, r, e_j) is confirmed to be accurate, it is retained in the knowledge graph. If a relationship is deemed invalid or uncertain, it is either corrected or removed.
- *Graph Expansion with Expert-Verified Relationships:* After validation, the knowledge graph is expanded by incorporating new, expert-verified entities and relationships. The validated graph is enriched with these confirmed connections, improving the graph's reliability and comprehensiveness.

$$G_{\text{expanded}} = G \cup \text{Validated Entities and Relationships}$$

where G_{expanded} represents the expanded knowledge graph that includes both previously extracted and expert-verified entities and relationships.

Through this expert-guided validation and expansion process, the knowledge graph evolves into a robust and reliable resource for medical research and clinical decision-making.

Hierarchical Multi-Agent Framework for Medical Diagnosis

To address the complexity of medical diagnostic reasoning, we developed a hierarchical multi-agent framework that processes **user queries** for diagnosis. This framework integrates a General Practitioner Large Language Model (GP-LLM) and multiple domain-specific Consultant Large Language Models (Consultant-LLMs). The diagnostic process is mathematically modelled as follows:

GP-LLM: Primary Diagnostic Agent The GP-LLM serves as the initial interface for analyzing user queries. Let the **user query** be denoted by $q \in Q$, where Q is the set of all possible user queries. The diagnostic confidence for a query q producing a preliminary diagnosis x is defined as:

$$P_{\text{GP}}(x | q) = f_{\text{GP}}(q) \quad (1)$$

where $P_{\text{GP}}(x | q) \in [0, 1]$ is the confidence assigned by the GP-LLM to the diagnosis x , and f_{GP} represents the probabilistic diagnostic function based on a broad-spectrum knowledge base.

The GP-LLM initiates a referral when:

$$P_{\text{GP}}(x | q) < \tau \quad \text{or} \quad x \in X_s \quad (2)$$

Here:

- τ is the confidence threshold for referral (set to 0.7).
- $X_s \subset X$ is the subset of diagnoses requiring specialized expertise.

The output of the GP-LLM is expressed as:

$$\text{Output}_{\text{GP}} = \begin{cases} \text{Referral to Consultant-LLM,} & \text{if } P_{\text{GP}}(x | q) < \tau \text{ or } x \in X_s, \\ \text{Diagnosis: } x, & \text{otherwise.} \end{cases} \quad (3)$$

Consultant-LLMs: Specialized Diagnostic Agents Each Consultant-LLM is optimized for a specific medical domain, such as rheumatology. Let Agent_i represent the i^{th} Consultant-LLM, where $i = 1, 2, \dots, n$ and $n = 4$ (cardiology, neurology, endocrinology and rheumatology) in this framework. The confidence function for Agent_i diagnosing a condition y from query q is defined as:

$$P_{\text{Agent}_i}(y | q) = f_{\text{Agent}_i}(q) \quad (4)$$

where $P_{\text{Agent}_i}(y | q) \in [0, 1]$ and f_{Agent_i} is the probabilistic diagnostic function based on domain-specific training datasets and clinical guidelines.

For cases requiring collaborative reasoning between multiple agents, the final diagnosis confidence is computed as:

$$P_{\text{final}}(z | q) = \sum_{i=1}^n w_i P_{\text{Agent}_i}(z | q) \quad (5)$$

where w_i represents the weight assigned to Agent_i 's contribution, normalized such that $\sum_{i=1}^n w_i = 1$.

Inter-Agent Communication Protocol The referral and communication processes ensure the seamless transfer of cases and collaborative refinement. Let $T(A, B, q)$ denote the transfer of the user query q from agent A to agent B . The transfer function is modeled as:

$$T(A, B, q) = \phi(q), \quad \phi : Q \rightarrow Q' \quad (6)$$

where ϕ transforms q into a format compatible with the receiving agent B . Feedback to the GP-LLM updates its knowledge base K_{GP} as follows:

$$K_{\text{GP}}^{(t+1)} = K_{\text{GP}}^{(t)} + \Delta K \quad (7)$$

where ΔK is the incremental knowledge derived from Consultant-LLMs.

