

GPT-RagAD: Two-layer Retrieval-Augmented Multilingual Diagnosis System

Xinyi Liu

University of Illinois Urbana-Champaign, 201 N Goodwin Ave, Urbana, IL 61801, USA

LIU323@ILLINOIS.EDU

Dachun Sun

University of Illinois Urbana-Champaign, 201 N Goodwin Ave, Urbana, IL 61801, USA

DSUN18@ILLINOIS.EDU

Yi R. Fung

Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong SAR

YRFUNG@UST.HK

Dilek Hakkani-Tür

University of Illinois Urbana-Champaign, 201 N Goodwin Ave, Urbana, IL 61801, USA

DILEK@ILLINOIS.EDU

Tarek Abdelzaher

University of Illinois Urbana-Champaign, 201 N Goodwin Ave, Urbana, IL 61801, USA

ZAHER@ILLINOIS.EDU

Abstract

We introduce GPT-RagAD, a multilingual, zero-shot automated diagnosis system that achieves high accuracy without relying on real patient data.¹ GPT-RagAD adopts a two-layer Retrieval-Augmented Generation (RAG) architecture: a knowledge graph-based retriever selects disease candidates from 1,058 conditions, and an LLM-based re-ranker applies prompt-based reasoning to refine predictions. Unlike traditional diagnostic models that require supervised training and large clinical datasets, GPT-RagAD is privacy-preserving, scalable, and language-agnostic. Extensive evaluations on three multilingual datasets (Chinese and English) show that GPT-RagAD achieves 40.6% Hit@1 and 56.7% NDCG@10 on the Symptom2Disease benchmark—substantially outperforming embedding-based and direct LLM baselines. Ablation and sensitivity analyses further validate its robustness. GPT-RagAD presents a practical, lightweight solution for clinical triage and pre-diagnosis support.

Keywords: Automated Diagnosis, Large Language Models, Retrieval-Augmented Generation, Knowledge Graph, Zero-shot Learning, Medical NLP, Heterogeneous Information Network

Data and Code Availability We build our disease-symptom knowledge graph from the public *Mayo Clinic Symptoms and Diseases* dataset (1,058

diseases), and evaluate on three existing benchmarks: DX Xu et al. (2019), IMCS21 Chen et al. (2023), and Symptom2Disease (Kaggle). All datasets are publicly available from their original providers; we do not use or release any private clinical records. An anonymized implementation of GPT-RagAD (graph construction, inference, and evaluation scripts) is included as supplementary material for review. If the paper is accepted, we will de-anonymize and publicly release the code.

Institutional Review Board (IRB) Our work uses only previously published, de-identified datasets and does not involve new data collection or access to identifiable health information. Under our institution’s policy, this study is not considered human subjects research and therefore did not require IRB review or approval.

1. Introduction

Automated diagnosis (AD) is a cornerstone of modern healthcare systems, enabling swift and accurate identification of potential diseases based on patient-reported symptoms. Traditional AD methods often rely on supervised learning models trained on vast amounts of patient data and diagnostic reports Hosseini et al. (2018); Wang et al. (2021); Shoham and Rappoport (2023); Tu et al. (2024). However, several challenges remain with this approach, including concerns about data privacy, the difficulty in collecting comprehensive and diverse patient data, the need

1. Code is available at <https://github.com/tracy3057/GPT-RagAD>.

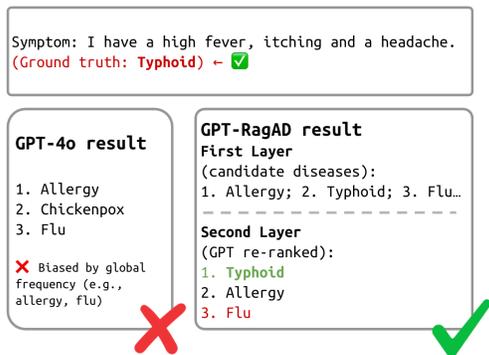


Figure 1: An example of GPT-4o making incorrect diagnosis for a patient, influenced by global priors bias.

61 for high-quality data labeling, and the challenge of
 62 acquiring sufficient patient cases for training. These
 63 limitations impede the effectiveness of conventional
 64 AD systems.

65 Recently, large language models (LLMs) have
 66 demonstrated impressive generalization abilities
 67 across a variety of tasks and have shown considerable
 68 promise in the AD domain. LLMs benefit from ex-
 69 tensive knowledge learned from large-scale real-world
 70 data, enabling them to perform automatic diagnoses
 71 with remarkable accuracy. The use of LLMs not only
 72 mitigates the challenges of data collection and pri-
 73 vacy but also facilitates more flexible diagnosis sys-
 74 tems. However, general-purpose LLMs are not explic-
 75 itly trained for medical diagnosis tasks. As a result,
 76 they may suffer from biases introduced by global pri-
 77 ors, which can lead to suboptimal or incorrect diag-
 78 noses. As illustrated in Figure 1, a general-purpose
 79 LLM such as GPT-4o can be influenced by global
 80 priors, leading to incorrect diagnoses.

81 One potential solution is fine-tuning domain-
 82 specific LLMs for AD tasks [Brown et al. \(2020\)](#); [Fan et al. \(2024\)](#); [Zhao et al. \(2023\)](#). However, this ap-
 83 proach is not cost-effective or feasible due to concerns
 84 over patient data privacy and the challenges related
 85 to the quality of available data. To address these
 86 concerns, Retrieval-Augmented Generation (RAG)
 87 methods have gained traction [Brown et al. \(2020\)](#);
 88 [Fan et al. \(2024\)](#); [Zhao et al. \(2023\)](#). RAG facilitates
 89 the dynamic retrieval of external knowledge sources,
 90 thereby introducing personalized information into the
 91 diagnostic process. This method helps mitigate the
 92 global prior bias. However, directly applying RAG to
 93 AD tasks presents two significant challenges. First,

95 relying on large-scale external medical resources, such
 96 as PubMed, UMLS, and clinical guidelines, incurs
 97 high computational costs, making these sources un-
 98 suitable for real-time diagnosis. This reliance on mas-
 99 sive datasets complicates the efficient use of essential
 100 data for medical diagnosis. Second, using patient
 101 case retrieval for augmentation raises privacy con-
 102 cerns, as it may inadvertently expose sensitive health
 103 data from previous patients.

104 To address these challenges, we propose GPT-
 105 RagAD, a novel and effective two-layer framework for
 106 knowledge-graph-based AD. Unlike traditional meth-
 107 ods that depend on large and computationally ex-
 108 pensive medical datasets, GPT-RagAD generates dis-
 109 ease knowledge graphs from medical encyclopedias,
 110 which focus on describing the core features of dis-
 111 eases. This allows for the efficient use of essen-
 112 tial, high-quality information without the computa-
 113 tional burden of vast medical datasets. By combin-
 114 ing these medical knowledge graphs with the powerful
 115 language understanding capabilities of LLMs, GPT-
 116 RagAD predicts diseases from a more comprehensive
 117 disease space while alleviating biases from global pri-
 118 ors. In this way, GPT-RagAD provides a more ac-
 119 curate, privacy-preserving, and scalable approach to
 120 automated diagnosis.

121 Figure 2 illustrates the overall architecture of GPT-
 122 RagAD, which integrates knowledge graph generation
 123 in the first layer and leverages the strengths of LLMs
 124 in the second layer to refine the diagnostic process.
 125 This two-layer framework not only addresses compu-
 126 tational cost and privacy concerns but also enhances
 127 the accuracy and scalability of AD systems.

128 We introduce GPT-RagAD, a multilingual, two-
 129 layer diagnosis framework that combines symbolic
 130 knowledge retrieval with prompt-based LLM reason-
 131 ing. The first layer performs coarse-grained disease
 132 retrieval using a disease-symptom knowledge graph,
 133 while the second layer refines predictions through a
 134 retrieval-augmented generation (RAG) strategy.

