

CoD, Towards an Interpretable Medical Agent using Chain of Diagnosis

Junying Chen^{1,2†}, Chi Gui^{2†}, Anningzhe Gao^{1,2*}, Ke Ji^{1,2}, Xidong Wang^{1,2}
Xiang Wan^{1,2}, Benyou Wang^{1,2*}

¹ Shenzhen Research Institute of Big Data

² The Chinese University of Hong Kong, Shenzhen

<https://github.com/FreedomIntelligence/Chain-of-Diagnosis>

Abstract

The field of medical diagnosis has undergone a significant transformation with the advent of large language models (LLMs), yet the challenges of interpretability within these models remain largely unaddressed. This study introduces **Chain-of-Diagnosis (CoD)** to enhance the interpretability of LLM-based medical diagnostics. CoD transforms the diagnostic process into a diagnostic chain that mirrors a physician’s thought process, providing a transparent reasoning pathway. Additionally, CoD outputs the disease confidence distribution to ensure transparency in decision-making. This interpretability makes model diagnostics controllable and aids in identifying critical symptoms for inquiry through the entropy reduction of confidences. With CoD, we developed **DiagnosisGPT**, capable of diagnosing 9,604 diseases. Experimental results demonstrate that DiagnosisGPT outperforms other LLMs on diagnostic benchmarks. Moreover, DiagnosisGPT provides interpretability while ensuring controllability in diagnostic rigor.

1 Introduction

In healthcare, disease diagnosis is pivotal yet complex, serving as a bridge between medical expertise and decision-making [1, 2]. The diagnosis task involves predicting a disease using explicit (self-reported) and implicit (inquired) symptoms [3–5]. As illustrated as Figure 1, the doctor either inquires further or makes a diagnosis, aiming for high accuracy with minimal inquiries. Large Language Models (LLMs) offer a promising path for automated diagnosis, due to their superior reasoning and dialogue abilities, coupled with extensive knowledge. These capabilities enable them to address a wide range of diseases and interact effectively with patients [6].

However, the application of LLMs in medical diagnosis encounters significant challenges of **interpretability** [7, 8]. Considering the issue of hallucinations, LLMs could arbitrarily make a diagnosis. Without interpretability, it is unclear if decisions meet sound analysis and ethical standards [9]. Although LLMs can offer rudimentary explanations for their decision, they lack a comprehensive process to explain why other potential diseases are excluded and to which extent of confidence it made such a decision. [10] This highlights the need for an interpretable LLM solution for diagnosis.

In response to these limitations, we propose the **Chain of Diagnosis (CoD)** to enhance the interpretability of LLMs. CoD provides transparency in both reasoning and decision-making. It transforms the black-box decision-making process into a diagnostic chain that mirrors a physician’s thinking process with five distinct steps. It reveals the series of thoughts behind each decision. For decision transparency, CoD outputs a confidence distribution, with higher confidence indicating a stronger belief in diagnosing a specific disease. This allows for control over the LLM’s decisions using a confidence threshold. Additionally, diagnostic uncertainty can be quantified by the entropy of the confidences. The goal of entropy reduction can aid in eliciting more effective symptoms for inquiry.

*Benyou and Anningzhe are the corresponding authors. The first two authors contributed to this work equally.

To implement CoD, we propose a method for constructing CoD training data from patient cases. Due to the restriction of real-world cases, this paper proposes generating synthetic cases from disease encyclopedias. This approach facilitates scalability in supported diseases without any concerns of patient privacy. With synthetic cases, we constructed a training dataset with 48,020 CoD instances, leading to the development of our model, **DiagnosisGPT**, capable of diagnosing 9,604 diseases. In out-of-domain settings, DiagnosisGPT outperforms other LLMs on public diagnostic datasets and a newly-created benchmark DxBench. Moreover, it achieves an accuracy exceeding 90% on all datasets with a diagnostic threshold of 0.55, highlighting the reliability of its confidence levels.

Our contributions are summarized as follows: 1) We introduce the Chain-of-Diagnosis (CoD) diagnostic method, designed to enhance interpretability of LLMs in disease diagnosis; 2) We propose a method to synthesize patient cases using disease encyclopedia data. This enables low-cost creation of CoD training data for various diseases while avoiding privacy and ethical concerns; 3) Using CoD, we built DiagnosisGPT that can support automatic diagnosis for 9,604 diseases. Experiments demonstrate that DiagnosisGPT surpasses other LLMs across diverse diagnostic datasets; 4) We present DxBench, a diagnostic benchmark with 1,148 real cases covering 461 diseases, manually verified and derived from public doctor-patient dialogues.

2 Problem Definition of Diagnosis

2.1 Problem definition

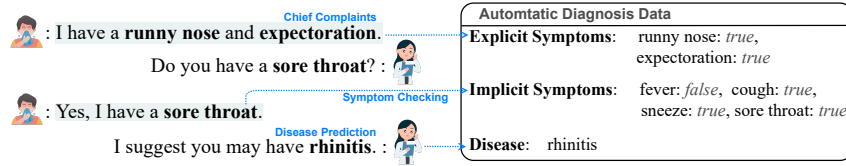


Figure 1: Example of the automatic diagnosis task, with sample data from 2.

In the diagnostic process, physicians navigate through **explicit symptoms** (\mathcal{S}_{exp}), self-reported by patients, and **implicit symptoms** (\mathcal{S}_{imp}), elicited through further inquiry, to predict a target disease (d_t) from a predefined list (\mathcal{D}). This process involves a sequence of decisions: physicians either probe for additional symptoms (**Symptom Inquiry**) or select the most probable disease from \mathcal{D} (**Disease Diagnosis**). The objective is to maximize diagnostic accuracy ($a \uparrow$) while minimizing the number of inquiries ($n \downarrow$), ensuring the inquiry count n does not exceed a patient’s patience threshold L .

2.2 The Challenge

The challenge lies in determining when and how to inquire about symptoms to enhance diagnostic accuracy. As stated in [11], LLMs don’t excel at questioning users. To explore this, we conducted a preliminary experiment on diagnostic benchmarks with two LLMs:

Table 1: Diagnostic results with two different LLMs, with the maximum number of inquiries of $L = 5$. The **green** and **red** fonts indicate the increase and decrease in accuracy after using symptom inquiries, respectively. The details of experiments can be found in Appendix D and C.