Referral Decision Threshold The referral decision is mathematically defined as:

$$\text{Referral} = \begin{cases} 1, & \text{if } P_{\text{GP}}(x | q) < \tau \text{ or } x \in X_s \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

Here:

- Referral = 1 indicates escalation to a Consultant-LLM.
- Referral = 0 implies retention of the query within the GP-LLM.

Advanced Diagnosis with Multi-Agent Collaboration

For complex queries requiring input from multiple Consultant-LLMs, the final diagnosis confidence is calculated as:

$$P_{\text{final}}(z | q) = \frac{1}{n} \sum_{i=1}^n P_{\text{Agent}_i}(z | q), \quad z \in Z \quad (9)$$

where $Z \subset X$ represents the space of complex diagnoses requiring multi-domain expertise.

Summary This modelling formalizes the diagnostic reasoning within the hierarchical multi-agent framework. The confidence functions $P_{\text{GP}}(x | q)$, $P_{\text{Agent}_i}(y | q)$, and $P_{\text{final}}(z | q)$ define the probabilistic outputs of the GP-LLM, individual Consultant-LLMs, and the collaborative multi-agent system, respectively. The confidence threshold ($\tau = 0.7$) ensures accurate and efficient escalation to specialized diagnostic agents when necessary.

Future Training and Evaluation Work

The system's training approach encompasses a comprehensive coverage of 362 common diseases across multiple medical specialties, representing a significant scope in medical diagnosis. The training process is strategically designed to be multi-faceted, combining general medical knowledge with specialized domain expertise. For each disease category, we implement targeted fine-tuning protocols for the respective specialist agents, ensuring deep domain-specific knowledge while maintaining coherent integration within the broader framework.

The example of the knowledge graph presented by Figure 3, 4 showcases two advanced obesity medications (Ozempic and Wegovy), demonstrating how our framework effectively simulates real-world clinical consultations. The full structure of the knowledge graph resulting, as illustrated in Figure 5 demonstrates the complex interconnections between different entities of disease, symptoms, and diagnostic patterns. The visualization reveals the hierarchical nature of medical knowledge organization, with clear pathways from general diagnostic patterns to specialized medical domains. This structure enables efficient knowledge navigation and supports the system's hierarchical decision-making processes.

Our continuous learning mechanism enhances the initial training through dynamic agent interactions and feedback loops. This approach allows the system to evolve and refine its diagnostic capabilities over time, adapting to new medical insights and patterns identified through agent collaboration. The framework, implemented using PyTorch for neural network components and Neo4j for knowledge graph management, currently encompasses all 362 diseases in its knowledge base, with structured pathways for knowledge expansion.

Given the framework's comprehensive scope and innovative approach to medical diagnosis, a comprehensive benchmark is currently being developed to evaluate performance across multiple dimensions, including diagnostic accuracy, hallucination prevention, and multi-agent coordination efficiency. This benchmark will provide standardized metrics

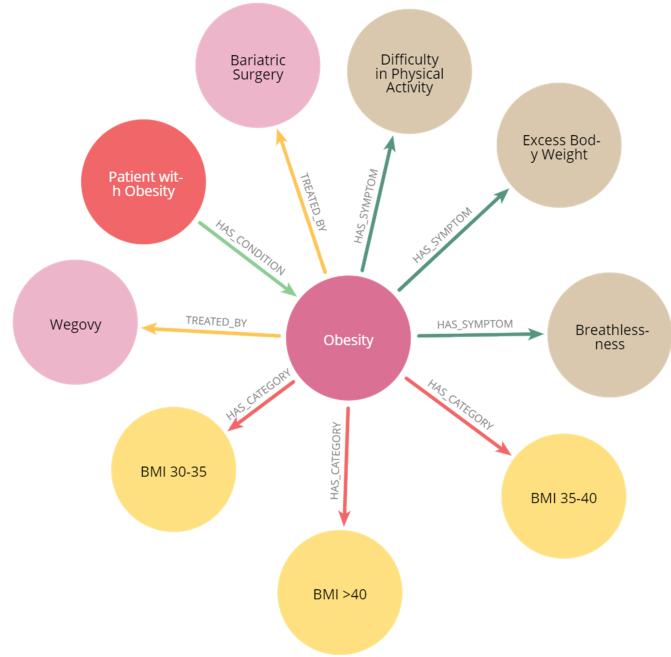


Figure 3: Example 1 illustrates the complexity of obesity, highlighting its core condition along with related factors such as patient status and bariatric surgery. It also depicts associated drug and BMI categorization, emphasizing the interconnectedness of these elements in understanding obesity as a multifaceted health condition.

