

MedKGI: Iterative Differential Diagnosis with Medical Knowledge Graphs and Information-Guided Inquiring

Anonymous ACL submission

Abstract

Recent advancements in Large Language Models (LLMs) have demonstrated significant promise in clinical diagnosis. However, current models struggle to emulate the iterative, diagnostic hypothesis-driven reasoning of real clinical scenarios. Specifically, current LLMs suffer from three critical limitations: (1) generating hallucinated medical content due to weak grounding in verified knowledge, (2) asking redundant or inefficient questions rather than discriminative ones that hinder diagnostic progress, and (3) losing coherence over multi-turn dialogues, leading to contradictory or inconsistent conclusions. To address these challenges, we propose MedKGI, a diagnostic framework grounded in clinical practices. MedKGI integrates a medical knowledge graph (KG) to constrain reasoning to validated medical ontologies, selects questions based on information gain to maximize diagnostic efficiency, and adopts an OSCE-format structured state to maintain consistent evidence tracking across turns. Experiments on clinical benchmarks show that MedKGI outperforms strong LLM baselines in both diagnostic accuracy and inquiry efficiency, improving dialogue efficiency by 30% on average while maintaining state-of-the-art accuracy.

1 Introduction

Large Language Models (LLMs) are increasingly demonstrating their value as powerful tools for clinical diagnosis (Wang et al., 2023; Singhal et al., 2025; Lin et al., 2025). However, real-world clinical reasoning is an iterative process in which doctors need to strategically construct diagnostic hypotheses and gather clinical information in a sequential manner to make the final decision. Current LLMs often struggle in such iterative, hypothesis-driven settings due to a fundamental discrepancy between their probabilistic, token-by-token generation and the systematic rigor required for clinical deduction (Hager et al., 2024).

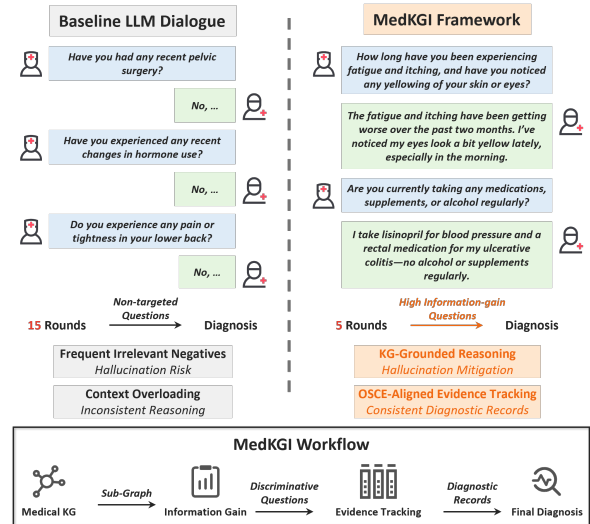


Figure 1: Comparison of diagnostic dialogues and the MedKGI workflow. Left: The dialogues from the baseline LLMs. Right: The dialogues from the proposed MedKGI framework. Bottom: The MedKGI workflow.

This gap leads to several critical limitations: a tendency to produce *hallucinations* by prioritizing plausible patterns over verified knowledge (Zhu et al., 2025); *ineffective questioning* due to the lack of an explicit reasoning framework (Li et al., 2025); and *context overloading* in multi-turn dialogues caused by their associative nature (Savage et al., 2024). As illustrated in Figure 1 (left), these limitations often result in redundant and non-strategic diagnostic dialogues by baseline LLMs, in contrast to the structured trajectory of rigorous medical reasoning.

To bridge this gap, we ground our method in established differential diagnosis frameworks rather than language patterns. *Differential diagnosis* is inherently an iterative process of systematically weighing competing hypotheses against clinical evidence (Hannigen, 2018), which is the opposite of associative LLM reasoning. Accordingly, we integrate three key principles:

064	• Knowledge-Anchored Hypothesis Generation: Inspired by the clinical practice of grounding initial differentials in established medical knowledge (Zuo et al., 2025), we integrate a medical knowledge graph (KG) to generate diagnostic candidates based on verified disease–symptom relationships.	113
065		114
066		115
067		116
068		117
069		118
070		119
071	• Strategic Uncertainty Reduction: Following the differential diagnosis principle of prioritizing high-yield findings, we adopt an information gain-based questioning strategy (Liu et al., 2025) to select the most discriminative questions, thereby minimizing diagnostic uncertainty.	120
072		121
073		122
074		123
075		124
076		125
077		125
078	• Iterative Evidence Refinement: To simulate a doctor’s belief updating process as new evidence emerges, we implement a state-tracking mechanism that maintains a coherent diagnostic record, enabling consistent hypothesis management while mitigating context overloading in long dialogues (Xu et al., 2024).	126
079		127
080		128
081		129
082		130
083		131
084		132
085	Building on these principles, we propose MedKGI, a diagnostic reasoning framework designed to emulate the systematic inquiry of human clinicians within the LLM paradigm. As shown in Figure 1, MedKGI integrates a medical knowledge graph to anchor all diagnostic reasoning in verified medical ontologies, thereby mitigating hallucinations. Building on this grounded knowledge (sub-graph), it employs an information gain-based question selection strategy. This strategy evaluates candidate questions by their expected reduction in diagnostic uncertainty, enabling MedKGI to prioritize the most discriminative inquiries and optimize diagnostic efficiency. Finally, MedKGI adopts the Objective Structured Clinical Examination (OSCE) format (Cushing et al., 2014) to maintain a structured diagnostic state. This state tracks and updates accumulated evidence across dialogue turns, which mitigates context overloading and ensures reasoning consistency.	133
086		134
087		135
088		136
089		137
090		138
091		139
092		140
093		141
094		142
095		143
096		144
097		145
098		146
099		147
100		148
101		149
102		150
103		151
104		152
105	Our key contributions are:	153
106	• A Systematic, Hypothesis-Driven Diagnostic Framework: We propose MedKGI, a novel framework that explicitly models the iterative, hypothesis-driven process of differential diagnosis, bridging the gap between LLMs’ generative nature and the analytical rigor of differential diagnosis.	154
107		155
108		156
109		157
110		158
111		159
112		160
		161
		162
	• Knowledge-Anchored & Strategically Optimized Reasoning: MedKGI uniquely integrates a medical knowledge graph to prevent hallucinations and employs an information gain-based questioning strategy to maximize diagnostic efficiency, grounding reasoning in verified ontologies.	113
		114
		115
		116
		117
		118
		119
	• Superior Empirical Performance: Extensive experiments show that MedKGI outperforms state-of-the-art baselines, achieving high diagnostic accuracy while improving dialogue efficiency by 30%.	120
		121
		122
		123
		124
	2 Related Work	125
	Clinical dialogue involves dynamic, multi-turn information exchange and hypothesis refinement (Nori et al., 2025). We categorize existing approaches into LLM-driven sequential diagnosis, knowledge-augmented frameworks, and agent-based clinical frameworks.	126
		127
		128
		129
		130
		131
	LLM-Driven Sequential Diagnosis. Early methods primarily leverage LLMs’ reasoning capabilities, enhanced through fine-tuning or reinforcement learning (RL) to improve medical diagnosis. Chain-of-Thought (CoT) prompting has been widely adopted to elicit diagnostic reasoning (Dai et al., 2025). Domain-specialized models for diagnosis like Huatuo (Wang et al., 2023) and Meditron (Chen et al., 2023) are pre-trained on medical corpora. AgentClinic simulates doctor–patient interactions but relies on static prompting without dynamic evidence tracking (Schmidgall et al., 2024). Recent works have focused on inquiry strategies: MedAgent (Kim et al., 2025b) formulates diagnosis as multi-agent collaboration while PATIENCE (Zhu et al., 2025) incorporates Bayesian active learning for interactive questioning. However, these model-centric approaches often struggle with hallucinations and struggle with precision in open-ended, multi-turn scenarios due to a lack of external grounding (Zuo et al., 2025).	132
		133
		134
		135
		136
		137
		138
		139
		140
		141
		142
		143
		144
		145
		146
		147
		148
		149
		150
		151
		152
	Knowledge-Augmented Approaches. To mitigate the limitations of pure LLM-based reasoning, recent work has integrated external knowledge. RAG-based methods like MRD-RAG (Sun et al., 2025) leverages the tree-structure medical KG for differential diagnosis, while ClinicalRAG (Lu et al., 2024) fuses structured and unstructured medical knowledge. Beyond retrieval, some methods explicitly model diagnostic reasoning over KGs using search or planning. For instance, Unit of Thought	153
		154
		155
		156
		157
		158
		159
		160
		161
		162