135 To efficiently represent complex and uncertain
 136 symptom-disease associations, we build a Heteroge-
 137 neous Information Network (HIN) and apply a Varia-
 138 tional Graph Autoencoder (VGAE) [Kipf and Welling \(2016\)](#) that produce deterministic representations,
 139 VGAE models uncertainty explicitly, helping reduce
 140 overconfidence in noisy matches. Compared to al-
 141 ternatives like GAT [Veličković et al. \(2017\)](#) and R-
 142 GCN [Schlichtkrull et al. \(2018\)](#), VGAE offers a bet-
 143 ter balance of robustness, scalability, and expressive-
 144 ness—especially when applied to loosely structured
 145

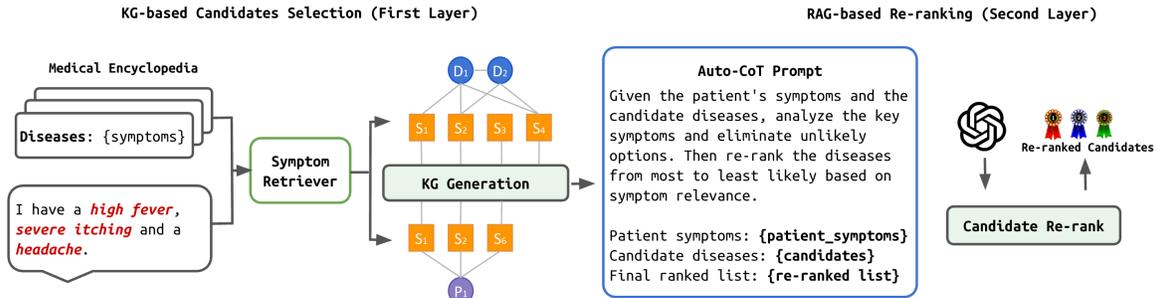


Figure 2: Overview of GPT-RagAD: Knowledge Graph-Based Disease Candidate Selection and RAG-based Re-ranking.

146 encyclopedia-derived graphs rather than curated clinical
147 datasets.

148 Finally, the learned HIN embeddings guide the
149 RAG module to focus on high-quality candidates and
150 consider the relative importance of different symp-
151 toms, enhancing the precision and interpretability of
152 downstream predictions.

153 Our contributions are as follows:

154 Contributions.

- 155 • A RAG-enhanced AD framework that
156 builds a disease-symptom HIN from cu-
157 rated encyclopedic data, encoding dis-
158 ease-symptom/disease-disease/symptom-symptom
159 relations.
- 160 • A two-stage system: fast HIN-based retrieval plus
161 LLM re-ranking, achieving scalability and high ac-
162 curacy with low inference overhead.
- 163 • Multilingual diagnosis over 1,058 diseases with
164 >70% Hit@10 on English and Chinese datasets,
165 outperforming state-of-the-art baselines.
- 166 • Unified symbolic retrieval and LLM reasoning that
167 avoids annotation, preserves privacy, and supports
168 practical deployment.

169 2. Related Work

170 2.1. Graph- and LLM-based Diagnosis

171 Graph-based AD encodes EHRs into heteroge-
172 neous information networks (HINs) to model dis-
173 ease-symptom-entity relations Hosseini et al. (2018);
174 Wang et al. (2021). These methods work well
175 with structured data but rely on large, clean EHRs,
176 limiting cross-domain and multilingual generaliza-
177 tion. LLM-based approaches Shoham and Rappoport
178 (2023); Wang et al. (2023); Tu et al. (2024) pre-

179 dict diseases or synthesize symptoms directly from
180 text; they are flexible and annotation-free but ex-
181 hibit global prior bias and lack explicit reasoning over
182 structured medical knowledge.

183 2.2. RAG for Clinical Reasoning and Our 184 Positioning

185 RAG grounds LLMs with retrieved context Lewis
186 et al. (2020); Ram et al. (2023); Shi et al. (2023),
187 yet most pipelines retrieve free text, risking hal-
188 lucination, higher latency, and privacy concerns
189 in clinical settings. **GPT-RagAD** retrieves from
190 a disease-symptom knowledge graph (public en-
191 cyclopedias), models uncertainty via VGAE, and
192 applies LLM re-ranking with optimized prompt-
193 ing—improving interpretability and scalability while
194 preserving privacy and enabling cross-lingual gener-
195 alization. See Appendix D for an extended survey.

196 3. GPT-RagAD Framework

197 The overview of GPT-RagAD is in Figure 2. We
198 firstly develop the first layer: KG-based Candidates
199 Selection, shown in Figure 3. Then, based on the se-
200 lected candidate diseases, we establish a RAG-based
201 candidate Re-ranking, introduced in Section 3.2. The
202 detailed procedures are shown in Algorithm 1.

203 3.1. First Layer: Graph-based Candidate 204 Disease Generation

205 In the first layer of GPT-RAGAD, we construct an
206 enhanced multi-relational disease-symptom knowl-
207 edge graph and use a graph-based similarity frame-
208 work to generate candidate diseases. The process in-

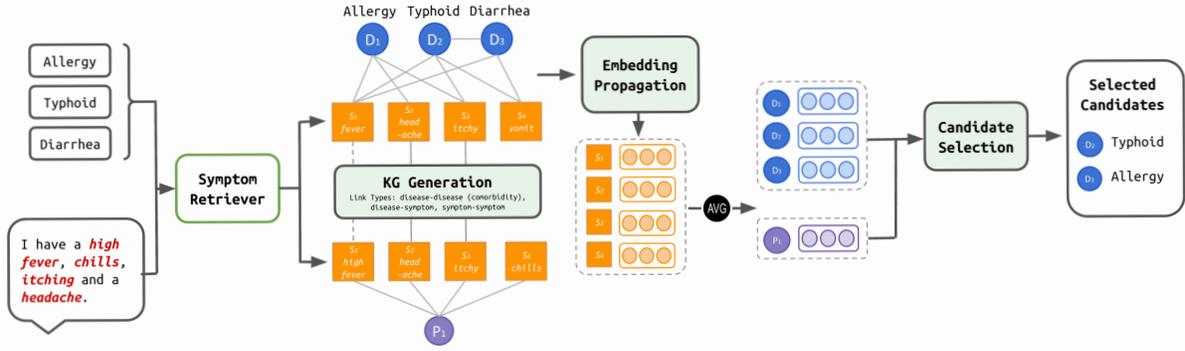


Figure 3: Detailed procedure of KG-based Candidates Selection (First Layer).

Algorithm 1: GPT-RagAD Disease Diagnosis Process

- 1: **Input:**
- 2: \mathcal{D} : Disease descriptions from medical encyclopedia
- 3: \mathcal{P} : Patient symptom inputs (free text)
- 4: **Output:** Final diagnosis $\hat{d} \in \mathcal{D}$ for each patient
- 5: **Step 1: Knowledge Graph Construction**
- 6: Extract symptom sets $S(D_i)$ from each $D_i \in \mathcal{D}$ via LLM
- 7: Build multi-relational graph $G = (V, E)$ with:
 - 8: E^{ds} : disease–symptom edges from encyclopedic links
 - 9: E^{sym} : symptom–symptom similarity edges (BERT similarity $\geq \tau$)
 - 10: E^{dd} : disease–disease comorbidity or Jaccard similarity edges
- 11: **Step 2: Patient Symptom Normalization**
- 12: For each $P_i \in \mathcal{P}$, extract raw symptoms via LLM
- 13: Normalize symptoms to KG vocabulary using cosine similarity $\geq \beta$
- 14: Obtain patient symptom set \hat{S}_i^P
- 15: **Step 3: Graph-based Embedding Learning**
- 16: Apply relation-aware VGAE on G to obtain node embeddings \mathbf{z}_v
- 17: Aggregate:
 - 18: $\mathbf{z}_{P_i} = \text{mean}(\mathbf{z}_{s_j})$ for $s_j \in \hat{S}_i^P$
 - 19: $\mathbf{z}_{D_j} = \text{mean}(\mathbf{z}_{s_k})$ for $s_k \in S(D_j)$
- 20: **Step 4: Candidate Disease Retrieval**
- 21: Compute similarity scores:
 - 22: $\text{sim}(P_i, D_j) = \cos(\mathbf{z}_{P_i}, \mathbf{z}_{D_j})$
 - 23: Select top- K diseases: $\mathcal{C}_i = \text{TopK}_{D_j}(\text{sim}(P_i, D_j))$
- 24: **Step 5: LLM-based Re-ranking**
- 25: Construct prompt $\mathcal{T}_i = \text{PromptTemplate}(P_i, \mathcal{C}_i)$
- 26: Use LLM to re-rank \mathcal{C}_i : $\mathcal{D}_i^{\text{final}} = \text{LLM}(\mathcal{T}_i)$
- 27: Output top prediction: $\hat{d}_i = \mathcal{D}_i^{\text{final}}[0]$

209 involves three key steps: (1) heterogeneous graph construction with multiple relation types, (2) relation-aware embedding propagation using VGAE, and (3) candidate disease selection via patient-disease similarity. The overall workflow is illustrated in Figure 3.