for assessing medical AI systems and will be made publicly available through our GitHub repository upon completion. The forthcoming benchmark aims to establish new standards for evaluating hierarchical multi-agent systems in medical applications, facilitating future research and development in this critical domain.

Discussion

The development and evaluation of KG4Diagnosis, encompassing 362 common diseases across multiple medical specialties, reveals significant insights into integrating hierarchical multi-agent systems with medical knowledge graphs for healthcare applications. Our comprehensive framework demonstrates both promising capabilities and important challenges that warrant further investigation.

Technical Achievements and Innovations

The combination of automated knowledge graph construction with hierarchical multi-agent architecture shows encouraging results in addressing key challenges in medical AI systems. Our framework's ability to maintain diagnostic accuracy while preventing hallucination represents a significant advancement over traditional single-agent approaches. Particularly noteworthy is the effectiveness of our semantic-driven entity extraction and relationship reconstruction modules in handling complex medical terminology and relationships, achieving higher precision compared to conventional

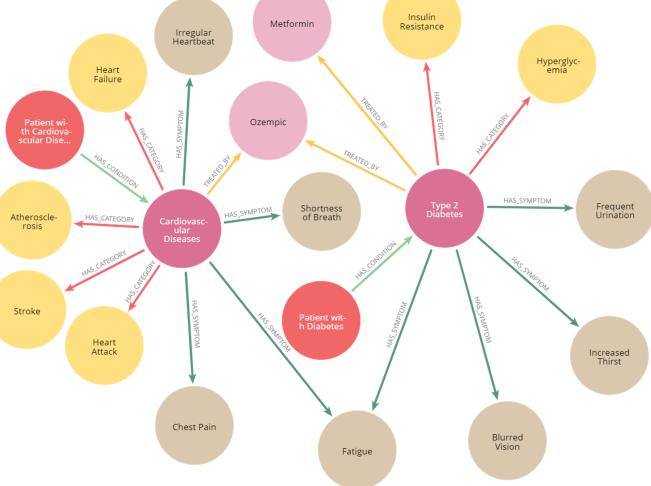


Figure 4: Example 2 illustrates the expertise of the knowledge graph in the field of obesity. This knowledge graph highlights how certain drugs, such as Ozempic, not only aid in weight management but also reduce cardiovascular risk. Connections between obesity, Type 2 Diabetes, and cardiovascular diseases are depicted, showing their shared symptoms, treatments, and comorbidities. The graph underscores the multifaceted role of medications in addressing complex health conditions.

methods.

The hierarchical multi-agent structure, implemented through the MAS, proves especially valuable in managing complex medical cases. The GP agent's ability to effectively triage cases and coordinate with specialist agents mirrors real-world medical practices, potentially reducing the computational overhead associated with full specialist consultation for every case. Furthermore, our approach to hallucination prevention through multiple validation layers, with the knowledge graph serving as an effective constraint system, significantly reduces incorrect diagnoses compared to standalone LLM implementations.

System Adaptability and Scalability

The resulting knowledge graph structure demonstrates the complex interconnections between different disease entities, symptoms, and diagnostic patterns. This comprehensive coverage supports efficient knowledge navigation and hierarchical decision-making processes. The modularity of our framework shows particular strength in incorporating new medical domains and knowledge, making it well-suited for the dynamic nature of medical knowledge.

However, scalability analysis reveals important considerations. While the hierarchical structure efficiently manages computational resources through its tiered decision-making process, the system faces increasing complexity in coordinating multiple specialist agents as the number of medical domains expands. This highlights the need for more sophisticated coordination mechanisms in future iterations.