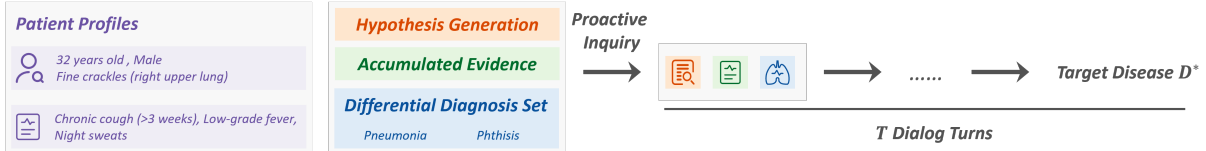


Figure 2: An illustration of the iterative hypothesis refinement process and its corresponding clinical differential diagnosis example.

(UoT) (Hu et al., 2024) decomposes clinical reasoning into discrete, verifiable knowledge units grounded in a KG. However, these approaches treat evidence retrieval statically, lacking dynamic state-tracking for handling diagnostic contexts (Wang et al., 2025).

Agent-based Clinical Frameworks. Multi-agent systems simulate clinical workflows by decomposing tasks across specialized agents for symptom collection, evidence retrieval, and reasoning. DDO (Jia et al., 2025b) uses a diagnosis agent, a strategy agent, and a patient agent for stage-specific inquiry, while MeDDxAgent (Rose et al., 2025) integrates a control agent, a history agent, and a knowledge agent to simulate clinical diagnostic processes with external knowledge. MEDIQ (Li et al., 2024) introduces a query-planning agent that prioritizes questions based on symptom severity, CoD (Chen et al., 2025a) coordinates diagnostic agents through a consensus-driven protocol, and DoctorAgent-RL (Feng et al., 2025) models consultations as an RL process under uncertainty. Despite these advances, existing multi-agent frameworks lack criteria for question selection, relying on heuristic role-playing rather than information-theoretic objectives to optimally reduce diagnostic uncertainty (Chen et al., 2025b).

Summary. Existing approaches provide flexibility, factual accuracy, and workflow simulation, there is no single existing approach that effectively unifies: (1) KG-grounded reasoning, (2) structured state tracking for context management, and (3) information-theoretic inquiry optimization. Our MedKGI framework addresses these challenges by integrating KGs with information gain-driven selection within a structured state tracking mechanism.

3 Problem Definition

The multi-step clinical diagnosis can be modeled as an iterative decision-making process that refines hypotheses over T dialogue turns (Figure 2). The process begins with the patient profile \mathcal{P} , which includes demographics and chief complaints. Over a

sequence of dialogue turns t , the framework maintains a dynamic clinical state. At turn t , given \mathcal{P} and accumulated evidence \mathcal{E}_t , the objective is to estimate the posterior probability for each candidate disease $D \in \mathcal{D}_t$, where $\mathcal{D}_t = \{D_1, D_2, \dots, D_n\}$ through evidence collection. The proactive symptom inquiry is defined as identifying the optimal inquiry s that maximizes the expected reduction in diagnostic uncertainty. This mechanism enables the framework to iteratively generate and refine \mathcal{D}_t until the target disease D^* is reached:

$$P(D | \mathcal{P}, \mathcal{E}_t) \rightarrow \delta(D, D^*)$$

4 Methodology

As illustrated in the differential diagnosis scenario (Figure 2), given patient profiles and symptoms, the objective is to narrow a differential diagnosis set toward the target disease D^* through sequential evidence collection. Unlike static classification, this scenario requires actively navigating a hypothesis space, simulating a doctor’s cognitive process (Polotskaya et al., 2024). Specifically, we formulate the iterative refinement process (Figure 3) as a knowledge-guided active diagnostic framework, where a doctor agent iteratively refines the differential diagnosis set by strategically gathering discriminative evidence.

To realize this iterative and knowledge-driven process, we design the MedKGI workflow integrating three components: (1) **Entity Extraction & Alignment:** maps the diagnosis input and generated hypothesis to a medical KG, constructing a diagnostic subgraph grounded in clinically validated disease-symptom relationships to mitigate hallucinations, (2) **Information Gain-Based Inquiry:** calculates information gain to identify discriminative symptoms, ensuring each clinical inquiry maximally reduces diagnostic uncertainty and improves diagnostic efficiency, (3) **OSCE-Aligned Diagnostic Record Management:** organizes accumulated evidence into an OSCE-format diagnostic record, maintaining a consistent state to prevent

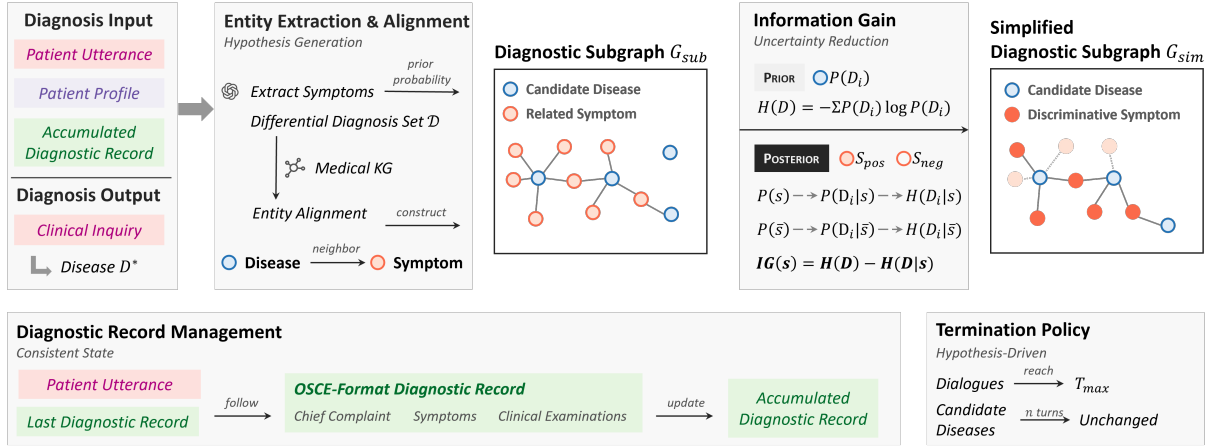


Figure 3: An overview of MedKGI framework. Given a patient’s chief complaint, MedKGI iteratively refines differential diagnosis through (1) medical knowledge graph alignment, (2) information gain–driven symptom inquiry to minimize diagnosis uncertainty, and (3) OSCE-aligned diagnostic records for coherent evidence tracking. A hypothesis-driven termination policy ensures diagnostic efficiency.

context overloading.