	Muzhi Dataset			Dxy Dataset		
	w/o inquiry a	w/ inquiry a	n	w/o inquiry a	w/ inquiry a	n
Gemma-7b-it	49.3	45.8 -3.5	5.0	33.7	31.7 -2.0	5.0
GPT-3.5-turbo-1106	56.3	55.6 -0.7	0.2	45.6	46.2 +0.6	0.4

The results shown in Table 1 reveals two possible issues for diagnosis: **Case I, Arbitrary Diagnosis (too small n):** LLMs diagnose immediately without any symptom inquiry, leading to minimal benefit from symptom inquiry; **Case II, Excessive Inquiries (too large n):** LLMs continuously inquire about symptoms until reaching the maximum limit $n = L$, resulting in inefficient diagnosis.

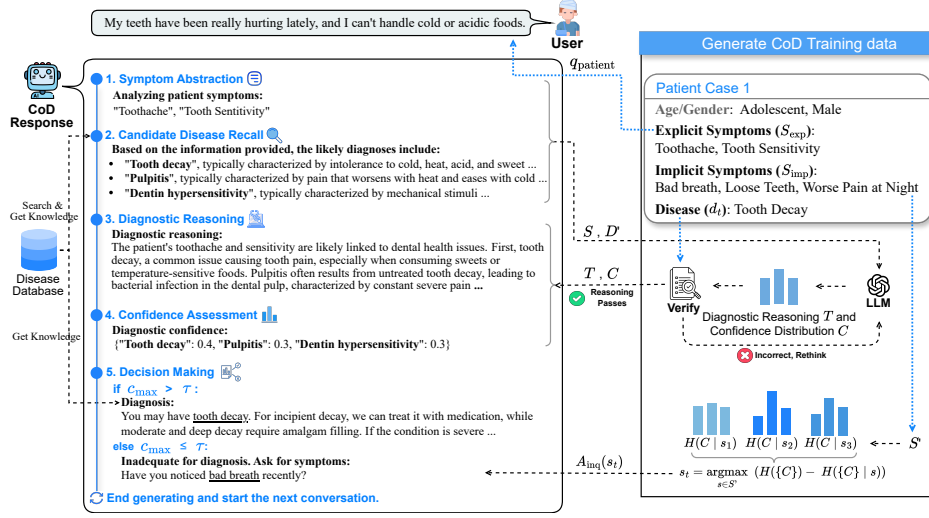


Figure 2: Left: Example of a CoD response. Right: Construction of CoD training data.

The objective of diagnosis is to enhance disease identification accuracy through symptom inquiry. However, there is a trade-off between efficiency and accuracy. Efficient diagnosis is faster with fewer inquiries, while accurate diagnosis requires more detailed inquiries. LLMs may struggle with symptom inquiry, particularly in determining when to ask questions. Additionally, symptom inquiry by LLMs does not significantly improve accuracy, indicating that LLMs might fail to ask for crucial symptom information.

3 Methodology: Chain of Diagnosis

As depicted on the left side of Figure 2, the CoD outputs a diagnostic chain that represents the thought process of LLM diagnosis. It mirrors a physician’s diagnostic thinking, which involves summarizing symptoms, identifying diseases, analyzing and deciding. This interpretability aids in diagnostic decision-making. To implement the CoD, we constructed CoD training data based on patient cases to fine-tune LLMs to perform CoD, as shown on the right side of Figure 2.

3.1 The Philosophy of CoD for Interpretability

Lipton [12] classified *interpretability* into two aspects: 1) transparency, i.e., *how does the model work?* and 2) post-hoc explanations, i.e., *what can the learned model tell us?* The two aspects inspire us design an interpretable diagnostic LLM, consisting of a CoD framework that includes Property 1, 2 for transparency and Property 3 for post-hoc explanations.

Transparency [12] encompasses two types: *decomposability* and *algorithmic transparency*. *Decomposability* refers to a system where each component can be individually interpreted. For CoD, this is demonstrated in Property 1.

Property 1 Decomposability with a pipeline chain: *CoD transform the black-box decision of diagnosis into a explainable chain, which can be viewed as a sequence of intermediate steps. Therefore, the designed CoD is a chain pipeline (see Sec. 3.2), each step of which has a clear functionality. The overall chain mimics a real diagnosis of a physician.*

Algorithmic transparency refers to the understanding of the learning algorithm itself, such as *whether it converges and why*. Regarding the challenge in Sec. 2.2, the algorithmic transparency of CoD can be understood from an entropy-reduction perspective: with more inquiries made, the uncertainty of the diagnosis estimate (quantitatively measured by entropy) will be reduced, see Property 2.

Property 2 Transparency with confidence-driven flow: *In the chain, CoD introduces a disease confidence distribution $C = \{c_d \mid d \in \mathcal{D}\}$, where decision-making is based on whether the highest*

confidence exceeds a given threshold τ . By using explicit confidence levels in CoD, one can easily observe that the accuracy becomes larger with increasing inquiries thanks to the improved certainty (i.e., reduced entropy), see Sec. 3.3. It converges when accuracy saturates with enough inquiries.

Post-hoc explanations [12] indicate the information and functions a model can offer to humans. The post-hoc explanations for CoD are illustrated in Property 3.

Property 3 *Diagnosis explainability*: CoD elucidates the diagnostic thinking process, providing physicians with a diagnostic pathway that supports their clinical decisions and ensures that the LLM’s decisions adhere to reasonable analysis and ethical standards.

3.2 The Diagnosis Chain

Here, we introduce the response methods and the construction approach of CoD, as illustrated in Figure 2. All prompts for building CoD training data are detailed in Appendix F.

Step 1: Symptom Abstraction The first step summarizes the symptoms \mathcal{S} of the patient’s question:

$$\mathcal{S} = f_1(q_{\text{patient}}) \quad (1)$$

It allow the model to focus on the refined symptoms and provide an understanding of patient’s query. For training data, the initial patient question is generated from $(\mathcal{S}_{\text{exp}})$ with the LLM.

Step 2: Disease Recall & Knowledge Integration

Next, CoD identifies the top-K potential diseases based on a disease retriever:

$$\mathcal{D}' = f_2(\mathcal{D}, \mathcal{S}, k) \quad (2)$$

where $\mathcal{D}' \subseteq \mathcal{D}$ and $|\mathcal{D}'| = k$. A smaller space \mathcal{D}' is necessary for subsequent analysis and reasoning, since analyzing all diseases is impractical (considering $|\mathcal{D}| = 9604$) and most irrelevant diseases can realistically be excluded. We use Dense Retrieval training methods [13, 14] to train this retriever, with the following training objective:

$$\mathcal{L}(\mathcal{S}_{\text{exp}}, \mathcal{S}_{\text{imp}}, d_t) = -\log \frac{e^{\text{sim}(E_S(\mathcal{S}_{\text{exp}} \cup \mathcal{S}_{\text{imp}}), E_D(d_t))}}{\sum_{d \in \mathcal{D}} e^{\text{sim}(E_S(\mathcal{S}_{\text{exp}} \cup \mathcal{S}_{\text{imp}}), E_D(d))}} \quad (3)$$

where sim denotes the cosine similarity, and E_S and E_D are the symptom and disease encoders, respectively. The performance of the disease retriever is detailed in Appendix I.