214 **Multi-Relational Knowledge Graph Construction** We begin by extracting symptoms from both disease encyclopedia entries and patient inputs. Let \mathcal{D} denote the set of diseases, and \mathcal{P} denote patient

218 cases. For each disease $D_i \in \mathcal{D}$, we use an LLM-based extractor to obtain its associated symptom set $S(D_i)$ from its encyclopedia description R_i :
219
220

$$S(D_i) = \text{LLM}(R_i). \quad (1)$$

Beyond the canonical disease–symptom edges, we introduce two additional relation types to enrich the graph semantics:

- **Symptom–Symptom Edges:** If two symptoms co-occur frequently across diseases or exhibit high textual embedding similarity, we connect them with a semantic edge. Formally, for symptoms $s_a, s_b \in S^D$,

$$(s_a, s_b) \in E^{\text{sym}} \quad \text{if} \quad \text{sim}(s_a, s_b) \geq \tau,$$

where similarity is computed using BERT embeddings and τ is a fixed threshold.

- **Disease–Disease Edges:** If two diseases share a high number of overlapping symptoms (i.e., $\text{Jaccard}(S(D_i), S(D_j)) \geq \gamma$), or are commonly mentioned in clinical literature (co-morbidity), we create a disease–disease edge.

Combining all relation types, we obtain a heterogeneous graph $G = (V, E)$ with node set $V = \mathcal{D} \cup S^D$ and edge set $E = E^{\text{ds}} \cup E^{\text{sym}} \cup E^{\text{dd}}$ for disease–symptom, symptom–symptom, and disease–disease relations respectively.

Each patient input $C_i \in \mathcal{P}$ is similarly processed using an LLM extractor, followed by BERT-based symptom alignment. The normalized symptom set for patient i is:

$$\hat{S}_i^P = \left\{ S_j \in S^D \mid \max_{S_k \in \text{LLM}(C_i)} \text{sim}(S_j, S_k) \geq \beta \right\}. \quad (2)$$

Multi-Relational Embedding Propagation To capture the semantics of multi-type interactions, we apply a relation-aware Variational Graph Autoencoder (VGAE). Let $\mathbf{X} \in \mathbb{R}^{|V| \times F}$ be the initial BERT-based node feature matrix, and $\mathbf{A}^{(r)}$ the adjacency matrix for relation $r \in \{\text{ds, sym, dd}\}$. We aggregate information across R relation types using a shared GCN encoder with relation-specific propagation:

$$\mathbf{H}^{(l)} = \gamma \left(\sum_{r=1}^R \tilde{\mathbf{A}}^{(r)} \mathbf{H}^{(l-1)} \mathbf{W}^{(r)} \right), \quad (3)$$

where $\tilde{\mathbf{A}}^{(r)}$ is the normalized adjacency matrix of relation r , $\mathbf{W}^{(r)}$ is the trainable weight for relation r , and γ is a non-linear activation (e.g., ReLU).

Each node’s latent representation is modeled as a Gaussian distribution:

$$q(\mathbf{z}_i | \mathbf{X}, \{\mathbf{A}^{(r)}\}) \sim \mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2)), \quad (4)$$

and trained via a standard reconstruction loss using inner product decoding.

Patient and Candidate Representation After embedding propagation, we compute patient and disease embeddings via mean aggregation over their respective symptom representations:

$$\begin{aligned} \mathbf{z}_{D_i} &= \frac{1}{|S(D_i)|} \sum_{s_j \in S(D_i)} \mathbf{z}_{s_j}, \\ \mathbf{z}_{P_i} &= \frac{1}{|\hat{S}_i^P|} \sum_{s_j \in \hat{S}_i^P} \mathbf{z}_{s_j}. \end{aligned} \quad (5)$$

We then compute cosine similarity between the patient embedding and each disease embedding:

$$\text{sim}(P_i, D_j) = \frac{\mathbf{z}_{P_i}^\top \mathbf{z}_{D_j}}{\|\mathbf{z}_{P_i}\| \cdot \|\mathbf{z}_{D_j}\|}. \quad (6)$$

The top- K most similar diseases form the candidate set:

$$\mathcal{C}_i = \text{Top-}K(\{\text{sim}(P_i, D_j) \mid D_j \in \mathcal{D}\}). \quad (7)$$

These candidates are passed to the second-layer LLM for semantic re-ranking (Section 3.2).

3.2. Second Layer: LLM-based Disease Re-ranking

After retrieving a top- K candidate disease list from the first-layer knowledge graph, we leverage the language understanding and reasoning capabilities of

Large Language Models (LLMs) to re-rank these candidates based on patient-specific information. This step introduces contextual and semantic reasoning into the diagnostic process, allowing the model to consider nuanced symptom-disease relationships beyond structural similarity.

Prompt Construction and Optimization To enable effective re-ranking, we design structured prompts that instruct the LLM to identify the most likely disease based on the patient’s symptoms and candidate list. A base version of the prompt includes the following components:

- **Patient Symptoms:** a natural language description of the patient’s reported symptoms.
- **Candidate Diseases:** the list of top- K diseases retrieved from the knowledge graph.

We explore three prompt strategies for re-ranking candidate diseases:

- **Vanilla Prompting:** A basic template that directly asks the model to re-rank diseases based on the given symptoms.
- **Chain-of-Thought (CoT) Prompting:** Enhances reasoning by instructing the model to analyze symptoms step-by-step and eliminate unlikely candidates.
- **Auto-CoT Prompting:** Automatically generates and ranks multiple CoT-style prompts using GPT-4, selecting the best-performing variant on a held-out validation set Zhang et al. (2022).

These strategies are designed to improve the alignment between the LLM’s reasoning process and the subtle symptom variations across languages and patient descriptions.

Candidate Re-ranking Given a patient P_i with symptom description \mathcal{S}_i , and a candidate disease set \mathcal{C}_i from the first layer, the re-ranking model selects a prompt \mathcal{T}_i^* from the optimized prompt set and uses it to re-order the candidates. Formally:

$$\mathcal{D}_i^{\text{final}} = \text{LLM}(\mathcal{T}_i^*(\mathcal{S}_i, \mathcal{C}_i)), \quad (8)$$

where $\mathcal{D}_i^{\text{final}}$ is the re-ranked list of diseases output by the LLM. The final prediction \hat{d}_i is selected from the top of this list. The right panel of Figure 2 provides an overview of this re-ranking process.

In our experiments (Section 4), we demonstrate that optimized prompts lead to consistent gains in diagnostic accuracy, particularly in multilingual and high-ambiguity cases. This confirms the importance of task-specific prompt design in complex reasoning tasks such as medical diagnosis.

4. Experiment

Further dataset statistics, training details, metrics, and full ablations are in Appendix A.

4.1. Datasets

Dataset Statistics. Table 2 summarizes sizes, languages, and diagnosis diversity across splits.

Disease Knowledge Base We use a structured medical encyclopedia to construct the disease-symptom knowledge graph used in the first layer of GPT-RagAD. we adopt the following dataset:

- **Mayo Clinic Symptoms and Diseases**²: A curated dataset consisting of textual symptom descriptions for 1058 diseases collected from Mayo Clinic’s public medical knowledge base.

Multilingual Patient Symptom Data To evaluate GPT-RagAD’s multilingual diagnosis capabilities, we construct three test sets from real-world patient-doctor dialogues, including two Chinese datasets and one English dataset. Descriptive statistics are summarized in Table 2.

- **DX (Chinese)** Xu et al. (2019): A medical dialogue dataset from dxy.com. We select 104 cases among 5 diseases.
- **IMCS21 (Chinese)** Chen et al. (2023): A Chinese clinical QA dataset from Muzhi, a Baidu health platform. We select 500 pneumonia cases among 10 diseases.
- **Symptom2Disease (English)**: An English dataset containing free-text symptom descriptions from 24 diseases, curated from Kaggle. We randomly select 500 samples.

4.2. Baseline Models

Recent domain-specific LLMs, such as Med-PaLM 2 Singhal et al. (2025), PMC-LLaMA Wu et al. (2024), and ClinicalCamel Toma et al. (2023), have shown strong performance in medical reasoning tasks by leveraging large-scale supervised training on curated clinical corpora. However, such models typically require access to sensitive patient records, expensive medical annotations, or proprietary biomedical datasets—conditions that limit their scalability and practical use in privacy-sensitive or low-resource environments.