4.1 Diagnosis Workflow of MedKGI

At any diagnostic turn t , the doctor agent receives the following as input:

- The *Patient Profile* (e.g., age, sex, and chief complaints) provided at the beginning;
- The *Patient’s Recent Utterance*, which may contain new symptoms or responses to prior questions;
- The *Accumulated Diagnostic Record*, a structured OSCE-aligned summary (detailed in Section 4.4) containing confirmed symptoms, medical history, and examination findings up to turn $t - 1$.

Based on this context, the doctor agent generates two key outputs:

- A natural language *Clinical Inquiry* to elicit discriminative evidence;
- A *Differential Diagnosis Set* outputted by doctor agent;

To ensure efficiency and diagnostic precision, diagnostic process concludes when one of the following termination conditions is met

- **Turn Limit.** If the dialogue reaches T_{max} turns, the doctor agent must issue a final diagnosis.

- **Stagnation Detection.** If the differential diagnosis set remains unchanged for n consecutive turns, doctor agent is prompted to seek evidence that could refute the current hypothesis. If no such symptom exists, the doctor agent outputs the final diagnosis.

4.2 Entity Extraction and Knowledge Graph Alignment

At each turn, the LLM first proposes a preliminary differential diagnosis set, which is then mapped to the medical KG. To align medical terms mentioned in the patient utterance with standardized entities in the KG, we implement a multi-stage alignment pipeline:

- **Exact Matching.** For standard medical terms in the dialogue, we query the KG for an entity that exactly matches the candidate disease name.
- **Edit-Distance Matching.** For minor spelling variations or errors, we apply the Levenshtein Distance (Levenshtein, 1965) to identify approximate matches, allowing a maximum edit distance of 3.
- **Semantic Embedding Matching.** For conceptually equivalent but lexically divergent expressions, we leverage the pre-trained PubMedBERT (Gu et al., 2021) to generate vector embeddings for the candidate disease name and all KG disease entity names. We calculate cosine similarity and retrieve the top-ranked entity, discarding matches with a similarity

score below a threshold $\tau = 0.85$ to ensure alignment quality.

4.3 Information Gain-Based Symptom Selection

Once the candidate diseases are mapped to the KG, we construct a task-specific diagnostic subgraph $G_{sub} = (V_{sub}, E_{sub})$, comprising the current differential diagnosis set and their directly connected symptom nodes. To strategically reduce diagnostic uncertainty, MedKGI selects symptom queries that maximize information gain (Quinlan, 1986). This selection is made over the diagnostic subgraph G_{sub} and is conditioned on the patient’s reported positive and negative symptoms (S_{pos}, S_{neg}).

4.3.1 Posterior Disease Probability Estimation

First, we establish the context by constructing a diagnostic subgraph $G_{sub} = (V_{sub}, E_{sub})$, where $V_{sub} = \bigcup_{i=1}^n \{D_i\} \cup N(D_i)$ consists of the differential diagnosis set $\mathcal{D} = \{D_1, D_2, \dots, D_n\}$ and all symptoms connected to them in the KG: $N(D_i) = \{s \in V_{KG} \vee s \in \mathcal{S} : (D_i, s) \in E_{KG}\}$. We initialize the prior probability of each candidate disease D_i based on its average semantic similarity to the confirmed symptoms S_{pos} extracted from the current dialogue:

$$P(D_i) = \frac{1}{|S_{pos}|} \sum_{s \in S_{pos}} \text{semantic_sim}(s, D_i)$$

where $\text{semantic_sim}(\cdot, \cdot)$ denotes cosine similarity between the PubMedBERT (Gu et al., 2021) embeddings of symptom $s \in S_{pos}$ and disease D_i .

Given the accumulated observed symptoms S_{pos} and S_{neg} , we update disease beliefs over candidate disease \mathcal{D} using Bayes’ theorem:

$$P(D_i | S_{pos}, S_{neg}) = \frac{P(S_{pos}, S_{neg} | D_i) \cdot P(D_i)}{P(S_{pos}, S_{neg})}$$

Assuming conditional independence among symptoms given a disease, the likelihood factorizes as:

$$P(S_{pos}, S_{neg}) = \sum_{j=1}^n P(S_{pos}, S_{neg} | D_j) \cdot P(D_j)$$

where we adopt a uniform conditional probability model: $P(s, | D_i) = 1/|N(D_i)|$ for all symptoms $s \in N(D_i)$.

4.3.2 Information Gain Computation

For the current differential disease set \mathcal{D} with posterior probabilities $P(D_i | S_{pos}, S_{neg})$, we compute the prior diagnostic uncertainty using the Shannon Entropy:

$$H(\mathcal{D}) = - \sum_{i=1}^n P(D_i) \log P(D_i)$$

For any symptom s , we compute its marginal probability and the resulting posterior disease distributions:

$$P(s) = \sum_{i=1}^n P(s | D_i) P(D_i),$$

$$P(D_i | s) = \frac{P(s | D_i) P(D_i)}{P(s)}$$

Finally, the Information Gain of asking about symptom s is defined as the expected reduction in entropy:

$$IG(s) = H(\mathcal{D}) - H(\mathcal{D} | s)$$

where $H(\mathcal{D} | s)$ represents the expected conditional entropy after observing symptom s :

$$H(\mathcal{D} | s) = P(s)H(\mathcal{D} | s^+) + P(\neg s)H(\mathcal{D} | s^-)$$

and $H(\mathcal{D} | s^+)$ and $H(\mathcal{D} | s^-)$ are the entropies if the symptom is observed positive or negative, respectively. The doctor agent then selects the top- k symptoms with the highest $IG(s)$ for strategic inquiry, ensuring that each subsequent question maximally reduces uncertainty. Compared to methods that rely on predefined question templates or fixed inquiry sequences, MedKGI dynamically adapts questions based on the evolving diagnostic hypothesis, enabling more targeted and efficient information gathering.

4.4 Consistent State by Diagnostic Record Management

To support coherent and consistent reasoning, we employ the LLM to generate and maintain a structured diagnostic record in JSON format, aligned with the Objective Structured Clinical Examination (OSCE) standard. At the beginning of each dialogue session, we initialize an empty diagnostic record following a predefined schema, including chief complaint, symptoms, and recent medical examinations.

At each turn, MedKGI takes the latest diagnostic record and patient utterance as input. Then, MedKGI outputs an updated diagnostic record that integrates new evidence while preserving clinical context. This accumulated diagnostic record prevents context overloading across turns, which commonly occurs when context windows accumulate redundant information in vanilla prompting methods (Schmidgall et al., 2024).

5 Evaluation

5.1 Experiment Setup

Datasets. We conducted experiments on two medical QA benchmarks: MedQA (Jin et al., 2021) and CMB (Wang et al., 2024; cme, 2023). To further assess multi-modal clinical reasoning, we additionally introduce a dataset of real-world cases from the NEJM Image Challenge¹. We followed the settings of AgentClinic (Schmidgall et al., 2024) to simulate the multi-agent medical consultation scenarios based on the cases in MedQA and CMB. We denote the processed datasets as agent-MedQA and agent-CMB.