Then, for each candidate disease $d \in \mathcal{D}'$, CoD retrieves corresponding disease knowledge from the disease database and integrates it into the output to enhance understanding of the disease. Similarly, other tools like RAG can also be utilized in this step to enhance reasoning.

Step 3: Diagnostic Reasoning In step 3, CoD generates the diagnostic reasoning process T :

$$T = f_3(\mathcal{S}, \mathcal{D}') \quad (4)$$

Similar to CoT, T is a thought process that carefully analyzes whether each disease in \mathcal{D}' corresponds to the patient’s symptoms. To build training data, we prompt a LLM to generate T .

Step 4: Confidence Assessment After generating T , CoD generates a confidence distribution:

$$\mathcal{C} = f_4(\mathcal{S}, \mathcal{D}', T) \quad (5)$$

\mathcal{C} satisfies $\sum_{d \in \mathcal{D}'} c_d = 1$. This distribution indicates the model’s tendency towards diagnosing a disease, mainly according to the analysis of T . According to f_3 , \mathcal{C} can be considered a posterior probability distribution:

$$\mathcal{C} = \{p_\theta(d|\mathcal{S}, \mathcal{D}') | d \in \mathcal{D}'\} \quad (6)$$

Here, p_θ represents the confidence distribution generated by the LLM θ . For constructing training data, we validate \mathcal{C} against the target disease d_t to ensure T and \mathcal{C} are reasonable. If $\max_{d \in \mathcal{D}' \setminus \{d_t\}} c_d \geq \tau$, the generated data is considered erroneous, i.e., the model assigns high confidence to an incorrect disease. If erroneous, we prompt the model to rethink and correct its reasoning until the distribution is verified. With \mathcal{C} , CoD can make decisions based on the confidence in its diagnosis.

Step 5: Decision Making In the last step, a confidence threshold τ is set to control the decision-making. The diagnostic task involves two decision types: 1) making a diagnosis $A_{\text{diag}}(d)$, where d is

the diagnosed disease, and 2) to inquiring about a symptom $A_{inq}(s)$, where s represents the symptom under inquiry. The next decision A_{next} of the CoD is defined as:

$$A_{next} = \begin{cases} A_{diag}(d_{\max}), & \text{if } c_{\max} > \tau \\ A_{inq}(s_t), & \text{if } c_{\max} \leq \tau \end{cases} \quad (7)$$

where $c_{\max} = \max_{d \in \mathcal{D}'} \{c_d\}$ and $d_{\max} = \operatorname{argmax}_{d \in \mathcal{D}'} \{c_d\}$. $A_{inq}(s_t)$ signifies the operation of querying about the symptom s_t that the CoD generates. Here, τ serves as a hyperparameter. A higher τ allows the model to perform more rigorous diagnoses (that achieving higher accuracy a but requiring more rounds of questioning, i.e., higher n). Conversely, a lower τ can reduce n but also lowers a .

3.3 CoD as an Entropy-reduction Process

Symptom inquiry is a key step in diagnosis, serving to gather additional patient information to clarify the diagnosis. This inquiry process can be viewed as a transition from diagnostic uncertainty to certainty. The uncertainty level can be captured by the entropy of confidence:

$$H(C) = - \sum_{d \in \mathcal{D}'} c_d \log c_d \quad (8)$$

Symptom inquiry is a process of entropy reduction. Given a symptom s , its post-inquiry entropy is:

$$H(C|s) = - \sum_{d \in \mathcal{D}'} p_{\theta}(d|S \cup \{s\}, D') \log p_{\theta}(d|S \cup \{s\}, D') \quad (9)$$

For the diagnostic task, it is preferable to minimize n . Hence, the objective of symptom inquiry can be formalized as maximizing the increase in diagnostic certainty to expedite the diagnosis. Accordingly, CoD selects the symptom to inquire about by maximizing the entropy reduction:

$$s_t = \operatorname{argmax}_{s \in \mathcal{S}'} (H(C) - H(C|s)) \quad (10)$$

where \mathcal{S}' represents the candidate symptoms for inquiry and s_t is the chosen symptom. $\mathcal{S}' = \mathcal{S}_{\text{imp}} \cup \{s_{\text{gen}}\}$, where s_{gen} is the symptom generated by the LLM and \mathcal{S}_{imp} comes from the training case data. Through entropy reduction, the CoD training data tuned the model to inquire about more crucial symptoms for diagnosis, thereby enhancing its querying capability.

3.4 Synthesizing Patient Cases

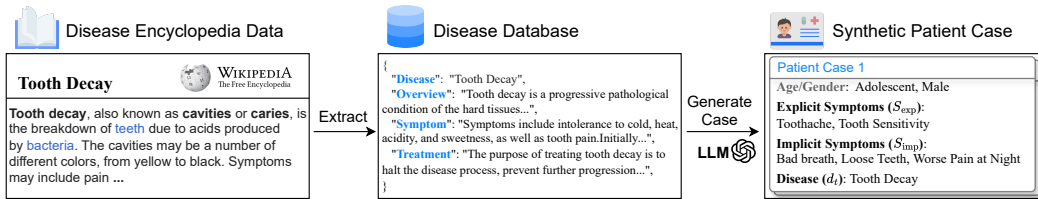


Figure 3: Schematic of constructing disease database and synthesizing patient cases.

CoD requires patient cases to build training data. However, due to privacy concerns, the collection of such data is significantly restricted. To address this, we propose generating synthetic case data in reverse from online disease encyclopedias, which provide comprehensive and reliable disease information. As illustrated in Figure 3, the synthesis process is a pipeline consists of two stages:

Stage 1: Constructing Disease Database The first step involves the extraction of essential information from the disease encyclopedia data. This process results in a knowledge base encompassing 9,604 diseases, each detailed with sections on "Overview," "Symptoms," and "Treatment". We use regular expression matching to identify and extract these key sections.

Stage 2: Synthesizing Patient Cases In disease diagnosis [3, 15], a patient can be abstracted into a triplet (S_{exp}, S_{imp}, d_t) . Using the LLM, we generate structured case data based on the disease

knowledge from the database. For each disease, we synthesize five distinct cases to ensure diversity. The prompt used for generation is provided in the Appendix E.