2. http://huggingface.co/datasets/celikmus/mayo_clinic_symptoms_and_diseases_v1

In contrast, GPT-RagAD is designed as a lightweight, training-free framework that builds a symbolic disease-symptom knowledge graph from publicly available medical encyclopedias (e.g., Mayo Clinic) and augments reasoning through prompt-based use of general-purpose LLMs. This hybrid design eliminates the need for fine-tuning on medical data, reduces hallucination risks through constrained candidate spaces, and generalizes across languages and LLM backbones. To assess its effectiveness, we compare GPT-RagAD against two categories of baseline systems:

- **Embedding Similarity Models:** These methods compute cosine similarity between patient symptom descriptions and disease texts using pre-trained sentence encoders. We include BERT Devlin et al. (2018), RoBERTa Liu et al. (2019), and MPNet Song et al. (2020). They serve as retrieval-only baselines without reasoning or external knowledge integration.
- **Direct Generation (LLM-only):** We evaluate GPT-4o in a direct zero-shot diagnosis setting using both vanilla and CoT-style prompts, without knowledge graph grounding or candidate filtering. This baseline reflects the upper bound of general-purpose LLM reasoning in the absence of structured symbolic support.

This comparison isolates the added value of graph-based candidate narrowing and prompt-based re-ranking in our two-layer architecture, showing how symbolic retrieval and prompt optimization complement LLM generation to improve both accuracy and reliability.

4.3. Metrics

We adopt two widely used ranking metrics:

- **Hit@K (H@K):** Measures whether the ground truth disease is among the top-K predicted results.
- **NDCG@K (N@K):** Evaluates the ranking quality of the disease list by rewarding correct predictions ranked higher.

4.4. Implementation Details

In the first layer, we generate 100 candidate diseases for each patient from the 1058-disease knowledge base using our graph-based retriever. In the second layer, we use three LLMs—LLaMA3, GPT-3.5, and GPT-4o—for re-ranking. Due to cost constraints, GPT-3.5 and GPT-4o are evaluated once on each dataset,

Table 1: Performance comparison across baseline models and our proposed GPT-RAGAD variants. Best results are in **bold**, second-best are underlined.

Dataset	Metric	BERT	RoBERTa	MPNet	GPT-RagAD (GPT-4o)		GPT-RagAD (Auto-CoT)		
					Vanilla	CoT	LLaMA3	GPT-3.5	GPT-4o
DX (Chinese)	H@1	0.1541	0.1497	0.1754	0.2663	0.2779	0.3074	0.2391	<u>0.3057</u>
	H@10	0.3054	0.2711	0.3502	0.7016	0.7141	0.7823	<u>0.7581</u>	0.7357
	H@20	0.4007	0.2741	0.3690	0.7624	0.7791	0.8075	<u>0.7973</u>	0.7901
	H@50	0.4537	0.2976	0.4321	0.7845	0.7812	0.8205	<u>0.8132</u>	0.8093
	N@1	0.1541	0.1497	0.1754	0.2663	0.2779	0.3074	0.2391	<u>0.3057</u>
	N@10	0.2169	0.2058	0.2246	0.4827	<u>0.4924</u>	0.4028	0.4125	0.5014
	N@20	0.2486	0.2114	0.2642	0.4865	<u>0.4973</u>	0.4793	0.4702	0.5104
	N@50	0.2674	0.2198	0.2819	0.4911	<u>0.5056</u>	0.4801	0.4289	0.5120
IMCS21 (Chinese)	H@1	0.1579	0.1520	0.1725	0.2754	0.2908	0.3116	0.2391	<u>0.2989</u>
	H@10	0.3158	0.2602	0.3450	0.7025	0.7189	<u>0.7625</u>	0.7743	0.7267
	H@20	0.3889	0.2690	0.3684	0.7592	0.7726	0.8039	<u>0.8013</u>	0.7846
	H@50	0.4678	0.2953	0.4415	0.7815	0.7663	<u>0.8157</u>	0.8182	0.8031
	N@1	0.1579	0.1520	0.1725	0.2754	0.2908	0.3116	0.2391	<u>0.2989</u>
	N@10	0.2366	0.2030	0.2489	0.4612	<u>0.4827</u>	0.4667	0.4001	0.4924
	N@20	0.2552	0.2051	0.2548	0.4705	<u>0.4823</u>	0.4763	0.4131	0.5066
	N@50	0.2706	0.2102	0.2703	0.4890	<u>0.4967</u>	0.4787	0.4248	0.5115
Symptom2Disease (English)	H@1	0.0588	0.0783	0.1110	0.3840	<u>0.3991</u>	0.2971	0.2824	0.4064
	H@10	0.2677	0.2481	0.2655	0.7291	0.7368	0.7566	0.7340	<u>0.7445</u>
	H@20	0.3286	0.2884	0.3025	0.8024	0.8107	0.7901	0.8171	<u>0.8159</u>
	H@50	0.3972	0.3493	0.3885	0.8116	0.8248	0.8193	0.8382	<u>0.8313</u>
	N@1	0.0588	0.0783	0.1110	0.3840	<u>0.3991</u>	<u>0.2971</u>	0.2824	0.4064
	N@10	0.1555	0.1617	0.1822	0.5412	<u>0.5615</u>	0.4650	0.4598	0.5677
	N@20	0.1709	0.1720	0.1915	0.5496	<u>0.5637</u>	0.4724	0.4791	0.5779
	N@50	0.1846	0.1838	0.2085	0.5684	<u>0.5727</u>	0.4790	0.4847	0.5835

Table 2: Dataset Statistics.

Dataset	#Diseases	#Samples	Language
DX	5	104	Chinese
IMCS21	10	500	Chinese
Symptom2Disease	24	500	English

while LLaMA3 is evaluated with three independent runs and averaged.

4.5. Main Results

Table 1 presents a comprehensive comparison of GPT-RagAD variants against embedding-based retrieval baselines (BERT, RoBERTa, MPNet) and GPT-4o direct generation (Vanilla/CoT prompting), across three multilingual datasets.

Substantial Performance Gains. GPT-RagAD consistently and significantly outperforms all baselines across datasets and evaluation metrics. On the Chinese **DX** dataset, our best variant (Auto-CoT + LLaMA3) achieves a Hit@10 of **78.2%**, more than doubling the best baseline (MPNet: 35.0%) and substantially improving over GPT-4o with vanilla prompt (70.2%). Similar patterns are observed in English: on the **Symptom2Disease** dataset, our method improves Hit@1 from 11.1% (MPNet)

to **40.6%** (Auto-CoT + GPT-4o), a nearly **3.7×** gain. This validates the core design of GPT-RagAD—combining symbolic retrieval with prompt-based neural re-ranking—as a powerful framework for multilingual, zero-shot diagnosis.

Prompting Strategy Matters. Across all datasets, prompting strategies significantly impact final performance. CoT-style prompting consistently outperforms vanilla instructions (e.g., IMCS21, NDCG@10 improves from 46.1% to 48.3% with GPT-4o), and our **Auto-CoT prompt generation** further improves results. Notably, Auto-CoT + GPT-4o achieves the best NDCG@10 across datasets (DX: **50.1%**, IMCS21: **49.2%**, Symptom2Disease: **56.8%**). This confirms the importance of context-aware reasoning and demonstrates that prompt optimization is a crucial component of LLM-based diagnosis.

LLM Portability and Scalability. While GPT-4o yields the highest overall accuracy, GPT-RagAD remains effective across diverse LLMs. The LLaMA3 variant achieves best-in-class performance on several metrics (e.g., Hit@10 = **78.2%** on DX, Hit@50 = **81.8%** on IMCS21), showing strong competitiveness despite being open-source. GPT-3.5 provides a lightweight alternative with reasonable performance

Table 3: Ablation Study: Impact of Second-Layer Re-ranking. Performance across three datasets comparing candidate retrieval (1st layer) and full re-ranking (2nd layer). Best results are in **bold**.

Dataset	Layer	Hit@K				NDCG@K			
		H@1	H@10	H@20	H@50	N@1	N@10	N@20	N@50
DX (Chinese)	1st Layer	0.0919	0.3024	0.4135	0.6445	0.0919	0.1883	0.2160	0.2615
	2nd Layer	0.3057	0.7357	0.7901	0.8093	0.3057	0.5014	0.5104	0.5120
IMCS21 (Chinese)	1st Layer	0.0893	0.2831	0.4039	0.6277	0.0893	0.1679	0.2057	0.2602
	2nd Layer	0.2989	0.7267	0.7846	0.8031	0.2989	0.4924	0.5066	0.5115
Symptom2Disease (English)	1st Layer	0.0667	0.3767	0.5301	0.7842	0.0667	0.1954	0.2339	0.2696
	2nd Layer	0.4064	0.7445	0.8159	0.8313	0.4064	0.5677	0.5779	0.5835

Table 4: Bucket-based Accuracy by Disease Type (3 Lowest Hit@1)

Disease Type	#Cases	Hit@1	Hit@10	NDCG@10
Cough	34	0.3219	0.6757	0.5039
Diarrhea	33	0.3342	0.6941	0.5227
Indigestion	28	0.3692	0.7046	0.5465

(e.g., NDCG@10 = 41.2% on DX), making the system scalable across computational budgets. This demonstrates that our framework is not tied to any specific proprietary model and can generalize across language model backbones.