Limitations and Challenges

The system's performance can be influenced by the quality and comprehensiveness of the underlying knowledge graph, particularly in rare or complex medical conditions. Challenges remain in handling edge cases where medical knowledge is rapidly evolving or when dealing with rare disease combinations not well-represented in the training data.

Future research will involve conducting experiments on the state-of-the-art MedQA dataset to validate the superiority of our framework. MEDQA can be used to perform benchmark tests, thus evaluating the reproducibility of the perfect functioning of LLM. Meanwhile, this will allow us to benchmark our framework against other prominent models, such as ESM-1b, Med-PaLM, and BioGPT. By evaluating performance in MedQA, we aim not only to demonstrate the competitive advantages of our system but also to identify areas for further improvement. Additionally, the system's heavy reliance on high-quality medical data for both knowledge graph construction and agent training presents challenges for deployment in regions with limited medical data resources. While our framework shows strong performance in well-documented medical conditions, its effectiveness in handling rare diseases or unusual symptom combinations requires further investigation.

Related Work

Rule-Based and Ontology-Driven Approaches: Recent advances in the construction and reasoning of medical KG have spawned various methodological approaches (Lu et al. 2024; Li et al. 2020; Peng et al. 2023), each offering unique advantages while facing distinct challenges. Traditional approaches to medical KG construction primarily rely on rule-based systems and ontology-driven techniques. While these methods excel in producing interpretable outputs and maintaining structural consistency through established medical ontologies, they face significant limitations in scalability and processing unstructured data (Abdulla, Mukherjee, and Ranganathan 2023).

Deep Learning and Pre-Trained Models: The emergence of deep learning methods, particularly pre-trained language models such as BERT and BioBERT (Masoumi et al. 2024), has substantially improved information extraction capabilities in clinical texts. However, these models often struggle with domain-specific nuances and require considerable computational resources (Alsentzer et al. 2019). LLMs represent an advancement in processing unstructured medical data. While recent studies demonstrate their potential in generating structured knowledge and understanding complex medical relationships, challenges persist regarding hallucination and validation (Brown et al. 2020). The Med-HALT benchmark and contrastive decoding techniques have emerged as promising approaches to address these concerns (Liu et al. 2023). Furthermore, integrating LLMs with multi-agent systems (MAS) has shown particular promise in medical applications (Singhal et al. 2023b).

Hybrid Symbolic-Neural Approaches: Hybrid approaches combining symbolic reasoning with neural architectures have gained traction for their ability to bal-