In the end, we developed a database containing **9,604** diseases and then synthesized **48,020** unique cases. The used LLM is **GPT-4-0125-preview**. Medical experts reviewed the quality of the synthetic cases, finding potential errors in only 6 out of 100 cases, which confirms the reliability of the synthetic cases. For more details, see Appendix K. Based on these synthetic cases, we constructed a training dataset for medical diagnosis, which consists of 48,020 instances with an average of 2.4 consultation rounds. This dataset is used to train an interpretable medical diagnosis model, **DiagnosisGPT**.

4 Experiments

4.1 Model Training & Setup

Utilizing the created CoD data, we fully fine-tuned the **Yi-34B-Base** [16] to develop **DiagnosisGPT**. To equip it with chat capabilities, ShareGPT data² is incorporated into the training data. Training parameters included a batch size of 64 and a learning rate of $2e-5$. For the retrieval model, we trained on the all-mpnet-base-v2 [17] model using DRhard [18], with a batch size of 256 and a learning rate of $2e-5$. The training was conducted on a GPU server with 8 NVIDIA A100 GPU cards.

4.2 Benchmarking Settings

Traditional baselines (Non-LLM)

Traditional supervised Automatic Diagnosis methods approach the diagnostic task as a decision-making task, where all symptoms and diseases are predefined. In traditional methods, we adhere to the original settings, which involve training on a training set of benchmarks. We compared four models: Basic DQN [15], HRL [19], Diaformer [20] and MTDiag [2].

LLM baselines Our comparison mainly focused on advanced LLMs including proprietary models like Gemini-Pro [21], ERNIE Bot(文心一言) [22], Claude-3-Opus [23], GPT-3.5 [24] (GPT-3.5-turbo-1106), and GPT-4 [25] (GPT-4-0125-preview), as well as open-source models including gemma-7b-it [26], Mixtral-8x7B-Instruct-v0.1 [27], and Yi-34B-Chat [16].

LLM Evaluation We simulate a patient using GPT-4, which presents both \mathcal{S}_{exp} and \mathcal{S}_{imp} . The simulation begins with \mathcal{S}_{exp} (chief complaints). When the evaluated LLM inquires about symptoms, the simulator can only respond with "yes" or "no" to prevent information leakage. Details of the LLM evaluation can be found in Appendix D. For the evaluated LLMs, we prompt them to perform an automated diagnosis task, which is detailed in Appendix C.

4.3 Benchmark

Public benchmarks To evaluate diagnostic performance, we used two publicly available benchmarks: Muzhi [15] and Dxy [4]. Both are based on real doctor-patient consultations. However, their data scale and disease variety are limited, as shown in Table 2.

DxBench To better assess diagnostic capabilities, we develop a larger dataset, DxBench. Using the MedDialog [28] dataset, which contains real doctor-patient dialogues, we filtered out 3,121 cases with clear dialogues and definitive diagnoses. Then GPT-4 is employed to extract \mathcal{S}_{exp} and \mathcal{S}_{imp} , and we manually refine this to 1,148 high-quality cases. Details are in Appendix G. DxBench includes over 1,000 real cases, covering 461 disease types from 15 departments and 5,038 symptoms. Considering the large number of diseases in DxBench, each case is provided with three candidate diseases.

4.4 Diagnosis Performance

As shown in Table 3, we evaluated diagnostic performance using two public benchmarks: Muzhi [15] and Dxy [4]. The results indicate that LLMs' performance approach traditional IID methods while requiring fewer inquiries. Among the LLMs, DiagnosisGPT achieves the highest accuracy improvement with symptom inquiries, demonstrating superior inquiry capabilities. In the Muzhi

²<https://huggingface.co/datasets/philschmid/sharegpt-raw>

Table 2: Comparison of DxBench with other datasets.

Dataset	# Disease	# Symptom	# Test Data	# Department
MuZhi	4	66	142	1
Dxy	5	41	104	1
DxBench	461	5038	1148	15

Table 3: Dxy and Muzhi Benchmark results. τ denotes the confidence threshold. **Acc.** stands for the accuracy of disease diagnosis (percentage). n represents the average number of inquiry rounds. All experiments were limited to $L = 5$ inquiry rounds. The **green** and **red** fonts indicate the increase and decrease in accuracy after using symptom inquiries.

	Muzhi Dataset (4 Candidate Diseases)			Dxy Dataset (5 Candidate Diseases)			Avg Acc.
	w/o inquiry	w/ inquiry	n	w/o inquiry	w/ inquiry	n	
	Acc.	Acc.		Acc.	Acc.		
Traditional Methods (Supervised Learning)							
Basic DQN	-	64.1	2.9	-	64.7	2.5	64.4
HRL	-	67.6	2.8	-	70.2	1.9	68.9
Diaformer	-	72.2	5.0	-	76.6	4.8	74.4
MTDiag	-	72.6	5.0	-	76.1	5.0	74.4
Large Language Models (Out-of-domain Setting)							
gemma-7b-it	49.3	45.8 _{-3.5}	5.0	33.7	31.7 _{- 2.0}	5.0	38.8
Yi-34B-Chat	52.1	50.7 _{-1.4}	0.4	52.9	50.5 _{- 2.4}	0.5	50.5
GPT-3.5-turbo-1106	56.3	55.6 _{-0.7}	0.2	45.6	46.2 _{+ 0.6}	0.4	50.9
Mixtral-8x7B-Instruct-v0.1	56.3	50.0 _{-6.3}	1.9	47.1	55.8 _{+ 8.7}	1.7	52.9
ERNIE Bot	61.3	57.0 _{-4.3}	0.4	51.9	51.9 _{+ 0.0}	0.8	54.5
Gemini-Pro	63.4	60.6 _{-2.8}	0.2	57.7	56.7 _{- 1.0}	0.1	58.7
GPT-4-0125-preview	59.2	56.3 _{-3.6}	0.4	62.5	65.4 _{+ 2.9}	0.6	60.8
Claude-3-Opus	57.7	56.3 _{-1.4}	1.8	62.5	73.1 _{+10.6}	1.8	64.7
DiagnosisGPT ($\tau = 0.5$)	62.0	64.1 _{+2.1}	1.4	60.5	72.6 _{+12.1}	1.5	68.4
DiagnosisGPT ($\tau = 0.6$)	62.0	65.5 _{+3.5}	2.4	60.5	75.4 _{+14.9}	2.8	70.5

dataset, DiagnosisGPT is the only model that improves accuracy through symptom inquiries. In the Dxy dataset, DiagnosisGPT required fewer inquiries than Claude-3-Opus but achieved a greater improvement. It can also be observed that Gemini-Pro and GPT-3.5-turbo-1106 tend not to ask questions. Overall, DiagnosisGPT ($\tau = 0.6$) delivered the best results across the datasets, owing to its inquiry and decision-making capabilities.