Baselines. We compared MedKGI against 12 baselines across four categories: (1) Dialog-Based Methods: AgentClinic (Schmidgall et al., 2024), CoT (Chain-of-Thought), Huatuo (Wang et al., 2023) and Meditron (Chen et al., 2023); (2) KG-Based Methods: MCTS-BT (Ding et al., 2025), MCTS-MV (Ding et al., 2025), UoT (Hu et al., 2024); (3) Agent-Based Methods: MEDIQ (Li et al., 2024), CoD (Chen et al., 2025a), DDO (Jia et al., 2025b), and MEDDxAgent (Rose et al., 2025); and (4) SFT-Based Methods including models fine-tuned on domain-specific dialogues. Specialized medical LLMs (e.g., Huatuo and Meditron) are not evaluated on the NEJM benchmark if they lack the ability of multi-modal analysis. Furthermore, SFT and SFT-GT are excluded from NEJM evaluation due to the lack of multimodal dialogue training data required for effective fine-tuning. A complete description of all individual models and their configurations is provided in Appendix B.

Implementation. Details of our knowledge graph integration are provided in the Appendix A. To simulate realistic clinical interactions, we implemented a multi-agent framework with three specialized agents, adapted from AgentClinic (Schmidgall et al., 2024):

- Doctor agent asks up to $T_{max} = 20$ questions.
- Patient agent responds only with symptom descriptions and never reveals diagnosis.
- Measurement agent simulates the outcome of laboratory tests or medical examinations.

We modified inquiry termination criteria and evidence-collection protocols to better align with clinical workflows. Detailed descriptions of the specific prompt engineering for each agent are provided in Appendix G.

Our experiment employed Qwen3-8B (Yang et al., 2025) and Llama3.1-8B-Instruct (AI@Meta, 2024) for agent-MedQA and agent-CMB, and Qwen3-VL-8B-Instruct (Yang et al., 2025) for NEJM. For specialized models (e.g., Huatuo and Meditron), we used their architectures.

Metrics. Diagnostic accuracy (*acc*): The diagnostic accuracy is quantified by the exact match between the final diagnostic output and the ground truth D^* . Higher values indicate a more robust alignment with clinical benchmarks. Dialogue rounds (*Rounds*): We also record the average dialogue turns to reach a diagnosis for each method. Fewer *Rounds* indicate more efficient diagnosis.

5.2 Main Result

Table 1 presents a comprehensive comparison of MedKGI against state-of-the-art baselines across three medical consultation benchmarks, agent-MedQA, agent-CMB, and NEJM, using multiple backbone LLMs. Our method demonstrates superior performance in both diagnostic accuracy and efficiency.

Overall Performance. MedKGI achieves superior accuracy across three benchmarks using comparable base models: 69.81% on agent-MedQA (Qwen3-8B), 68.21% on agent-CMB (Qwen3-8B), and 69.09% on NEJM (Qwen3-VL-8B-Instruct). Notably, these results are obtained with the fewest dialogue rounds, 9.11, 9.13, and 10.53 out of a maximum of 20 rounds respectively.

Comparison with LLM-based Methods. Compared to general LLM-based approaches, MedKGI outperforms methods like AgentClinic and CoT across all benchmarks. It also surpasses specialized medical LLMs (marked with *). For instance, on agent-CMB, MedKGI using Qwen3-8B achieves higher accuracy (68.21%) than Medical-CoT with MediTron-7B (66.23%), while doing so in significantly fewer dialogue rounds.

¹<https://www.nejm.org/image-challenge>

Table 1: Comprehensive evaluation of MedKGI across three benchmarks agent-MedQA, agent-CMB, and NEJM using Qwen3-8B, Llama3.1-8B-Instruct, Qwen3-VL-8B-Instruct as base language models. All methods are evaluated with a maximum dialogue round of 20. Reported metrics including average dialogue rounds (*Rounds* ↓) and diagnostic accuracy (*Acc (%)* ↑). Methods marked with * employ specialized LLMs (i.e. HuatuoGPT-o1, Meditron-7B, and DiagnosisGPT-7B) rather than the base LLM used in our unified evaluation. Best results per column are **bolded**; second-best are underlined.

Baselines	agent-MedQA				agent-CMB				NEJM	
	Qwen3-8B		Llama3.1-8B-Instruct		Qwen3-8B		Llama3.1-8B-Instruct		Qwen3-VL-8B-Instruct	
	<i>Rounds</i>	<i>Acc (%)</i>	<i>Rounds</i>	<i>Acc (%)</i>	<i>Rounds</i>	<i>Acc (%)</i>	<i>Rounds</i>	<i>Acc (%)</i>	<i>Rounds</i>	<i>Acc (%)</i>
LLM-Based										
AgentClinic	11.32	59.43	10.37	50.00	10.00	58.28	11.23	59.60	13.28	54.55
CoT	18.92	24.52	17.46	34.90	18.32	43.05	16.33	50.33	16.46	50.91
Huatuo* (HuatuoGPT-o1)	16.70	56.60	-	-	16.56	62.25	-	-	-	-
Medical-CoT* (Meditron-7B)	18.94	60.37	-	-	18.34	<u>66.23</u>	-	-	-	-
KG-Based										
MCTS-BT	14.20	45.28	14.98	39.62	13.78	54.97	14.21	42.38	<u>11.50</u>	56.36
MCTS-MV	14.63	53.77	12.86	<u>52.83</u>	14.77	60.26	12.91	56.95	13.95	67.27
UoT	11.47	54.71	11.01	51.89	10.71	64.28	10.96	58.94	12.32	54.55
Agent-Based										
MediQ	13.80	<u>61.32</u>	13.85	49.06	13.76	65.56	15.01	60.93	14.48	65.45
DDO	17.27	<u>61.32</u>	18.39	50.94	17.38	63.58	18.01	60.93	17.91	70.73
MEDDxAgent	16.44	60.38	16.02	49.06	16.09	61.59	16.48	57.62	16.36	65.45
CoD* (DiagnosisGPT-7B)	13.32	56.60	-	-	11.99	28.50	-	-	-	-
SFT-Based										
SFT	11.51	51.89	11.21	49.06	12.04	59.60	<u>9.95</u>	52.98	-	-
SFT-GT	<u>9.35</u>	50.94	<u>10.27</u>	40.57	<u>9.93</u>	55.63	10.26	51.66	-	-
Ours	9.11	69.81	10.20	53.77	9.13	68.21	9.72	<u>60.26</u>	10.53	<u>69.09</u>

Comparison with KG-based and Agent-based Methods. Among KG-based methods, MedKGI surpasses even competitive approaches like MCTS-MV and UoT. For instance, on agent-CMB with Qwen3-8B, it achieves 68.21% accuracy, exceeding UoT’s 64.28%. In contrast to agent-based approaches such as MediQ, DDO, and MEDDxAgent, MedKGI also demonstrates superior performance. Furthermore, compared to the state-of-the-art method, MedKGI achieves comparable or better accuracy while typically requiring far fewer rounds across all three datasets.

Comparison with SFT-based Methods. While SFT-based methods achieve competitive dialogue efficiency, their accuracy lags behind that of our method. For example, on agent-MedQA using Qwen3-8B, SFT-GT achieves comparable effi-

ciency (9.35 average rounds vs. our 9.11) but its accuracy (50.94%) is significantly lower than ours (69.81%).

5.3 Analysis of Superior Performance

The performance of MedKGI generalizes across different backbone LLMs. For example, with Llama3.1-8B-Instruct, our method achieves the highest accuracy on agent-MedQA (53.77%) and the second-highest on agent-CMB (60.26%), while consistently requiring the fewest dialogue rounds. The superiority of MedKGI across diverse benchmarks and LLMs stems from three key factors: knowledge grounding, context-aware reasoning, and efficient inquiry. First, unlike methods that rely solely on pre-trained LLM knowledge (e.g., Agent-Clinic, CoT), which may lack structured clinical

Table 2: Ablation experiments on agent-MedQA using Qwen3-8B, demonstrating the contribution of each component to diagnostic performance.