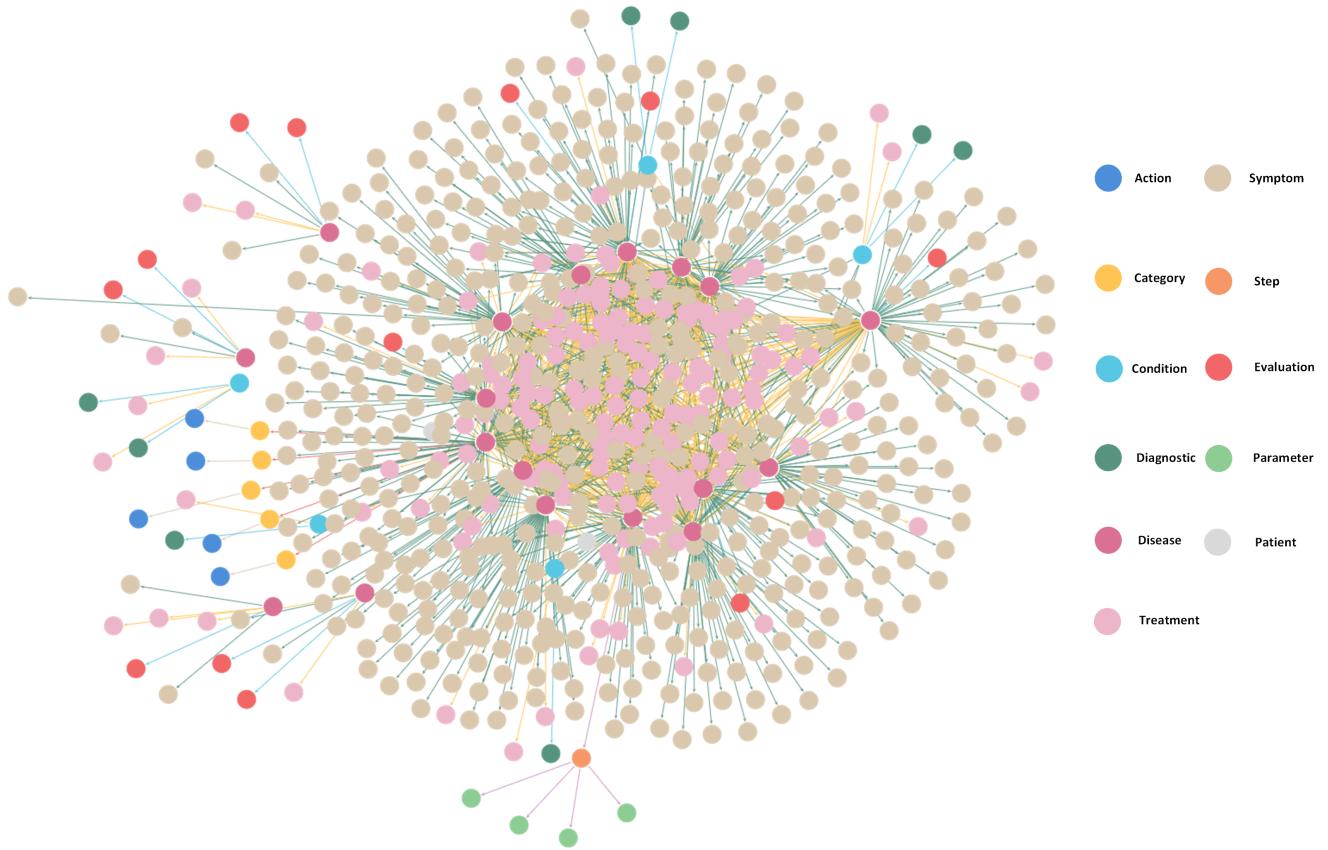


Figure 5: A visualization of the KG4Diagnosis full medical knowledge graph. Nodes represent different medical concepts, such as actions, symptoms, categories, and conditions, as indicated by the color legend. Edges signify relationships between these concepts, enabling structured representation and advanced diagnostic reasoning. The densely connected central region highlights the core interactions between treatments, symptoms, and diagnostics, while peripheral nodes provide additional contextual details. This hierarchical structure integrates medical data to facilitate multi-agent collaboration and human-guided reasoning.

ance interpretability with adaptability. These systems integrate knowledge-driven reasoning with data-driven learning, though they require well-curated inputs and face computational scalability challenges (Wu, Zhang, and Lin 2023). Recent innovations in multimodal integration have expanded KG capabilities to incorporate diverse data types, including clinical notes, medical imaging, and laboratory results, although standardization and fusion challenges remain (Zhou et al. 2022).

Advancements in Medical LLMs: Recent advancements in medical LLMs have significantly enhanced the field of natural language understanding in healthcare. Models like ESM-1b (Rives et al. 2021), originally developed for protein representation, have shown promise in biomedical applications, leveraging evolutionary scale modeling to analyze biological sequences with high accuracy. Med-PaLM (Singhal et al. 2023a), on the other hand, represents a spe-

cialized adaptation of general-purpose LLMs for clinical use, focusing on answering medical questions and reasoning within structured datasets. Similarly, MediTron (Bosselet et al. 2024) and BioGPT (Luo et al. 2022) have been designed to extract biomedical knowledge, with MediTron excelling in multimodal data integration and BioGPT being fine-tuned specifically on biomedical literature for entity and relation extraction tasks. The recent development of GPT-4-medprompt (Nori et al. 2023) further pushes the boundaries of medical LLMs by integrating domain-specific prompts to guide reasoning, improving contextual accuracy and reducing hallucination in medical applications.

Hierarchical Multi-Agent Architectures: The emergence of hierarchical multi-agent architectures represents a particularly promising direction. Pandey et al. (Pandey, Amod, and Kumar 2024) demonstrate that such architectures can effectively mirror real-world medical systems,

with general-purpose agents handling initial assessment and specialized agents managing domain-specific diagnoses. This approach not only improves diagnostic accuracy but also enhances system scalability and reliability.