The results of DxBench are shown on the Figure 4. DiagnosisGPT achieved the highest accuracy with the $\tau = 0.6$ setting. At $\tau = 0.5$, it outperformed Claude-3-Opus by 3.7% with the same inquiry round n . Since DiagnosisGPT supports diagnoses for 9604 diseases, it can perform open-ended consultations without relying on candidate diseases. In this scenario, DiagnosisGPT’s consultation process significantly improved diagnostic accuracy from 34.7% to 44.2%. Although the accuracy is only 44.2%, DiagnosisGPT demonstrates the potential of LLMs to identify the correct diagnosis from a large number of diseases. The right side of Figure 4 shows the accuracy across different departments, indicating that DiagnosisGPT performs well across all departments.

5 More Analysis

5.1 Interpretability on Confidence Levels

To assess the interpretability of the diagnostic confidence, we examined diagnostic accuracy at various confidence thresholds. The results, depicted in Figure 5, show that increasing the threshold indeed enhances accuracy. With $\tau = 0.55$, the model achieves over 90% accuracy across three datasets. This demonstrates that the model’s confidence in disease prediction is dependable and aligns with the expected accuracy rates. However, higher thresholds reduce success rates, indicating that the model becomes more stringent to make a diagnosis.

	DxBench		
	w/o inquiry Acc.	w/ inquiry Acc.	n
With Candidate Diseases (Given 3, 1 is correct)			
gemma-7b-it	37.5	37.2 -0.3	5.0
Mixtral-8x7B-Instruct-v0.1	42.1	41.2 -0.9	1.4
GPT-3.5-turbo-1106	52.1	52.3 $+0.2$	0.1
Yi-34B-Chat	52.4	54.6 $+2.2$	0.2
Gemini-Pro	55.3	55.1 -0.2	0.1
ERNIE Bot	58.4	59.9 $+1.5$	0.4
Claude-3-Opus	59.4	62.1 $+2.7$	1.8
GPT-4-0125-preview	60.1	62.5 $+2.4$	0.6
DiagnosisGPT ($\tau = 0.5$)	56.9	62.4 $+5.5$	1.8
DiagnosisGPT ($\tau = 0.6$)	56.9	64.5 $+7.6$	3.1
Without Candidate Diseases (Diagnose from 9604 diseases)			
DiagnosisGPT ($\tau = 0.5$)	34.7	41.4 $+6.7$	1.6
DiagnosisGPT ($\tau = 0.6$)	34.7	44.2 $+9.5$	3.3

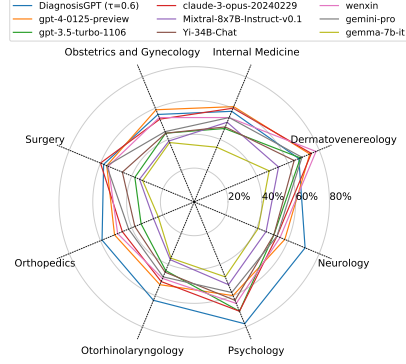


Figure 4: **Left:** Results of DxBench. **Acc.** indicates the accuracy of diagnosis. n denotes the average number of inquiry rounds, with limitation of $L = 5$. **Right:** Diagnostic accuracy in top eight departments of DxBench with DiagnosisGPT set to $\tau = 0.6$.

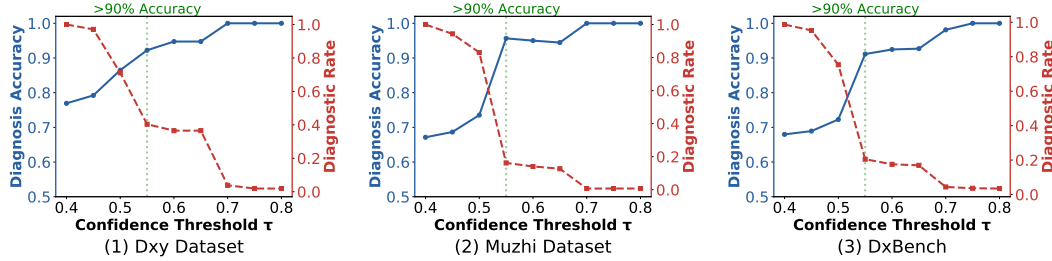


Figure 5: Relationship between diagnosis confidence and accuracy. We provided all symptoms ($\mathcal{S}_{\text{exp}} \cup \mathcal{S}_{\text{imp}}$) to DiagnosisGPT for direct disease diagnosis (without symptom inquiry). **Diagnosis Accuracy** represents the accuracy of diagnoses exceeding the threshold τ . **Diagnostic Rate** indicates the proportion of data that exceed τ , i.e., the proportion of cases where the model makes a diagnosis.

Table 4: The impact of τ on DiagnosisGPT. Experiments were conducted on the DxBench with $L = 5$ setting. a represents diagnostic accuracy. n denotes the number of queries.

τ	0	0.4	0.5	0.6	0.7	0.8
$a \uparrow$	56.9	61.3 $+4.1$	62.4 $+5.5$	64.5 $+7.6$	65.3 $+8.4$	64.9 $+8.0$
$n \downarrow$	0	0.7	1.8	3.1	4.2	4.7

Table 5: The effect of conversation rounds on entropy for DiagnosisGPT ($\tau = 0.6$) on DxBench.

Number of Rounds b	1	2	3	4	5
Average Entropy H	1.467	1.396	1.380	1.371	1.369

We further tested the performance of Diagnosis on DxBench under various τ values. As shown in Table 4, as τ increases, the number of iterations n rises, and the accuracy a also improves. It can be seen that τ serves as an effective control for balancing the trade-off in diagnosis. Table 5 illustrates the process of entropy reduction in Diagnosis as the number of inquiries increases.

5.2 Ablation Study

Table 6 presents the ablation study results for CoD. The performance of DiagnosisGPT_{baseline} indicates that directly training the target disease prediction on CoD training data does not significantly

Table 6: Ablation results for DiagnosisGPT. All ablation models are retrained using the CoD training dataset. *w/o Confidence for Decision* signifies that the model directly generates the subsequent decision, akin to other LLMs. DiagnosisGPT_*baseline* denotes the models that directly learn disease prediction without symptom inquiry from the CoD training data.

Model	DxBench		Muzhi Dataset		Dxy Dataset	
	Acc.	<i>n</i>	Acc.	<i>n</i>	Acc.	<i>n</i>
DiagnosisGPT ($\tau = 0.5$)	62.4	1.8	64.1	1.4	72.6	1.5
w/o Knowledge Integration	61.5	1.7	64.8	1.4	72.1	1.5
w/o Diagnostic Reasoning	60.5	1.4	61.2	1.4	68.2	1.4
w/o Confidence for Decision	59.2	0.8	59.1	0.6	65.3	0.7
DiagnosisGPT_ <i>baseline</i>	55.2	0.0	58.4	0.0	61.5	0.0

enhance diagnosis. The results show that diagnostic reasoning benefits the model’s performance, implying that providing diagnostic analysis is effective. Without diagnostic confidence, the model’s performance drops significantly, suggesting that model decisions based on confidence are valuable. Furthermore, knowledge integration does not appear to have a substantial effect, probably due to that the model grasps disease knowledge through CoD training data. We leave this for future research.