Method	Rounds	Acc (%)
w/o Knowledge Graph	13.44	44.34
w/o Clinical Record	12.09	57.55
Random node selection	19.25	31.13
Degree-based node selection	17.47	47.17
Ours	9.11	69.81

reasoning, MedKGI integrates a medical knowledge graph (KG). This external grounding enables precise inference and provides a structured hypothesis space for active querying (Jia et al., 2025a). Our ablation study (Table 2) confirms that removing the KG leads to a significant drop in accuracy (-25.47% on agent-MedQA). Second, compared to other KG-based methods that often rely on heuristic metrics for symptom selection, MedKGI selects questions based on information gain that accounts for patient-specific context. This approach avoids both random noise and popularity bias, leading to more discriminative queries (Kim et al., 2025a). Third, in contrast to agent-based methods (e.g., DDO, MED-DxAgent), our framework minimizes redundant interactions by dynamically pruning the candidate symptom set based on information gain and maintaining diagnostic records. This enables MedKGI to achieve diagnosis in fewer rounds while maintaining high symptom coverage.

5.4 Ablation Experiments

We tested three variants for the ablation study: (1) removing KG integration; (2) disabling the Clinical Record module; and (3) replacing the information gain-based symptom pruning strategy with random or frequency-based alternatives. As shown in Table 2, the full framework consistently achieves the highest diagnostic accuracy. Removing KG integration leads to a significant performance drop, underscoring the critical role of structured external knowledge. Meanwhile, omitting dialogue history summarization results in incomplete patient records, which impairs contextual coherence over multi-turn interactions. Finally, both random and frequency-based symptom pruning strategies result in lower accuracy than our information-gain approach, confirming that targeted, discriminative symptom selection is essential. Together, these

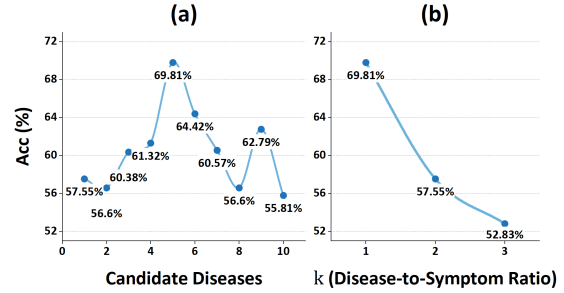


Figure 4: (a) Impact of candidate disease settings on accuracy. (b) Effect of k on Accuracy where k is defined as the ratio of the number of candidate diseases to the number of related symptoms from the KG.

findings validate the necessity of each key component in our design.

5.5 Hyper parameter Selection

We performed controlled experiments to determine the optimal settings for two key hyperparameters. First, we examined how the number of candidate diseases affects accuracy. Figure 4(a) shows that accuracy peaks at 5 candidates and declines with more, as low-relevance candidates introduce noise. In practice, we recommend finding the optimal value by testing on a small sampled dataset. Second, we tested symptom sampling by defining k as the average symptoms per candidate disease. Figure 4(b) indicates optimal performance at $k = 1$, implying that focused symptom selection maximizes discrimination while avoiding redundancy.

6 Conclusion

In this work, we present MedKGI, a framework that formalize multi-step clinical diagnosis as an active, knowledge-guided, and iterative refinement process. By integrating a medical KG for hypothesis grounding, an information gain-driven inquiry strategy for diagnostic uncertainty reduction, and a structured diagnostic record aligned with clinical standards, MedKGI enables systematic and efficient differential diagnosis in multi-turn dialogues. Unlike existing approaches that rely on static retrieval, heuristic questioning, or ungrounded LLM reasoning, our framework explicitly models the evolving clinical state and optimizes each diagnostic step toward maximal discriminative power. Experimental results demonstrate that MedKGI achieves both superior diagnostic accuracy and dialogue efficiency.

548 Limitations

549 While MedKGI demonstrates promising perfor-
550 mance in differential diagnosis, there are several
551 limitations that warrant discussion. First, our pa-
552 tient simulation relies on LLM-generated case de-
553 scriptions that may not fully capture the ambiguity
554 of real patient narratives. Critically, our patient
555 agent assumes cooperative and coherent symptom
556 reporting, reflecting an idealized clinical interac-
557 tion. In reality, patients often exhibit cognitive or
558 linguistic biases: underreporting stigmatized symp-
559 toms, inaccurate recall, or anxiety-driven concerns
560 rather than physiological reasoning. In addition,
561 our information gain computations assume condi-
562 tional independence among symptoms given a
563 disease and employ uniform likelihood over symp-
564 tom–disease edges in KG. This assumption may
565 lead to suboptimal question selection when dis-
566 eases are distinguished primarily by complex symp-
567 tom patterns.

568 Ethical Consideration

569 The application of AI in diagnosis support has
570 raised ethical concerns that we carefully acknowl-
571 edge. MedKGI is proposed as a diagnostic reason-
572 ing assistant rather than a replacement for licensed
573 doctors. All outputs must be validated by human
574 doctors before any clinical action is taken. Our
575 evaluation data are derived from publicly available
576 and anonymized datasets agent-MedQA, agent-
577 CMB and NEJM Image Challenge. No real patient
578 records are used, ensuring compliance with privacy
579 standards.

580 References

581 2023. Cmb: Chinese medical benchmark. [https://](https://github.com/FreedomIntelligence/CMB)
582 github.com/FreedomIntelligence/CMB.

583 AI@Meta. 2024. Llama 3 model card.

584 Payal Chandak, Kexin Huang, and Marinka Zitnik.
585 2023. Building a knowledge graph to enable pre-
586 cision medicine. *Scientific Data*, 10(1):67.

587 Junying Chen, Chi Gui, Anningzhe Gao, Ke Ji, Xidong
588 Wang, Xiang Wan, and Benyou Wang. 2025a. Cod,
589 towards an interpretable medical agent using chain of
590 diagnosis. In *Findings of the Association for Computa-*
591 *tional Linguistics: ACL 2025*, pages 14345–14368.

592 Xi Chen, Huahui Yi, Mingke You, Weizhi Liu, Li Wang,
593 Hairui Li, Xue Zhang, Yingman Guo, Lei Fan, Gang
594 Chen, and 1 others. 2025b. Enhancing diagnostic
595 capability with multi-agents conversational large lan-
596 guage models. *npj Digital Medicine*, 8(1):159.