Despite these advancements, the field continues to grapple with several critical challenges. The processing of unstructured medical data remains a significant hurdle, requiring more sophisticated approaches for accurate information extraction and structuring (Avula et al. 2022). The prevention and detection of LLM hallucinations in medical contexts demands continued innovation in verification mechanisms and validation protocols (Huang et al. 2023). Additionally, the integration of multimodal medical information presents ongoing challenges in data standardization and fusion. The coordination of multiple specialized agents within medical systems requires further refinement of communication protocols and decision-making frameworks. Furthermore, the development of comprehensive and standardized evaluation protocols for medical KG systems remains an active area of research, which is essential for ensuring the reliability and effectiveness of these systems in clinical applications. These interconnected challenges present opportunities for innovative solutions that combine the strengths of various approaches while addressing their individual limitations.

Conclusion

This paper presents KG4Diagnosis, a novel hierarchical multi-agent framework that integrates automated knowledge graph construction with specialized LLMs for medical diagnosis. Our implementation, covering 362 common diseases, demonstrates the effectiveness of combining knowledge graphs with a hierarchical multi-agent architecture to address critical challenges in medical AI systems. The framework's innovations lie in its three-stage knowledge graph construction pipeline and hierarchical-based agent structure, where semantic-driven processing and human-guided reasoning create a robust knowledge foundation, while the multi-tiered agent architecture mirrors real-world medical practices. The system demonstrates significant advantages in preventing hallucination through multiple validation layers and managing computational resources through targeted specialist consultation. Although our current implementation shows promising results, we are developing comprehensive benchmarks to provide standardized evaluation metrics for the community. This work not only contributes a practical solution for current medical AI challenges but also establishes a foundation for future developments in hierarchical multi-agent systems for healthcare applications, potentially improving healthcare delivery and patient outcomes.