5.3 Case Study

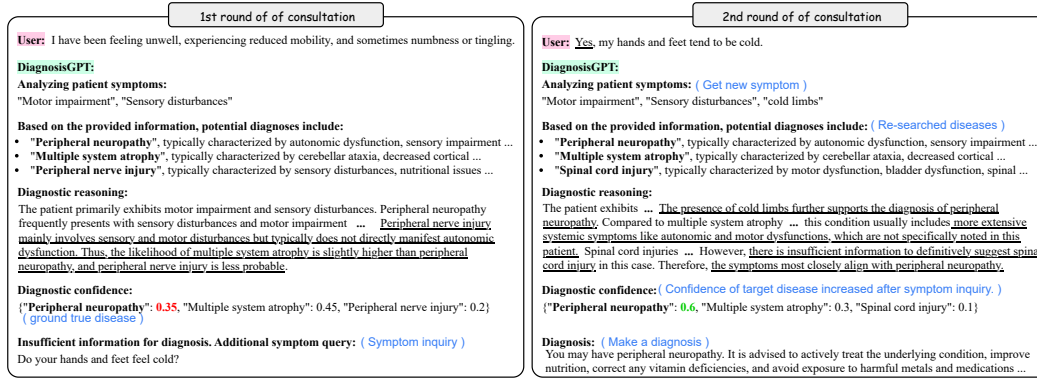


Figure 6: A diagnosis case by DiagnosisGPT, which made the correct diagnosis by inquiring symptoms. The underline indicates the rationale for the disease confidence.

Figure 6 presents a diagnostic case using DiagnosisGPT. When queried about diseases, DiagnosisGPT systematically outputs its diagnostic reasoning process. It first summarizes the user’s symptom information, then identifies potential diseases, and begins the diagnostic analysis, ultimately providing a diagnostic confidence level. As shown in the first round of replies, the highest confidence level is 0.45, below the threshold, prompting the model to inquire about symptoms. When the patient responds to the symptom inquiry, the probability of the target disease significantly increases, leading DiagnosisGPT to confirm and makes a correct diagnosis.

6 Conclusion

In this paper, we propose the **Chain of Diagnosis (CoD)** to enhance the interpretability of large language models (LLMs) for disease diagnosis. Using CoD, we developed DiagnosisGPT, an LLM that supports the diagnosis of 9,604 diseases. Distinct from other LLMs, DiagnosisGPT can provide diagnostic confidence and relies on its own disease database for diagnostic reasoning. Experiments show that the diagnostic capabilities of DiagnosisGPT surpass those of other LLMs. Furthermore, higher accuracy can be achieved by adjusting the diagnostic threshold values. This means that CoD can control the trade-off between effectiveness and efficiency in diagnosis. CoD offers a novel solution for medical diagnosis. We believe that the data, models, and methods from this work can advance the field of medical LLMs.

Acknowledgement

This work was supported by the Shenzhen Science and Technology Program (JCYJ20220818103001002), Shenzhen Doctoral Startup Funding (RCBS20221008093330065), Tianyuan Fund for Mathematics of National Natural Science Foundation of China (NSFC) (12326608), Shenzhen Key Laboratory of Cross-Modal Cognitive Computing (grant number ZDSYS20230626091302006), and Shenzhen Stability Science Program 2023, Shenzhen Key Lab of Multi-Modal Cognitive Computing.

References

- [1] Yuan Xia, Jingbo Zhou, Zhenhui Shi, Chao Lu, and Haifeng Huang. Generative adversarial regularized mutual information policy gradient framework for automatic diagnosis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 1062–1069, 2020.
- [2] Zhenyu Hou, Yukuo Cen, Ziding Liu, Dongxue Wu, Baoyan Wang, Xuanhe Li, Lei Hong, and Jie Tang. Mtdiag: an effective multi-task framework for automatic diagnosis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 14241–14248, 2023.
- [3] Chaitanya Shivade, Preethi Raghavan, Eric Fosler-Lussier, Peter J Embi, Noemie Elhadad, Stephen B Johnson, and Albert M Lai. A review of approaches to identifying patient phenotype cohorts using electronic health records. *Journal of the American Medical Informatics Association*, 21(2):221–230, 2014.
- [4] Lin Xu, Qixian Zhou, Ke Gong, Xiaodan Liang, Jianheng Tang, and Liang Lin. End-to-end knowledge-routed relational dialogue system for automatic diagnosis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7346–7353, 2019.
- [5] Arsene Fansi Tchango, Rishab Goel, Zhi Wen, Julien Martel, and Joumana Ghosn. Ddxplus: A new dataset for automatic medical diagnosis. *Advances in Neural Information Processing Systems*, 35:31306–31318, 2022.
- [6] Junying Chen, Xidong Wang, Anningzhe Gao, Feng Jiang, Shunian Chen, Hongbo Zhang, Dingjie Song, Wenya Xie, Chuyi Kong, Jianquan Li, et al. Huatuogpt-ii, one-stage training for medical adaption of llms. *arXiv preprint arXiv:2311.09774*, 2023.
- [7] Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath Rangan, and Jonathan H Chen. Diagnostic reasoning prompts reveal the potential for large language model interpretability in medicine. *NPJ Digital Medicine*, 7(1):20, 2024.
- [8] Chandan Singh, Jeevana Priya Inala, Michel Galley, Rich Caruana, and Jianfeng Gao. Rethinking interpretability in the era of large language models, 2024.
- [9] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions, 2023.
- [10] Chandan Singh, Armin Askari, Rich Caruana, and Jianfeng Gao. Augmenting interpretable models with large language models during training. *Nature Communications*, 14(1):7913, 2023.
- [11] Zhiyuan Hu, Chumin Liu, Xidong Feng, Yilun Zhao, See-Kiong Ng, Anh Tuan Luu, Junxian He, Pang Wei Koh, and Bryan Hooi. Uncertainty of thoughts: Uncertainty-aware planning enhances information seeking in large language models, 2024.
- [12] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.
- [13] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- [14] Junying Chen, Qingcai Chen, Dongfang Li, and Yutao Huang. Sedr: Segment representation learning for long documents dense retrieval. *arXiv preprint arXiv:2211.10841*, 2022.