Cross-lingual Robustness. GPT-RagAD performs robustly across both Chinese and English datasets without any language-specific tuning. On the Chinese **IMCS21** dataset, Auto-CoT + GPT-4o achieves 81.5% Hit@50, and on the English **Symptom2Disease** dataset, it reaches 83.1%, demonstrating high recall across languages. The model’s ability to generalize from structured KG retrieval to LLM reasoning makes it well-suited for multilingual deployment in low-resource settings.

Summary. Together, these results validate the efficacy of our two-layer design and highlight several key takeaways: (1) symbolic KG-based retrieval drastically improves LLM inference, (2) prompt engineering—especially via Auto-CoT—is critical for accurate re-ranking, and (3) GPT-RagAD generalizes well across models, languages, and data domains. Our system achieves state-of-the-art results across all benchmarks without requiring access to private patient data or task-specific fine-tuning.

4.6. Ablation Summary

We ablate (i) knowledge graph retrieval depth, (ii) re-ranker architecture, and (iii) reasoning style (Vanilla

vs. CoT). Results indicate the two-layer design contributes the largest gains, and CoT yields consistent improvements at higher K. See Appendix A for protocol details.

5. Error Analysis

To better understand failure modes in GPT-RagAD, we conduct qualitative and quantitative error analysis on the DX (Chinese) dataset. Table 10 summarizes representative failure cases, and Table 9 quantifies the proportion of errors attributed to retrieval-stage (top-100 candidate miss) and re-ranking-stage (top-10 miss) failures. We categorize four primary error sources:

(1) Semantic Overlap. Diseases with highly overlapping symptoms (e.g., *pneumonia* vs. *asthma*) often confuse both stages. Although symptom nodes are correctly matched in the knowledge graph, the LLM struggles to disambiguate without more precise contextual cues, leading to misclassification.

(2) Global Prior Bias. Even with candidate inclusion from the knowledge graph, LLM re-ranking remains biased toward globally frequent diseases. For instance, a case with canonical *herpes* symptoms was incorrectly ranked below *influenza*, reflecting residual training bias. Notably, 19.4% of errors occur despite the correct label being in the top-100 candidates—indicating re-ranking failures.

(3) Ambiguous Input. Patient descriptions with vague or nonspecific phrasing (e.g., “belly discomfort”) often lead to incorrect semantic interpretation by the LLM. These input-level ambiguities account for a notable portion of re-ranking errors, where lack of discriminative features undermines accurate reasoning.

(4) Retrieval Miss. In 12.9% of failures, the ground-truth disease was absent from the top-100

521 candidates, limiting re-ranking effectiveness regard-
 522 less of LLM capability. This is often due to symp-
 523 tom normalization mismatch or incomplete extrac-
 524 tion during KG construction—highlighting the need
 525 for improved recall in the first layer.

526 bucket-level analysis (Table 4) reveals lower ac-
 527 curacy for generalized symptom clusters like **cough**
 528 (Hit@1: 32.2%), **diarrhea** (33.4%), and **indigestion**
 529 (36.9%). These vague disease categories introduce
 530 high lexical and semantic variability, complicating
 531 both graph matching and LLM interpretation.

532 **Summary.** GPT-RagAD’s two-stage design miti-
 533 gates some limitations of end-to-end generation, but
 534 failures persist due to residual bias, retrieval incom-
 535 pleteness, and language ambiguity. Future directions
 536 include contrastive training to reduce global prior re-
 537 liance, symptom paraphrase expansion for KG recall,
 538 and more robust prompt formatting to reduce mis-
 539 alignment.

540 6. Conclusion

541 We presented **GPT-RagAD**, a two-layer retrieval-
 542 augmented diagnosis system: a graph-based genera-
 543 tor maximizes recall and an LLM re-ranker improves
 544 precision, with no task-specific training. Across DX,
 545 IMCS21, and Symptom2Disease, it outperforms neu-
 546 ral and prompting baselines on HIT@K/NDGC@K;
 547 ablations show the two-layer design yields most gains,
 548 and CoT helps at larger K (e.g., $K=100$ is a good
 549 trade-off). Remaining issues are overlapping symp-
 550 toms, sparse input, and occasional retrieval misses;
 551 next steps include stronger recall under noise, safer
 552 reasoning traces, and broader multilingual coverage.
 553 This system is for research, not clinical use.

554 7. Acknowledgement

555 This research was supported in part by NSF CNS-
 556 20-38817, DARPA HR0011-24-3-0325 (BRIES), and
 557 The Boeing Company. It also received partial sup-
 558 port from ACE, one of the seven centers in JUMP
 559 2.0, a Semiconductor Research Corporation (SRC)
 560 program sponsored by DARPA. The views and con-
 561 clusions expressed are those of the authors and do
 562 not necessarily reflect the official policies of the U.S.
 563 Government. The U.S. Government is authorized to
 564 reproduce and distribute reprints for governmental
 565 purposes notwithstanding any copyright notice.

References

- 566 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
 567 Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 568 Arvind Neelakantan, Pranav Shyam, Girish Sas-
 569 try, Amanda Askell, et al. Language models are
 570 few-shot learners. *Advances in neural information*
 571 *processing systems*, 33:1877–1901, 2020. 572
- 573 Wei Chen, Zhiwei Li, Hongyi Fang, Qianyuan Yao,
 574 Cheng Zhong, Jianye Hao, Qi Zhang, Xuanjing
 575 Huang, Jiajie Peng, and Zhongyu Wei. A bench-
 576 mark for automatic medical consultation system:
 577 frameworks, tasks and datasets. *Bioinformatics*,
 578 39(1):btac817, 2023.
- 579 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
 580 Kristina Toutanova. Bert: Pre-training of deep
 581 bidirectional transformers for language under-
 582 standing. *arXiv preprint arXiv:1810.04805*, 2018.
- 583 Wenqi Fan, Shijie Wang, Jiani Huang, Zhikai Chen,
 584 Yu Song, Wenzhuo Tang, Haitao Mao, Hui Liu,
 585 Xiaorui Liu, Dawei Yin, et al. Graph machine
 586 learning in the era of large language models (llms).
 587 *arXiv preprint arXiv:2404.14928*, 2024.
- 588 Anahita Hosseini, Ting Chen, Wenjun Wu, Yizhou
 589 Sun, and Majid Sarrafzadeh. Heteromed: Hetero-
 590 geneous information network for medical diagno-
 591 sis. In *Proceedings of the 27th ACM International*
 592 *Conference on Information and Knowledge Man-*
 593 *agement*, pages 763–772, 2018.
- 594 Thomas N Kipf and Max Welling. Semi-supervised
 595 classification with graph convolutional networks.
 596 *arXiv preprint arXiv:1609.02907*, 2016.
- 597 Patrick Lewis, Ethan Perez, Aleksandra Piktus,
 598 Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
 599 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim
 600 Rocktäschel, et al. Retrieval-augmented genera-
 601 tion for knowledge-intensive nlp tasks. *Advances in*
 602 *Neural Information Processing Systems*, 33:9459–
 603 9474, 2020.
- 604 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du,
 605 Mandar Joshi, Danqi Chen, Omer Levy, Mike
 606 Lewis, Luke Zettlemoyer, and Veselin Stoyanov.
 607 Roberta: A robustly optimized bert pretraining ap-
 608 proach. *arXiv preprint arXiv:1907.11692*, 2019.
- 609 Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhl-
 610 gay, Amnon Shashua, Kevin Leyton-Brown, and