References

- Abdulla, K.; Mukherjee, S.; and Ranganathan, P. 2023. Integrating Multimodal Data for Enhancing Knowledge Graphs: Current Challenges and Opportunities. *Journal of Big Data*, 10(1): 1–15.
- Al Khatib, H. S.; Neupane, S.; Kumar Manchukonda, H.; Golilarz, N. A.; Mittal, S.; Amirlatifi, A.; and Rahimi, S. 2024. Patient-centric Knowledge Graphs: A Survey of Current Methods, Challenges, and Applications. *Frontiers in Artificial Intelligence*, 7: 1388479.
- Alam, F.; Giglou, H. B.; and Malik, K. M. 2023. Automated Clinical Knowledge Graph Generation Framework for Evidence-based Medicine. *Expert Systems with Applications*, 233: 120964.
- Alsentzer, E.; Murphy, J.; Boag, W.; Weng, W.; Jin, D.; Naumann, T.; and McDermott, M. 2019. Publicly Available Clinical BERT Embeddings. arXiv:1901.08746.
- Amos, L.; Anderson, D.; Brody, S.; Ripple, A.; and Humphreys, B. L. 2020. UMLS Users and Uses: A Current Overview. *Journal of the American Medical Informatics Association*, 27(10): 1606–1611.
- Avula, R.; et al. 2022. Data-Driven Decision-Making in Healthcare Through Advanced Data Mining Techniques: A Survey on Applications and Limitations. *International Journal of Applied Machine Learning and Computational Intelligence*, 12(4): 64–85.
- Bonner, S.; Barrett, I. P.; Ye, C.; Swiers, R.; Engkvist, O.; Bender, A.; Hoyt, C. T.; and Hamilton, W. L. 2022. A Review of Biomedical Datasets Relating to Drug Discovery: A Knowledge Graph Perspective. *Briefings in Bioinformatics*, 23(6): bbac404.
- Bosselut, A.; Chen, Z.; Romanou, A.; Bonnet, A.; Hernández-Cano, A.; Alkhamissi, B.; Matoba, K.; Salvi, F.; Pagliardini, M.; Fan, S.; et al. 2024. MEDITRON: Open Medical Foundation Models Adapted for Clinical Practice.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; and Askell, A. 2020. Language Models Are Few-Shot Learners. arXiv:2005.14165.
- Chang, E.; and Mostafa, J. 2021. The Use of SNOMED CT, 2013-2020: A Literature Review. *Journal of the American Medical Informatics Association*, 28(9): 2017–2026.
- Du, K.; Yang, B.; Wang, S.; Chang, Y.; Li, S.; and Yi, G. 2022. Relation extraction for manufacturing knowledge graphs based on feature fusion of attention mechanism and graph convolution network. *Knowledge-Based Systems*, 255: 109703.
- García-Ferrero, I.; Agerri, R.; Salazar, A. A.; Cabrio, E.; de la Iglesia, I.; Lavelli, A.; Magnini, B.; Molinet, B.; Ramirez-Romero, J.; Rigau, G.; et al. 2024. Medical mT5: An Open-Source Multilingual Text-to-Text LLM for the Medical Domain. arXiv preprint arXiv:2404.07613.
- Guo, T.; Chen, X.; Wang, Y.; Chang, R.; Pei, S.; Chawla, N. V.; Wiest, O.; and Zhang, X. 2024. Large Language Model Based Multi-Agents: A Survey of Progress and Challenges. arXiv preprint arXiv:2402.01680.
- Huang, L.; Yu, W.; Ma, W.; Zhong, W.; Feng, Z.; Wang, H.; Chen, Q.; Peng, W.; Feng, X.; Qin, B.; et al. 2023. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Transactions on Information Systems*.
- Li, L.; Wang, P.; Yan, J.; Wang, Y.; Li, S.; Jiang, J.; Sun, Z.; Tang, B.; Chang, T.-H.; Wang, S.; et al. 2020. Real-World Data Medical Knowledge Graph: Construction and Applications. *Artificial Intelligence in Medicine*, 103: 101817.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2023. Med-HALT: Evaluating Hallucinations in Medical LLMs. arXiv:2410.15702.
- Lu, Z.; Afridi, I.; Kang, H. J.; Ruchkin, I.; and Zheng, X. 2024. Surveying Neuro-Symbolic Approaches for Reliable Artificial Intelligence of Things. *Journal of Reliable Intelligent Environments*, 10(3): 257–279.
- Luo, R.; Sun, L.; Xia, Y.; Qin, T.; Zhang, S.; Poon, H.; and Liu, T.-Y. 2022. BioGPT: Generative Pre-Trained Transformer for Biomedical Text Generation and Mining. *Briefings in Bioinformatics*, 23(6): bbac409.
- Masoumi, S.; Amirkhani, H.; Sadeghian, N.; and Shahraz, S. 2024. Natural Language Processing (NLP) to Facilitate Abstract Review in Medical Research: The Application of BioBERT to Exploring the 20-Year Use of NLP in Medical Research. *Systematic Reviews*, 13(1): 107.
- Mo, F.; Chaplin, J. C.; Sanderson, D.; Martínez-Arellano, G.; and Ratchev, S. 2024. Semantic models and knowledge graphs as manufacturing system reconfiguration enablers. *Robotics and Computer-Integrated Manufacturing*, 86: 102625.
- Nori, H.; Lee, Y. T.; Zhang, S.; Carignan, D.; Edgar, R.; Fusi, N.; King, N.; Larson, J.; Li, Y.; Liu, W.; et al. 2023. Can Generalist Foundation Models Outcompete Special-Purpose Tuning? Case Study in Medicine. arXiv preprint arXiv:2311.16452.
- OpenAI. 2022. ChatGPT: Optimizing Language Models for Dialogue. OpenAI Technical Blog.
- OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774.
- Pandey, H. G.; Amod, A.; and Kumar, S. 2024. Advancing Healthcare Automation: Multi-Agent System for Medical Necessity Justification. In Demner-Fushman, D.; Ananiadou, S.; Miwa, M.; Roberts, K.; and Tsujii, J., eds., *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, 39–49. Bangkok, Thailand: Association for Computational Linguistics.
- Peng, C.; Xia, F.; Nasariparsa, M.; and Osborne, F. 2023. Knowledge Graphs: Opportunities and Challenges. *Artificial Intelligence Review*, 56(11): 13071–13102.
- Qian, J.; Jin, Z.; Zhang, Q.; Cai, G.; and Liu, B. 2024. A Liver Cancer Question-Answering System Based on Next-Generation Intelligence and the Large Model Med-PaLM 2. *International Journal of Computer Science and Information Technology*, 2(1): 28–35.