Zeming Chen, Alejandro Hern A Ndez Cano, An- 597
gelika Romanou, Antoine Bonnet, Kyle Matoba, 598
Francesco Salvi, Matteo Pagliardini, Simin Fan, An- 599
dreas Köpf, Amirkeivan Mohtashami, and 1 others. 600
2023. Meditron-70b: Scaling medical pretraining 601
for large language models. *arXiv e-prints*, page 602
arXiv:2311.16079. 603

A. M. Cushing, J. S. Ker, P. Kinnersley, P. Mckeown, 604
J. Silverman, J. Patterson, and O. M. R. Westwood. 605
2014. Objective structured clinical examination. 606
APA PsycTests Database Record. 607

Guangxin Dai, Xiang Li, Lizhou Fan, and Xin Ma. 2025. 608
Enhancing medical diagnostic reasoning with chain- 609
of-thought in large language models. In *2025 Inter-* 610
national Conference on Mechatronics, Robotics, and 611
Artificial Intelligence (MRAI), pages 294–299. 612

Hongxin Ding, Baixiang Huang, Yue Fang, Weibin Liao, 613
Xinke Jiang, Zheng Li, Junfeng Zhao, and Yasha 614
Wang. 2025. ProMed: Shapley Information Gain 615
Guided Reinforcement Learning for Proactive Medi- 616
cal LLMs. *arXiv e-prints*, page arXiv:2508.13514. 617

Yichun Feng, Jiawei Wang, Lu Zhou, Zhen Lei, 618
and Yixue Li. 2025. Doctoragent-rl: A multi- 619
agent collaborative reinforcement learning system 620
for multi-turn clinical dialogue. *arXiv e-prints*, page 621
arXiv:2505.19630. 622

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto 623
Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng 624
Gao, and Hoifung Poon. 2021. Domain-specific lan- 625
guage model pretraining for biomedical natural lan- 626
guage processing. *ACM Trans. Comput. Healthcare*, 627
3(1). 628

Paul Hager, Friederike Jungmann, Robbie Holland, Ku- 629
nal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob 630
Vielhauer, Marcus Makowski, Rickmer Braren, Geor- 631
gios Kaissis, and Daniel Rueckert. 2024. Evalua- 632
tion and mitigation of the limitations of large lan- 633
guage models in clinical decision-making. *Nature* 634
Medicine, 30(9):2613–2622. 635

Sarah Hannigen. 2018. Differential diagnosis. In Jef- 636
frey S. Kreutzer, John DeLuca, and Bruce Caplan, 637
editors, *Encyclopedia of Clinical Neuropsychology*, 638
Encyclopedia of Clinical Neuropsychology, pages 639
1148–1148. Springer International Publishing, Cham. 640

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan 641
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and 642
Weizhu Chen. 2021. Lora: Low-rank adaptation 643
of large language models. *arXiv e-prints*, page 644
arXiv:2106.09685. 645

Zhiyuan Hu, Chumin Liu, Xidong Feng, Yilun Zhao, 646
See-Kiong Ng, Anh Tuan Luu, Junxian He, Pang Wei 647
Koh, and Bryan Hooi. 2024. Uncertainty of thoughts: 648
Uncertainty-aware planning enhances information 649
seeking in llms. In *Advances in Neural Information* 650
Processing Systems, volume 37, pages 24181–24215. 651

652	Mingyi Jia, Junwen Duan, Yan Song, and Jianxin Wang.	heterogeneous knowledge retrieval. Proceedings of	708
653	2025a. medikal: Integrating knowledge graphs as	the 1st Workshop on Towards Knowledgeable Lan-	709
654	assistants of llms for enhanced clinical diagnosis on	guage Models (KnowLLM 2024), pages 64–68.	710
655	emrs . In <i>Proceedings of the 31st International Con-</i>		
656	<i>ference on Computational Linguistics</i> , pages 9278–		
657	9298, Abu Dhabi, UAE. Association for Computa-	Harsha Nori, Mayank Daswani, Christopher Kelly, Scott	711
658	tional Linguistics.	Lundberg, Marco Tulio Ribeiro, Marc Wilson, Xi-	712
		aoxuan Liu, Viknesh Sounderajah, Jonathan Carlson,	713
659	Zhihao Jia, Mingyi Jia, Junwen Duan, and Jianxin Wang.	Matthew P. Lungren, and 1 others. 2025. Sequential	714
660	2025b. Ddo: Dual-decision optimization for llm-	diagnosis with language models . <i>arXiv e-prints</i> , page	715
661	based medical consultation via multi-agent collabora-	arXiv:2506.22405.	716
662	tion . In <i>Proceedings of the 2025 Conference on</i>		
663	<i>Empirical Methods in Natural Language Processing</i> ,	Kristina Polotskaya, Carlos S. Muñoz-Valencia, Alejandro	717
664	pages 26380–26397.	Rabasa, Jose A. Quesada-Rico, Domingo Orozco-	718
		Beltrán, and Xavier Barber. 2024. Bayesian networks	719
665	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,	for the diagnosis and prognosis of diseases: A scop-	720
666	Hanyi Fang, and Peter Szolovits. 2021. What disease	ing review . <i>Machine Learning and Knowledge Ex-</i>	721
667	does this patient have? a large-scale open domain	<i>traction</i> , 6(2):1243–1262.	722
668	question answering dataset from medical exams . <i>Ap-</i>		
669	<i>plied Sciences</i> , 11(14).	J. R. Quinlan. 1986. Induction of decision trees . <i>Ma-</i>	723
		<i>chine Learning</i> , 1(1):81–106.	724
670	Jonathan Kim, Anna Podlasek, Kie Shidara, Feng Liu,	Daniel Philip Rose, Chia-Chien Hung, Marco Lepri,	725
671	Ahmed Alaa, and Danilo Bernardo. 2025a. Limita-	Israa Alqassem, Kiril Gashkevski, and Carolin	726
672	tions of large language models in clinical problem-	Lawrence. 2025. Meddxagent: A unified modular	727
673	solving arising from inflexible reasoning . <i>Scientific</i>	agent framework for explainable automatic differ-	728
674	<i>Reports</i> , 15(1):39426.	ential diagnosis . In <i>Proceedings of the 63rd An-</i>	729
		<i>Annual Meeting of the Association for Computational</i>	730
675	Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu	<i>Linguistics (Volume 1: Long Papers)</i> , pages 13803–	731
676	Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee,	13826.	732
677	Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won	Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath	733
678	Park. 2025b. Mdagents: An adaptive collaboration	Rangan, and Jonathan H. Chen. 2024. Diagnostic	734
679	of llms for medical decision-making . In <i>Proceeed-</i>	reasoning prompts reveal the potential for large lan-	735
680	<i>ings of the 38th International Conference on Neural</i>	guage model interpretability in medicine . <i>npj Digital</i>	736
681	<i>Information Processing Systems</i> .	<i>Medicine</i> , 7(1):20.	737
682	Vladimir I. Levenshtein. 1965. Binary codes capable of	Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo	738
683	correcting deletions, insertions, and reversals. <i>Soviet</i>	Reis, Jeffrey Jopling, and Michael Moor. 2024.	739
684	<i>physics. Doklady</i> , 10:707–710.	Agentclinic: A Multimodal Agent Benchmark to	740
		Evaluate Ai in Simulated Clinical Environments .	741
685	Shuyue Stella Li, Vidhisha Balachandran, Shangbin	<i>arXiv e-prints</i> , page arXiv:2405.07960.	742
686	Feng, Jonathan S. Ilgen, Emma Pierson, Pang Wei		
687	Koh, and Yulia Tsvetkov. 2024. Mediq: Question-	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres,	743
688	asking llms and a benchmark for reliable interactive	Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin	744
689	clinical reasoning . In <i>Advances in Neural Informa-</i>	Clark, Stephen R Pfohl, Heather Cole-Lewis, and	745
690	<i>tion Processing Systems</i> , volume 37, pages 28858–	1 others. 2025. Toward expert-level medical ques-	746
691	28888.	tion answering with large language models . <i>Nature</i>	747
692	Shuyue Stella Li, Jimin Mun, Faeze Brahman, Pedram	<i>Medicine</i> , 31(3):943–950.	748
693	Hosseini, Bryceton G. Thomas, Jessica M. Sin, Bing	Penglei Sun, Yixiang Chen, Xiang Li, and Xiaowen Chu.	749
694	Ren, Jonathan S. Ilgen, Yulia Tsvetkov, and Maarten	2025. The multi-round diagnostic rag framework for	750
695	Sap. 2025. Alfa: Aligning Llms to ask good ques-	emulating clinical reasoning . <i>arXiv e-prints</i> , page	751
696	tions a case study in clinical reasoning . In <i>Second</i>	arXiv:2504.07724.	752
697	<i>Conference on Language Modeling</i> .		
698	Yanna Lin, Shaojie Xu, Wenshuo Zhang, Yushi Sun,	Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang,	753
699	Zixin Chen, Yanjie Zhang, and Rui Sheng. 2025. A	Sendong Zhao, Bing Qin, and Ting Liu. 2023. Hu-	754
700	survey of llm-based multi-agent systems in medicine .	atuo: Tuning llama model with chinese medical	755
		knowledge . <i>arXiv e-prints</i> , page arXiv:2304.06975.	756
701	Fenglin Liu, Hongjian Zhou, Boyang Gu, Xinyu Zou,	Xi Wang, Procheta Sen, Ruizhe Li, and Emine Yilmaz.	757
702	Jinfa Huang, Jinge Wu, Yiru Li, Sam S. Chen, Yining	2025. Adaptive retrieval-augmented generation for	758
703	Hua, Peilin Zhou, and 1 others. 2025. Application of	conversational systems . In <i>Findings of the Associ-</i>	759
704	large language models in medicine . <i>Nature Reviews</i>	<i>ation for Computational Linguistics: NAACL 2025</i> ,	760
705	Bioengineering , 3(6):445–464.	pages 491–503, Albuquerque, New Mexico. Associa-	761
706	Yuxing Lu, Xukai Zhao, and Jinzhuo Wang. 2024. Clin-	<i>ation for Computational Linguistics</i> .	762
707	icalrag: Enhancing clinical decision support through		