- 611 Yoav Shoham. In-context retrieval-augmented lan- 656
 612 guage models. *Transactions of the Association for* 657
 613 *Computational Linguistics*, 11:1316–1331, 2023.
- 614 Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, 658
 615 Rianne Van Den Berg, Ivan Titov, and Max 659
 616 Welling. Modeling relational data with graph 660
 617 convolutional networks. In *The semantic web: 15th* 661
 618 *international conference, ESWC 2018, Heraklion,* 662
 619 *Crete, Greece, June 3–7, 2018, proceedings 15,* 663
 620 pages 593–607. Springer, 2018.
- 621 Weijia Shi, Sewon Min, Michihiro Yasunaga, Min- 664
 622 joon Seo, Rich James, Mike Lewis, Luke Zettle- 665
 623 moyer, and Wen-tau Yih. Replug: Retrieval- 666
 624 augmented black-box language models. *arXiv* 667
 625 *preprint arXiv:2301.12652*, 2023.
- 626 Ofir Ben Shoham and Nadav Rappoport. Cpllm: 668
 627 Clinical prediction with large language models. 669
 628 *arXiv preprint arXiv:2309.11295*, 2023.
- 629 Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, 670
 630 Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin 671
 631 Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. 672
 632 Toward expert-level medical question answering 673
 633 with large language models. *Nature Medicine*, 674
 634 pages 1–8, 2025.
- 635 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and 675
 636 Tie-Yan Liu. Mpnet: Masked and permuted pre- 676
 637 training for language understanding. *Advances in* 677
 638 *neural information processing systems*, 33:16857– 678
 639 16867, 2020.
- 640 Augustin Toma, Patrick R Lawler, Jimmy Ba, 679
 641 Rahul G Krishnan, Barry B Rubin, and Bo Wang. 680
 642 Clinical camel: An open-source expert-level medi- 681
 643 cal language model with dialogue-based knowledge 682
 644 encoding. *CoRR*, 2023.
- 645 Tao Tu, Anil Palepu, Mike Schaekermann, Khaled 683
 646 Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, 684
 647 Brenna Li, Mohamed Amin, Nenad Tomasev, 685
 648 et al. Towards conversational diagnostic ai. *arXiv* 686
 649 *preprint arXiv:2401.05654*, 2024.
- 650 Petar Veličković, Guillem Cucurull, Arantxa 687
 651 Casanova, Adriana Romero, Pietro Lio, and 688
 652 Yoshua Bengio. Graph attention networks. *arXiv* 689
 653 *preprint arXiv:1710.10903*, 2017.
- 654 Huimin Wang, Wai-Chung Kwan, Kam-Fai Wong, 690
 655 and Yefeng Zheng. Coad: Automatic diagnosis 691
 through symptom and disease collaborative gen- 692
 eration. *arXiv preprint arXiv:2307.08290*, 2023.
- Zifeng Wang, Rui Wen, Xi Chen, Shilei Cao, Shao- 693
 Lun Huang, Buyue Qian, and Yefeng Zheng. On- 694
 line disease diagnosis with inductive heterogeneous 695
 graph convolutional networks. In *Proceedings of* 696
the Web Conference 2021, pages 3349–3358, 2021.
- Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, 697
 Ya Zhang, Weidi Xie, and Yanfeng Wang. Pmc- 698
 llama: toward building open-source language mod- 699
 els for medicine. *Journal of the American Medical* 700
Informatics Association, 31(9):1833–1843, 2024.
- Lin Xu, Qixian Zhou, Ke Gong, Xiaodan Liang, Jian- 701
 heng Tang, and Liang Lin. End-to-end knowledge- 702
 routed relational dialogue system for automatic di- 703
 agnosis. In *Proceedings of the AAAI conference on* 704
artificial intelligence, volume 33, pages 7346–7353, 705
 2019.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex 706
 Smola. Automatic chain of thought prompt- 707
 ing in large language models. *arXiv preprint* 708
arXiv:2210.03493, 2022.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, 709
 Xiaolei Wang, Yupeng Hou, Yingqian Min, Be- 710
 ichen Zhang, Junjie Zhang, Zican Dong, et al. A 711
 survey of large language models. *arXiv preprint* 712
arXiv:2303.18223, 2023.

Appendix A. Additional Experimental Details

A.1. Datasets

Disease Knowledge Base We utilize a structured medical encyclopedia to construct the disease-symptom knowledge graph used in the first layer of GPT-RagAD. Specifically, we adopt the following dataset:

- **Mayo Clinic Symptoms and Diseases**³: A curated dataset consisting of textual symptom descriptions for 1058 diseases collected from Mayo Clinic’s public medical knowledge base.

Multilingual Patient Symptom Data To evaluate GPT-RagAD’s multilingual diagnosis capabilities, we construct three test sets from real-world patient-doctor dialogues, including two Chinese datasets and one English dataset. Descriptive statistics are summarized in Table 2.

- **DX (Chinese)** Xu et al. (2019): A medical dialogue dataset from dxy.com. We select 104 cases among 5 diseases.
- **IMCS21 (Chinese)** Chen et al. (2023): A Chinese clinical QA dataset from Muzhi, a Baidu health platform. We select 500 pneumonia cases among 10 diseases.
- **Symptom2Disease (English)**: An English dataset containing free-text symptom descriptions from 24 diseases, curated from Kaggle. We randomly select 500 samples.

A.2. Baseline Models

Recent domain-specific LLMs, such as Med-PaLM 2 Singhal et al. (2025), PMC-LLaMA Wu et al. (2024), and ClinicalCamel Toma et al. (2023), have shown strong performance in medical reasoning tasks by leveraging large-scale supervised training on curated clinical corpora. However, such models typically require access to sensitive patient records, expensive medical annotations, or proprietary biomedical datasets—conditions that limit their scalability and practical use in privacy-sensitive or low-resource environments.

In contrast, GPT-RagAD is designed as a lightweight, training-free framework that builds a symbolic disease-symptom knowledge graph from publicly available medical encyclopedias (e.g., Mayo

3. http://huggingface.co/datasets/celikmus/mayo_clinic_symptoms_and_diseases_v1

Clinic) and augments reasoning through prompt-based use of general-purpose LLMs. This hybrid design eliminates the need for fine-tuning on medical data, reduces hallucination risks through constrained candidate spaces, and generalizes across languages and LLM backbones. To assess its effectiveness, we compare GPT-RagAD against two categories of baseline systems:

- **Embedding Similarity Models**: These methods compute cosine similarity between patient symptom descriptions and disease texts using pre-trained sentence encoders. We include BERT Devlin et al. (2018), RoBERTa Liu et al. (2019), and MPNet Song et al. (2020). They serve as retrieval-only baselines without reasoning or external knowledge integration.
- **Direct Generation (LLM-only)**: We evaluate GPT-4o in a direct zero-shot diagnosis setting using both vanilla and CoT-style prompts, without knowledge graph grounding or candidate filtering. This baseline reflects the upper bound of general-purpose LLM reasoning in the absence of structured symbolic support.

This comparison isolates the added value of graph-based candidate narrowing and prompt-based re-ranking in our two-layer architecture, showing how symbolic retrieval and prompt optimization complement LLM generation to improve both accuracy and reliability.

A.3. Metrics

We adopt two widely used ranking metrics:

- **Hit@K (H@K)**: Measures whether the ground truth disease is among the top-K predicted results.
- **NDCG@K (N@K)**: Evaluates the ranking quality of the disease list by rewarding correct predictions ranked higher.

A.4. Implementation Details

In the first layer, we generate 100 candidate diseases for each patient from the 1058-disease knowledge base using our graph-based retriever. In the second layer, we use three LLMs—LLaMA3, GPT-3.5, and GPT-4o—for re-ranking. Due to cost constraints, GPT-3.5 and GPT-4o are evaluated once on each dataset, while LLaMA3 is evaluated with three independent runs and averaged.

Table 5: Ablation Study: Impact of Second-Layer Re-ranking. Performance across three datasets comparing candidate retrieval (1st layer) and full re-ranking (2nd layer). Best results are in **bold**.