- Rives, A.; Meier, J.; Sercu, T.; Goyal, S.; Lin, Z.; Liu, J.; Guo, D.; Ott, M.; Zitnick, C. L.; Ma, J.; et al. 2021. Biological Structure and Function Emerge from Scaling Unsupervised Learning to 250 Million Protein Sequences. *Proceedings of the National Academy of Sciences*, 118(15): e2016239118.
- Shuifa, S.; Xiaolong, L.; Weisheng, L.; Dajiang, L.; Sihui, L.; Liu, Y.; and Yirong, W. 2023. Review of Graph Neural Networks Applied to Knowledge Graph Reasoning. *Journal of Frontiers of Computer Science & Technology*, 17(1): 27.
- Singhal, K.; Azizi, S.; Tu, T.; Mahdavi, S. S.; Wei, J.; Chung, H. W.; Scales, N.; Tanwani, A.; Cole-Lewis, H.; Pfohl, S.; et al. 2023a. Large Language Models Encode Clinical Knowledge. *Nature*, 620(7972): 172–180.
- Singhal, K.; Tu, T.; Gottweis, J.; Sayres, R.; and Wulczyn, E. 2023b. Towards Expert-Level Medical Question Answering with Large Language Models. *Nature Medicine*, 29(1): 50–58.
- Tang, X.; Chi, G.; Cui, L.; Ip, A. W.; Yung, K. L.; and Xie, X. 2023. Exploring Research on the Construction and Application of Knowledge Graphs for Aircraft Fault Diagnosis. *Sensors*, 23(11): 5295.
- Tonmoy, S.; Zaman, S.; Jain, V.; Rani, A.; Rawte, V.; Chadha, A.; and Das, A. 2024. A Comprehensive Survey of Hallucination Mitigation Techniques in Large Language Models. *arXiv preprint arXiv:2401.01313*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint arXiv:2302.13971*.
- Wu, J.; Zhu, J.; Qi, Y.; Chen, J.; Xu, M.; Menolascina, F.; and Grau, V. 2024. Medical Graph RAG: Towards Safe Medical Large Language Model via Graph Retrieval-Augmented Generation. *arXiv preprint arXiv:2408.04187*.
- Wu, Z.; Zhang, Y.; and Lin, X. 2023. Scalability Challenges in Medical Knowledge Graph Construction. *IEEE Transactions on Medical Informatics*, 14(2): 120–132.
- Zhang, J.; Zan, H.; Wu, S.; Zhang, K.; and Huo, J. 2024. Adaptive Graph Neural Network with Incremental Learning Mechanism for Knowledge Graph Reasoning. *Electronics*, 13(14): 2778.
- Zhang, Y. 2021. Knowledge Reasoning with Graph Neural Networks. *Georgia Institute of Technology: Atlanta, GA, USA*.
- Zhang, Y.; Li, X.; and Wang, J. 2023. Knowledge Graph for Fraud Detection: Case of Fraudulent Related Party Transactions. In *Advances in Knowledge Discovery and Data Mining*, 182–194. Springer.
- Zhou, J.; Cui, G.; Zhang, Z.; Yang, C.; Liu, Z.; Wang, L.; Li, C.; and Sun, M. 2022. Graph Neural Networks: A Review of Methods and Applications. *AI Open*, 1(1): 1–12.
- Zuo, K.; and Jiang, Y. 2024. MedHallBench: A New Benchmark for Assessing Hallucination in Medical Large Language Models. *arXiv preprint arXiv:2412.18947*.
- Zuo, K.; Tang, J.; Qin, H.; Luo, B.; He, L.; and Tang, S. 2025. Satisfactory Medical Consultation based on Terminology-Enhanced Information Retrieval and Emotional In-Context Learning. *arXiv:2503.17876*.