763 Xidong Wang, Guiming Chen, Song Dingjie, Zhang
764 Zhiyi, Zhihong Chen, Qingying Xiao, Junying Chen,
765 Feng Jiang, Jianquan Li, Xiang Wan, and 1 others.
766 2024. [Cmb: A comprehensive medical benchmark
767 in Chinese](#). In *Proceedings of the 2024 Conference
768 of the North American Chapter of the Association
769 for Computational Linguistics: Human Language
770 Technologies (Volume 1: Long Papers)*, pages 6184–
771 6205. Association for Computational Linguistics.

772 Kaishuai Xu, Yi Cheng, Wenjun Hou, Qiaoyu Tan, and
773 Wenjie Li. 2024. [Reasoning like a doctor: Improving
774 medical dialogue systems via diagnostic reasoning
775 process alignment](#). In *Findings of the Association for
776 Computational Linguistics: ACL 2024*, pages 6796–
777 6814.

778 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
779 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
780 Chengen Huang, Chenxu Lv, and 1 others. 2025.
781 [Qwen3 Technical Report](#). *arXiv e-prints*, page
782 arXiv:2505.09388.

783 Jiayuan Zhu, Jiazhen Pan, Yuyuan Liu, Fenglin Liu,
784 and Junde Wu. 2025. [Ask Patients with Patience:
785 Enabling LLMs for Human-Centric Medical Dia-
786 logue with Grounded Reasoning](#). *arXiv e-prints*,
787 page arXiv:2502.07143.

788 Zhihong Zhu, Yunyan Zhang, Xianwei Zhuang, Fan
789 Zhang, Zhongwei Wan, Yuyan Chen, Qingqing Long,
790 Yefeng Zheng, and Xian Wu. 2025. [Can we trust ai
791 doctors? a survey of medical hallucination in large
792 language and large vision-language models](#). In *Find-
793 ings of the Association for Computational Linguistics:
794 ACL 2025*, pages 6748–6769.

795 Kaiwen Zuo, Yirui Jiang, Fan Mo, and Pietro Lio. 2025.
796 [Kg4diagnosis: A hierarchical multi-agent llm frame-
797 work with knowledge graph enhancement for medical
798 diagnosis](#). In *Proceedings of The First AAAI Bridge
799 Program on AI for Medicine and Healthcare*, volume
800 281 of *Proceedings of Machine Learning Research*,
801 pages 195–204. PMLR.

A KG Implementation

To provide a clinically grounded foundation for our framework, we utilize PrimeKG (Precision Medicine Knowledge Graph) (Chandak et al., 2023), which is a comprehensive resource that integrates over 20 high-quality primary sources, including Orphanet, Mayo Clinic, and DrugBank.

PrimeKG comprises:

- **Nodes:** Approximately 17,000 disease entities and 1,300 symptom entities.
- **Edges:** We prioritize two primary relationship types:
 - **Disease-Symptom:** Indicating clinical manifestations associated with specific pathologies.
 - **Disease-Disease:** Representing comorbidity links and hierarchical relationships (e.g., “is-a” or “part-of” relations) that assist in differential grouping.

B Detailed Description of Baselines

We provide a comprehensive overview of the 12 baseline methods used in our experiments:

- **Dialog-Based Methods.** **AgentClinic** (Schmidgall et al., 2024): Simulates clinician–patient dialogues for diagnosis. **CoT** (Chain-of-Thought): Appends “Let’s think step by step” to encourage explicit reasoning. **Huatuo** (Wang et al., 2023) and (Chen et al., 2023): Representative specialized medical LLMs.
- **KG-Based Methods.** **MCTS-BT:** Uses Monte Carlo Tree Search with backtracking for hypothesis refinement. **MCTS-MV:** Extends MCTS by ranking symptom queries based on contextual informativeness. **UoT** (Unit of Thought) (Hu et al., 2024): Constructs symptom-centric “units” around confirmed positive symptoms and prioritizes structural importance.
- **KG-Based Methods.** **MEDIQ** (Li et al., 2024): A diagnostic agent implementing multiple diagnostic strategies through sequential dialogues. **CoD** (Chen et al., 2025a): Selects questions by maximizing information entropy over candidate diseases. **DDO** (Jia et al., 2025b): a multi-agent framework that

dynamically chooses symptoms using diverse strategies. **MEDDxAgent** (Rose et al., 2025): Adapts questioning strategy based on diagnostic uncertainty.

- **SFT-Based Methods.** **SFT / SFT-GT:** Fine-tuning on Qwen3-8B using generated dialogues by **AgentClinic** (Schmidgall et al., 2024), with or without ground-truth labels respectively.

C Dataset and Case Generation

C.1 Prompt for OSCE Case Generation

We use the following prompt to generate standardized Objective Structured Clinical Examination (OSCE) scenarios for evaluation. The prompt instructs the LLM to produce a structured JSON containing patient demographics, symptom history, physical findings, test results, and the ground-truth diagnosis while providing only the clinical objective to the Doctor Agent.

Figure 5: Prompt Template for OSCE Case Generation.

```
Please generate a sample Objective Structured Clinical Examination (OSCE) for the patient actor and the doctor, including what the correct diagnosis should be as a structured JSON.

Only provide the doctor with the objective and provide "test results" as a separate category. Provide these for a primary care doctor exam.

Generate an OSCE for the following case study. Please read the answer category for the correct diagnosis. Here is an example of correcting the OSCE format {example}. Please create a new one here:
```

An example output is shown in Figure 6.