Dataset	Layer	Hit@K				NDCG@K			
		H@1	H@10	H@20	H@50	N@1	N@10	N@20	N@50
DX (Chinese)	1st Layer	0.0919	0.3024	0.4135	0.6445	0.0919	0.1883	0.2160	0.2615
	2nd Layer	0.3057	0.7357	0.7901	0.8093	0.3057	0.5014	0.5104	0.5120
IMCS21 (Chinese)	1st Layer	0.0893	0.2831	0.4039	0.6277	0.0893	0.1679	0.2057	0.2602
	2nd Layer	0.2989	0.7267	0.7846	0.8031	0.2989	0.4924	0.5066	0.5115
Symptom2Disease (English)	1st Layer	0.0667	0.3767	0.5301	0.7842	0.0667	0.1954	0.2339	0.2696
	2nd Layer	0.4064	0.7445	0.8159	0.8313	0.4064	0.5677	0.5779	0.5835

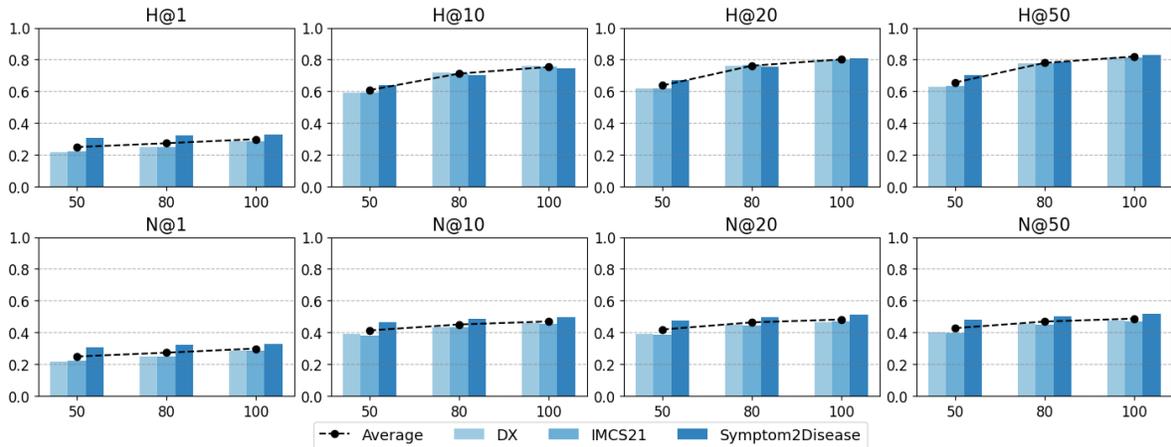


Figure 4: Impact of Candidate Set Size on Retrieval Performance Across Datasets and Metrics.

Table 6: Error Attribution Between Candidate Retrieval and Re-ranking

Error Source	Failure Rate (%)
1st-layer recall miss	12.93
2nd-layer re-ranking error	19.41

773 A.5. Ablation and Sensitivity

774 To assess the impact of candidate set size on final diagnosis performance, we conduct a sensitivity
 775 analysis by varying the number of retrieved disease
 776 candidates from 50 to 80 and 100. Results across
 777 three multilingual datasets and three LLM backbones
 778 (LLaMA3, GPT-3.5, GPT-4o) are shown in Figure 4.
 779

780 Larger Candidate Sets Yield Consistent Gains.

781 We observe consistent performance improvements as
 782 the candidate pool expands. For instance, on the Chi-
 783 nese **DX** dataset using LLaMA3, **Hit@1** increases
 784 from 17.7% (50 candidates) to 30.7% (100 candi-
 785 dates), a relative gain of 73%. Similarly, **NDCG@50**

786 rises from 36.9% to 48.0%, underscoring the benefit
 787 of enhanced candidate coverage.

788 Improvements Generalize Across Models and 789 Languages.

790 This trend holds across both English and Chinese
 791 datasets, and across all LLMs. On the English
 792 **Symptom2Disease** dataset, **GPT-4o**'s **Hit@10**
 793 improves from 65.5% to 74.4%, and **NDCG@10**
 794 increases from 52.1% to 56.8%. On the Chinese
 795 **IMCS21** dataset, **GPT-3.5**'s **Hit@50** im-
 796 proves from 63.5% to 81.8%, and **NDCG@50** in-
 797 creases from 39.8% to 42.5%. These results confirm
 798 that broader candidate pools facilitate deeper and
 799 more context-aware reasoning by the LLM.

799 Top-K Precision Metrics Benefit Most.

800 Metrics such as **Hit@1** and **NDCG@10**, which are criti-
 801 cal in high-stakes clinical scenarios, improve most
 802 significantly with larger candidate sets. For exam-
 803 ple, on Symptom2Disease with GPT-4o, **Hit@1** im-
 804 proves from 37.8% to 40.6%, and **NDCG@10** in-
 805 creases from 52.1% to 56.8%. Even modest improve-

Table 7: Qualitative Error Analysis: Sample Failure Cases in DX (Chinese) Dataset

True Label	Top-1 Prediction	Error Type	Observation
Pneumonia	Asthma	Semantic Overlap	Both diseases mention “shortness of breath” and “cough”, causing confusion.
Herpes	Influenza	Global Prior Bias	GPT favors flu due to higher frequency in corpus despite clear herpes symptoms.
Appendicitis	Gastritis	Ambiguous Input	LLM misinterpreted pain location due to vague phrasing: “belly discomfort”.
HFMD	Common Cold	Insufficient KG Match	First-layer candidate set lacked the ground-truth disease.

806 ments in top-ranked predictions can substantially en-
807 hance clinical reliability.

808 **Why 100 Candidates is the Default.** While ex-
809 panding the candidate set increases computational
810 cost, performance continues to improve up to 100
811 candidates with no sign of saturation. For example, on
812 DX with GPT-4o, **Hit@1** improves from 24.5% to
813 30.6% as the candidate size grows from 50 to 100—a
814 25% relative gain. Thus, we select 100 as the default
815 candidate size to ensure robust performance while
816 maintaining computational feasibility.

817 **Summary.** This analysis underscores the impor-
818 tance of candidate set size in enabling accurate and
819 robust multilingual medical diagnosis. The consis-
820 tent improvements observed across datasets and mod-
821 els demonstrate the scalability of GPT-RagAD’s two-
822 layer architecture, and support our choice of 100 can-
823 didates as a practical and effective design default.

824 To isolate the contribution of each component in
825 our two-layer GPT-RagAD framework, we conduct
826 an ablation study comparing the performance of (1)
827 the first-layer graph-based candidate selection alone
828 and (2) the full two-layer system, where an LLM re-
829 ranks the candidate diseases. Results are summarized
830 in Table 5 across three datasets and multiple evalua-
831 tion metrics.

832 **Second Layer Re-ranking Is Critical for High-
833 Precision Diagnosis.** Across all datasets and met-
834 rics, the second-layer LLM re-ranking yields sub-
835 stantial performance improvements over the first-
836 layer candidate retrieval alone. On the Chinese **DX**
837 dataset, **Hit@1** improves from 9.2% (1st layer) to
838 30.6% (2nd layer), a 3.3× relative gain. Similarly,
839 **NDCG@10** increases from 18.8% to 50.1%, high-
840 lighting that the re-ranking module not only improves
841 recall but significantly enhances the ordering of pre-
842 dictions, which is crucial for clinical use cases.

843 **Consistent Gains Across Languages.** This pat-
844 tern holds across datasets and languages. On the Chi-
845 nese **imcs21** dataset, **Hit@10** increases from 28.3%
846 to 72.7%, while **NDCG@50** grows from 26.0% to
847 51.1%, nearly doubling ranking quality. On the En-
848 glish **Symptom2Disease** dataset, the impact is even
849 more dramatic: **Hit@1** improves from 6.7% to 40.6%,
850 and **NDCG@10** from 19.5% to 56.8%. These con-
851 sistent gains indicate that the language model is able
852 to contextualize and reason over symptoms beyond
853 what structural similarity alone can capture.

854 **LLMs Resolve Structural Limitations of the
855 Graph Layer.** The first layer relies on structural
856 similarity between symptom nodes and diseases in the
857 knowledge graph, which is effective for broad filtering
858 but limited in semantic nuance. For example, the first
859 layer often ranks diseases with overlapping but non-
860 definitive symptoms near the top. The second-layer
861 LLM leverages natural language reasoning and global
862 medical knowledge to refine these predictions, better
863 capturing subtle linguistic cues and symptom impor-
864 tance. The boost from 53.0% to 81.6% in **Hit@20**
865 on Symptom2Disease directly reflects this capability.

866 **Conclusion.** The ablation study clearly demon-
867 strates that the second-layer LLM re-ranking is essen-
868 tial for achieving state-of-the-art performance in mul-
869 tilingual medical diagnosis. It dramatically enhances
870 both recall and ranking fidelity, making GPT-RagAD
871 more reliable in real-world applications. These find-
872 ings validate our two-layer architecture and highlight
873 the importance of combining symbolic retrieval with
874 neural language reasoning.

Table 8: Candidate Size Sensitivity Analysis on Three Datasets with Different LLM Backbones. **Bold** indicates best and underline indicates second-best performance.