- [15] Zhongyu Wei, Qianlong Liu, Baolin Peng, Huaixiao Tou, Ting Chen, Xuan-Jing Huang, Kam-Fai Wong, and Xiang Dai. Task-oriented dialogue system for automatic diagnosis. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 201–207, 2018.
- [16] 01. AI. :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024.
- [17] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *CoRR*, abs/1908.10084, 2019.
- [18] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. Optimizing dense retrieval model training with hard negatives. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1503–1512, 2021.
- [19] Kangenbei Liao, CHENG ZHONG, Wei Chen, Qianlong Liu, Baolin Peng, Xuanjing Huang, et al. Task-oriented dialogue system for automatic disease diagnosis via hierarchical reinforcement learning. 2021.
- [20] Junying Chen, Dongfang Li, Qingcai Chen, Wenxiu Zhou, and Xin Liu. Diaformer: Automatic diagnosis via symptoms sequence generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 4432–4440, 2022.
- [21] Machel Reid, Nikolay Savinov, Denis Teplyashin, and et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024.
- [22] Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, Yanbin Zhao, Yuxiang Lu, et al. Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation. *arXiv preprint arXiv:2107.02137*, 2021.
- [23] anthropic. Introducing the next generation of claude. <https://www.anthropic.com/news/claude-3-family>, 2024.
- [24] OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt>, 2022.
- [25] OpenAI. Gpt-4 technical report, 2023.
- [26] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- [27] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mixtral of experts, 2024.
- [28] Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, et al. Meddialog: Large-scale medical dialogue datasets. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9241–9250, 2020.
- [29] Honglin Xiong, Sheng Wang, Yitao Zhu, Zihao Zhao, Yuxiao Liu, Qian Wang, and Dinggang Shen. Doctorglm: Fine-tuning your chinese doctor is not a herculean task. *arXiv preprint arXiv:2304.01097*, 2023.
- [30] Ming Xu. Medicalgpt: Training medical gpt model. <https://github.com/shibing624/MedicalGPT>, 2023.

- [31] Wenya Xie, Qingying Xiao, Yu Zheng, Xidong Wang, Junying Chen, Ke Ji, Anningzhe Gao, Xiang Wan, Feng Jiang, and Benyou Wang. Llms for doctors: Leveraging medical llms to assist doctors, not replace them. *arXiv preprint arXiv:2406.18034*, 2024.
- [32] Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Jianquan Li, Guiming Chen, Xiangbo Wu, Zhiyi Zhang, Qingying Xiao, et al. Huatuoogpt, towards taming language model to be a doctor. *arXiv preprint arXiv:2305.15075*, 2023.
- [33] Junying Chen, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, Guangjun Yu, et al. Huatuoogpt-vision, towards injecting medical visual knowledge into multimodal llms at scale. *arXiv preprint arXiv:2406.19280*, 2024.
- [34] Xidong Wang, Nuo Chen, Junyin Chen, Yan Hu, Yidong Wang, Xiangbo Wu, Anningzhe Gao, Xiang Wan, Haizhou Li, and Benyou Wang. Apollo: Lightweight multilingual medical llms towards democratizing medical ai to 6b people. *arXiv preprint arXiv:2403.03640*, 2024.
- [35] Dongfang Li, Baotian Hu, Qingcai Chen, Weihua Peng, and Anqi Wang. Towards medical machine reading comprehension with structural knowledge and plain text. In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, pages 1427–1438, 2020.
- [36] Ke Ji, Junying Chen, Anningzhe Gao, Wenya Xie, Xiang Wan, and Benyou Wang. Llms could autonomously learn without external supervision. *arXiv preprint arXiv:2406.00606*, 2024.
- [37] Kai-Fu Tang, Hao-Cheng Kao, Chun-Nan Chou, and Edward Y Chang. Inquire and diagnose: Neural symptom checking ensemble using deep reinforcement learning. In *NIPS workshop on deep reinforcement learning*, 2016.
- [38] Lin Xu, Qixian Zhou, Ke Gong, Xiaodan Liang, Jianheng Tang, and Liang Lin. End-to-end knowledge-routed relational dialogue system for automatic diagnosis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7346–7353, 2019.
- [39] Kangerbei Liao, Qianlong Liu, Zhongyu Wei, Baolin Peng, Qin Chen, Weijian Sun, and Xuanjing Huang. Task-oriented dialogue system for automatic disease diagnosis via hierarchical reinforcement learning. *arXiv preprint arXiv:2004.14254*, 2020.
- [40] Zhengxu Hou, Bang Liu, Ruihui Zhao, Zijing Ou, Yafei Liu, Xi Chen, and Yefeng Zheng. Imperfect also deserves reward: Multi-level and sequential reward modeling for better dialog management. In *NAACL*, 2021.
- [41] Milene Santos Teixeira, Vinícius Maran, and M. Dragoni. The interplay of a conversational ontology and ai planning for health dialogue management. *Proceedings of the 36th Annual ACM Symposium on Applied Computing*, 2021.
- [42] Zhihong Chen, Feng Jiang, Junying Chen, Tiannan Wang, Fei Yu, Guiming Chen, Hongbo Zhang, Juhao Liang, Chen Zhang, Zhiyi Zhang, et al. Phoenix: Democratizing chatgpt across languages. *arXiv preprint arXiv:2304.10453*, 2023.
- [43] Fei Yu, Anningzhe Gao, and Benyou Wang. Ovm, outcome-supervised value models for planning in mathematical reasoning, 2024.
- [44] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [45] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [46] Maciej Besta, Nils Blach, Ales Kubicek, and Gerstenberger. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- [47] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023.

A Related Work

LLMs for Medical Scenarios The success of models like ChatGPT [24] has inspired research into their application in healthcare, resulting in medical-specific LLMs such as DoctorGLM [29], MedicalGPT [30], DotaGPT [31], HuatuoGPT [32, 6, 33], and Apollo [34]. Despite their focus on medical knowledge, these models have limited capabilities in automating medical diagnoses.

Automated Diagnosis Task Medical diagnosis, a key AI application in healthcare [20, 35, 36], has predominantly utilized reinforcement learning (RL). Pioneering works include [37], who introduced neural symptom checking using RL. Subsequent advancements include hierarchical RL for diagnostic and contextual decisions [19], Deep Q-networks for symptom collection from patient interactions [15], and incorporation of medical knowledge into RL policy learning [38]. Two-level hierarchical RL [39], policy gradient frameworks with Generative Adversarial Networks [1], and customization of RL models using multi-level rewards and dialogue data [40, 41] have further enhanced diagnostic accuracy. [20] and [2] conceptualizes automatic diagnosis as a sequence generation task. However, these models are limited by predefined symptoms and diseases, and cannot support open-ended consultations.