D Implementation Details

D.1 Base Model and Inference Configuration

- **Base models:** Qwen3-8B, Meta-Llama-3.1-8B-Instruct, and
- **Inference temperature:** 0.05
- **Max tokens:** 2048

D.2 Knowledge Graph Statistics

We use PrimeKG (Chandak et al., 2023), which contains:

- 877 - 17,080 disease nodes
- 878 - 3,357 symptom nodes
- 879 - 1,361 disease–disease relationships
- 880 - 11,072 disease–symptom relationships

881 D.3 LoRA Fine-Tuning Hyperparameters

882 We fine-tune the base LLM using Low-Rank Adap-
883 tation (LoRA) (Hu et al., 2021) with the following
884 configuration:

- 885 - **Learning rate:** $2e-4$
- 886 - **Batch size:**
 - 887 - `per_device_train_batch_size` = 2
 - 888 - `gradient_accumulation_steps` = 4
 - 889 - Effective batch size = $2 \times 4 = 8$
- 890 - **LoRA rank (r):** 8
- 891 - **Target modules:** ["q_proj", "k_proj",
892 "v_proj", "o_proj"]
- 893 - **LoRA alpha:** 32
- 894 - **Dropout:** 0.1
- 895 - **Training epochs:** 3
- 896 - **Warmup steps:** 1,000

897 E Prompt for Entity Extraction and 898 Evaluation

899 **Symptom Entity Extraction.** For symptom ex-
900 traction from patient utterances, we employ the
901 following prompt:

902 **Diagnostic Accuracy Judgment.** To evaluate
903 whether the Doctor Agent’s final diagnosis matches
904 the ground truth, we use the judgment prompt:

Figure 7: Prompt Template for Diagnosis Accuracy Check.

You are responsible for determining if the correct diagnosis and the doctor's diagnosis are the same disease. Please respond only with Yes or No. Nothing else. Here is the correct diagnosis: `{ground truth diagnosis}` Here was the doctor diagnosis: `{doctor diagnosis}` Are these the same?

905 F Prompt for Diagnostic Record 906 Initialization and Update

907 To ensure consistency of Diagnostic Record
908 throughout the diagnosis process, we use a prompt
909 that guides the LLM to perform evidence-based
910 updates:

Figure 6: Prompt Template for Symptom Entity Extraction.

You are a helpful assistant with expertise in medical symptom identification. Please identify and extract all disease and symptom entities from the following sentence.

Each entity must be no longer than 5 characters.

Rules:

1. Include symptoms that are affirmed (positive) in the "positive" list
2. Include symptoms that are explicitly denied (negative) in the "negative" list
3. Pay special attention to negation words like "no", "not", "don't", "haven't", "can still", which typically indicate negative symptoms

Return the result in valid JSON format as shown below (without any markdown formatting and explanation):

```
{
  "positive": ["Symptom 1", "Symptom 2", ...],
  "negative": ["Symptom 3", "Symptom 4", ...]
}
```

Sentence: `{sentence}`.

Figure 8: Prompt Template for Diagnostic Record Initialization and Update.

You are an experienced medical scribe. Your task is to read the patient's latest utterance and incrementally update the structured summary below.

Rules:

1. Add or revise only facts confirmed in the CURRENT utterance.
2. Preserve all existing information that is not contradicted.
3. Use the exact JSON schema that was provided (do not create new keys).
4. Return only the updated JSON object, with no extra commentary (without any markdown formatting).

Schema: `{schema}`
Current structured summary: `{current diagnostic record}`
Latest patient utterance: `{sentence}`

G Prompt for Agent Initialization and Diagnosis

To simulate realistic clinical interactions, we implement a multi-agent diagnostic framework comprising three specialized agents: doctor agent, Patient Agent, and Measurement Agent with each guided by prompts to enforce specific behaviors and constraints.

The doctor agent adopts a constrained prompt specifying question limits T_{max} and tracks the count of asked questions $t_{current}$.

Figure 9: Prompt Template for doctor agent Initialization.

You are a doctor named Dr. Agent who only responds in the form of dialogue. You are inspecting a patient whom you will ask questions in order to understand their disease.

You are allowed to ask `{T_max}` questions total before you must make a decision, and have asked `{t_current}` questions so far.

Additionally, during the doctor agent's differential diagnosis, we employ the following prompt template, which integrates patient demographics, recent dialogue history, structured clinical findings, and relevant medical knowledge extracted from the knowledge graph.

Figure 10: Prompt Template for doctor agent Differential Diagnosis.

Age: `{age}`
 Gender: `{gender}`
 Chief Complaint: `{chief complaint}`

Recent dialogue history:
`{3 recent dialogue turns}`

Medical record
 Chief Complaint: `{chief complaint}`
 Symptoms: `{symptoms}`
 Recent medical examinations:
`{recent medical examinations}`

Knowledge Graph Context
 Relevant medical knowledge represented as triples: `{simplified subgraph triples}`

The Patient Agent prevents patients from revealing diagnostic results directly, forcing the doctor agent to make diagnoses through symptom reasoning.

Figure 11: Prompt Template for Patient Agent Initialization.

You're a patient in a clinic. The doctor will ask questions or order exams to figure out your illness.

Below is all of your information:
`{patient profile}`

Never name your disease and only describe symptoms naturally: how they feel, when they flare, or how they affect you.

The Measurement Agent adopts a standardized result output format (RESULTS: [results here]), ensuring parseability of medical examination results.

Figure 12: Prompt Template for Measurement Agent Initialization.

You are a measurement reader who responds with medical test results. Please respond in the format "RESULTS: [results here]".

Below is all of the information you have:
`{medical examinations}`. If the requested results are not in your data, then you can respond with NORMAL READINGS.

Figure 13: An example of OSCE Case.

```

{
  "OSCE Examination": {
    "Objective for Doctor": "Assess and diagnose the patient presenting with acute abdominal pain .",
    "Patient Actor": {
      "Demographics": "30-year-old female",
      "History": "The patient complains of sudden onset of sharp, right lower quadrant abdominal pain since last night. The pain has progressively worsened over the last 12 hours. She mentions that she felt nauseous this morning but has not vomited. No recent changes in bowel habits or urinary symptoms have been reported.",
      "Symptoms": {
        "Primary Symptom": "Sharp, right lower quadrant abdominal pain",
        "Secondary Symptoms": ["Nausea", "No vomiting", "No change in bowel habits", "No urinary symptoms"]
      },
      "Past Medical History": "No significant past medical history. No previous surgeries.",
      "Social History": "Non-smoker, occasional alcohol use. Works as a software developer.",
      "Review of Systems": "Denies fever, vomiting, diarrhea, dysuria, or flank pain."
    },
    "Physical Examination Findings": {
      "Vital Signs": {
        "Temperature": "37.2 °C (99 °F)",
        "Blood Pressure": "120/75 mmHg",
        "Heart Rate": "78 bpm",
        "Respiratory Rate": "16 breaths/min"
      },
      "Abdominal Examination": {
        "Inspection": "No distension or visible masses.",
        "Auscultation": "Normal bowel sounds.",
        "Percussion": "Tympanic throughout, no shifting dullness.",
        "Palpation": "Tenderness in the right lower quadrant. No guarding or rebound tenderness. Rovsing's sign positive, suggesting peritoneal irritation."
      }
    },
    "Test Results": {
      "Complete Blood Count": {
        "WBC": "12,000 /μL (elevated)",
        "Hemoglobin": "13.5 g/dL",
        "Platelets": "250,000 /μL"
      },
      "Urinalysis": {
        "Appearance": "Clear",
        "WBC": "2-5 /HPF",
        "RBC": "0-2 /HPF",
        "Nitrites": "Negative",
        "Leukocyte Esterase": "Negative"
      },
      "Imaging": {
        "Ultrasound Abdomen": {
          "Findings": "Enlarged appendix with wall thickening and fluid collection. No evidence of ovarian cyst or ectopic pregnancy."
        }
      }
    },
    "Correct Diagnosis": "Acute Appendicitis"
  }
}

```