Dataset	Metric	LLaMA3			GPT-3.5			GPT-4o		
		50	80	100	50	80	100	50	80	100
DX (Chinese)	H@1	0.1769	0.2389	0.3074	0.2278	0.2344	0.2391	0.2454	0.2699	<u>0.3057</u>
	H@10	0.5907	0.7052	0.7823	0.5974	0.7293	<u>0.7581</u>	0.5907	0.7125	0.7357
	H@20	0.6057	0.7532	0.8075	0.6239	0.7754	<u>0.7973</u>	0.6245	0.7596	0.7901
	H@50	0.6179	0.7698	0.8205	0.6421	0.7894	<u>0.8132</u>	0.6303	0.7698	0.8093
	N@1	0.1769	0.2389	0.3074	0.2278	0.2344	0.2391	0.2454	0.2699	<u>0.3057</u>
	N@10	0.3641	0.4007	0.4028	0.3829	0.4096	0.4702	0.4197	<u>0.4901</u>	0.5014
	N@20	0.3659	0.4512	0.4793	0.3904	0.4047	0.4125	0.4231	<u>0.4894</u>	0.5104
	N@50	0.3685	0.4503	0.4801	0.4073	0.4201	0.4289	0.4312	<u>0.4927</u>	0.5120
IMCS21 (Chinese)	H@1	0.1893	0.2419	0.3116	0.2239	0.2322	0.2391	0.2478	0.2760	<u>0.2989</u>
	H@10	0.5853	0.7064	<u>0.7625</u>	0.6047	0.7357	0.7743	0.5874	0.7038	0.7267
	H@20	0.6121	0.7701	0.8039	0.6166	0.7711	<u>0.8013</u>	0.6290	0.7522	0.7846
	H@50	0.6345	0.7754	<u>0.8157</u>	0.6345	0.7807	0.8182	0.6345	0.7665	0.8031
	N@1	0.1893	0.2419	0.3116	0.2239	0.2322	0.2391	0.2478	0.2760	<u>0.2989</u>
	N@10	0.3556	0.4320	0.4667	0.3801	0.3944	0.4001	0.4128	<u>0.4715</u>	0.4924
	N@20	0.3605	0.4436	0.4763	0.3870	0.4091	0.4131	0.4128	<u>0.4847</u>	0.5066
	N@50	0.3666	0.4452	0.4787	0.3976	0.4131	0.4248	0.4257	<u>0.4883</u>	0.5115
Symptom2Disease (English)	H@1	0.2735	0.3003	0.2971	0.2747	0.2759	0.2824	0.3784	<u>0.3854</u>	0.4064
	H@10	0.6299	0.7103	0.7566	0.6264	0.6848	0.7340	0.6550	0.7141	<u>0.7445</u>
	H@20	0.6645	0.7600	0.7901	0.6719	0.7591	0.8171	0.6795	0.7461	<u>0.8159</u>
	H@50	0.7008	0.7804	0.8193	0.7008	0.7928	0.8382	0.7008	0.7842	<u>0.8313</u>
	N@1	0.2735	0.3003	0.2971	0.2747	0.2759	0.2824	0.3784	<u>0.3854</u>	0.4064
	N@10	0.4402	0.4575	0.4650	0.4385	0.4534	0.4598	0.5213	<u>0.5422</u>	0.5677
	N@20	0.4473	0.4676	0.4724	0.4519	0.4710	0.4791	0.5272	<u>0.5496</u>	0.5779
	N@50	0.4557	0.4728	0.4790	0.4597	0.4794	0.4847	0.5317	<u>0.5575</u>	0.5835

Table 9: Error Attribution Between Candidate Retrieval and Re-ranking

Error Source	Failure Rate (%)
1st-layer recall miss	12.93
2nd-layer re-ranking error	19.41

Appendix B. Extended Error Analysis

Appendix C. Error Analysis

To better understand failure modes in GPT-RagAD, we conduct qualitative and quantitative error analysis on the DX (Chinese) dataset. Table 10 summarizes representative failure cases, and Table 9 quantifies the proportion of errors attributed to retrieval-stage (top-100 candidate miss) and re-ranking-stage (top-10 miss) failures. We categorize four primary error sources:

(1) Semantic Overlap. Diseases with highly overlapping symptoms (e.g., *pneumonia* vs. *asthma*) often confuse both stages. Although symptom nodes are correctly matched in the knowledge graph, the

LLM struggles to disambiguate without more precise contextual cues, leading to misclassification.

(2) Global Prior Bias. Even with candidate inclusion from the knowledge graph, LLM re-ranking remains biased toward globally frequent diseases. For instance, a case with canonical *herpes* symptoms was incorrectly ranked below *influenza*, reflecting residual training bias. Notably, 19.4% of errors occur despite the correct label being in the top-100 candidates—indicating re-ranking failures.

(3) Ambiguous Input. Patient descriptions with vague or nonspecific phrasing (e.g., “belly discomfort”) often lead to incorrect semantic interpretation by the LLM. These input-level ambiguities account for a notable portion of re-ranking errors, where lack of discriminative features undermines accurate reasoning.

(4) Retrieval Miss. In 12.9% of failures, the ground-truth disease was absent from the top-100 candidates, limiting re-ranking effectiveness regardless of LLM capability. This is often due to symptom normalization mismatch or incomplete extrac-

Table 10: Qualitative Error Analysis: Sample Failure Cases in DX (Chinese) Dataset

True Label	Top-1 Prediction	Error Type	Observation
Pneumonia	Asthma	Semantic Overlap	Both diseases mention “shortness of breath” and “cough”, causing confusion.
Herpes	Influenza	Global Prior Bias	GPT favors flu due to higher frequency in corpus despite clear herpes symptoms.
Appendicitis	Gastritis	Ambiguous Input	LLM misinterpreted pain location due to vague phrasing: “belly discomfort”.
HFMD	Common Cold	Insufficient KG Match	First-layer candidate set lacked the ground-truth disease.

911 tion during KG construction—highlighting the need
912 for improved recall in the first layer.

913 Additionally, bucket-level analysis (Table 4) reveals
914 lower accuracy for generalized symptom clusters like
915 **cough** (Hit@1: 32.2%), **diarrhea** (33.4%), and **in-**
916 **digestion** (36.9%). These vague disease categories
917 introduce high lexical and semantic variability, com-
918 plicating both graph matching and LLM interpreta-
919 tion.

920 **Summary.** GPT-RagAD’s two-stage design miti-
921 gates some limitations of end-to-end generation, but
922 failures persist due to residual bias, retrieval incom-
923 pleteness, and language ambiguity. Future directions
924 include contrastive training to reduce global prior re-
925 liance, symptom paraphrase expansion for KG recall,
926 and more robust prompt formatting to reduce mis-
927 alignment.

928 Appendix D. Extended Related Work

929 D.1. Graph-based Diagnosis

930 Graph-based AD systems [Hosseini et al. \(2018\)](#);
931 [Wang et al. \(2021\)](#) typically encode medical records
932 into Heterogeneous Information Networks (HINs)
933 to capture structured interactions among diseases,
934 symptoms, and other clinical entities. While effective
935 in structured-data scenarios, these systems depend on
936 high-quality, large-scale EHRs, which constrains gen-
937 eralization across domains or languages.

938 D.2. LLM-based Diagnosis

939 LLM-based approaches [Shoham and Rappoport](#)
940 [\(2023\)](#); [Wang et al. \(2023\)](#); [Tu et al. \(2024\)](#) lever-
941 age pretrained transformers for direct disease predic-
942 tion or symptom synthesis. Although flexible and
943 annotation-free, such methods are prone to global

944 prior bias and often lack explicit reasoning over struc-
945 tured medical knowledge. Prior work commonly fine-
946 tunes LLMs on curated datasets, which limits scal-
947 ability and can reduce transparency.

948 D.3. Retrieval-Augmented Generation (RAG)

949 RAG enhances LLMs by retrieving external knowl-
950 edge to ground generation, and has proven effec-
951 tive in various NLP tasks [Lewis et al. \(2020\)](#); [Ram](#)
952 [et al. \(2023\)](#); [Shi et al. \(2023\)](#). However, most RAG
953 pipelines retrieve free-text passages, which may intro-
954 duce hallucinations, increase latency, and raise pri-
955 vacy concerns—issues that are especially salient in
956 clinical contexts.
957

958 D.4. Positioning of Our Work

959 **GPT-RagAD** advances RAG-based diagnosis by
960 introducing a symbolic retrieval layer over a dis-
961 ease-symptom knowledge graph constructed from
962 public medical encyclopedias, rather than raw text
963 or private records. The graph is embedded with a
964 Variational Graph Autoencoder (VGAE) to model
965 uncertainty in symptom-disease associations, and
966 candidates are re-ranked with an LLM via op-
967 timized prompting. This design balances inter-
968 pretability, scalability, and accuracy, while support-
969 ing privacy preservation and cross-lingual general-
970 ization—capabilities under-explored in prior AD or
971 RAG literature.