Reasoning of LLMs LLMs show promise in complex tasks such as mathematical reasoning [42, 43]. To harness their reasoning abilities, CoT[44] is proposed with intermediate steps, and Tree-of-Thought (ToT)[45] using DFS/BFS for enhanced reasoning paths. Graph of Thoughts (GoT) [46] is introduced for intricate problems. ReAct [47] combines reasoning with actions. Uncertainty of Thoughts (UoT) [11] improves decision-making by simulating multiple requests for information gain.

B DxBench Distribution

The data distribution in DxBench dataset is illustrated in Figure 7. We categorize the data distribution according to the medical departments responsible for diagnosing the diseases. The data shows a relatively balanced distribution across different departments. Notably, the Dermatovenereology department has the highest number of entries with 121 cases, while the Infectious Diseases and Immunology department has the fewest, with 27 cases.

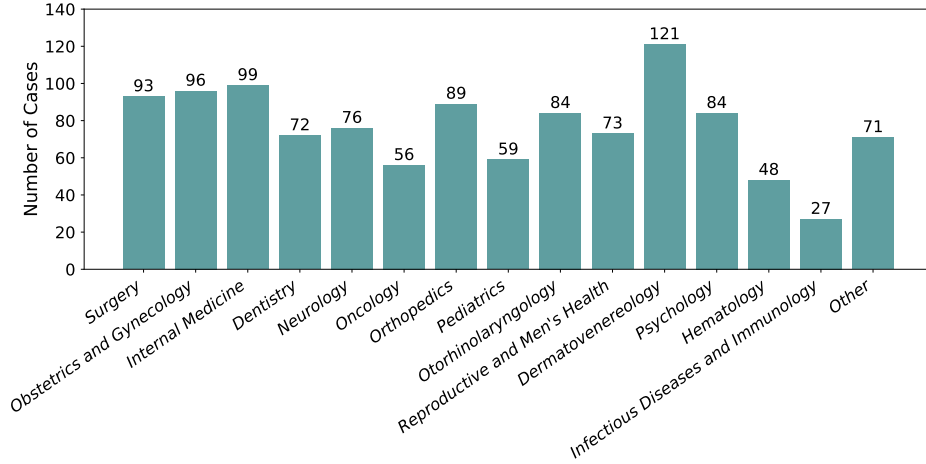


Figure 7: Data distribution across different departments in the DxBench dataset.

C The prompt for LLM Diagnosis

The prompt for LLM diagnosis is shown in Table 8. We instruct the LLMs to determine whether a diagnosis can be made. If a diagnosis is possible, the LLMs output the diagnosed disease. Otherwise, the LLMs query the user with questions regarding a specific symptom.

The prompt for LLM Diagnosis

Initial input:

You are a professional physician tasked with diagnosing a patient based on their symptom information. I will provide you with information on possible diseases, and you will need to carefully consider which of the candidate diseases the patient might have.

Patient symptom information is marked by <symptoms>, and candidate diseases by <candidate_diseases>.

<symptoms> {Known_symptoms} <symptoms>

<candidate_diseases> {candidate_diseases} <candidate_diseases>

If you believe a diagnosis can be made, select the most likely disease from <candidate_diseases> (choose only one). Example output:

```
{"judge": true, "disease": "common cold"}
```

If you believe the information on symptoms is insufficient, ask the patient for more symptom information, noting that you can only ask about one symptom. Example output:

```
{"judge": false, "symptom": "Do you have a lack of appetite?"}
```

Please output in JSON format.

Input after patient response:

The patient's response will be marked by <Patient>. The hints I give you are marked by <Hit>.

<Patient> {patient_response} <Patient>

<Hit>Please, based on the patient's response, decide now whether a diagnosis can be made.

If you believe a diagnosis can be made, select the most likely disease from <candidate_diseases> (choose only one). Example output:

```
{"judge": true, "disease": "common cold"}
```

If you believe the information on symptoms is insufficient, ask the patient for more symptom information, noting that you can only ask about one symptom. Example output:

```
{"judge": false, "symptom": "Do you have a lack of appetite?"}
```

Please output in JSON format.<Hit>

Figure 8: The prompt for LLM Diagnosis. `Known_symptoms` represents the symptoms currently known by the LLM. `candidate_diseases` represents the list of candidate diseases D . `candidate_diseases` represents the response of the patient.

D Patient Simulator for Evaluation

To evaluate the automatic diagnostic capabilities of LLMs, we instruct GPT-4 to play the role of a patient. Initially, we provide explicit symptoms S_{exp} as input for the model to diagnose. If the LLMs ask questions, the patient GPT will respond using a simulated patient prompt, as shown in Figure 9.

Simulated patient prompt

You are a patient, here are your symptom details:
{Symptoms}

Your actual disease is {disease}.

You need to answer the doctor's question:
{LLM_query}

Please answer the doctor's question based on your symptom information and disease, simply reply with "yes" or "no", and do not include any other content.

Figure 9: Simulated patient prompt for responding to questions posed by LLMs. The Symptoms represents all the symptom information of the case, $S_{exp} \cup S_{imp}$. The disease indicates the true disease of the case, d_t .

E Prompt of Data Synthesis

We constructed a disease database encompassing 9,604 diseases. Each disease entry includes four fields: "disease name", "overview", "symptoms", and "treatment". For each disease, we used the prompt shown in Figure 10 to generate five patient cases with GPT-4, ensuring that each case study exhibits distinct typical characteristics.

A Knowledge-Base-Driven Approach to Medical Case Generation

Disease: {disease name}

Overview of this disease: {overview}

Common symptoms of this disease include: {symptoms}

Please complete the following tasks based on the description above.

1. First, generate basic demographic information about the population affected by this disease: gender and age.
2. You need to construct five real cases concerning this disease. In these five cases, two should have only one main symptom, two should have two main symptoms, and one should have more than three main symptoms (main symptoms are the most noticeable ones). Each case should include 2-4 implicit symptoms (generally, symptoms that can be elicited by a doctor's questioning). Ensure each case is a typical example of this disease.

Output in JSON format, and only output the JSON content, do not output anything else. The example output is:

```
{
  "Basic Information": {
    "Gender": "Female",
    "Age": "Child"
  },
  "Case 1": {
    "Main Symptoms": ["Symptom 1", "Symptom 2"],
    "Implicit Symptoms": ["Symptom 3", "Symptom 4", "Symptom 5"]
  },
  "Case 2": "..."
}
```

Figure 10: The prompt of synthesizing patient cases. {disease name}, {overview}, and {symptom} represent the corresponding information for diseases in the database.

F Prompt of CoD

To generate CoD training data, we prompt GPT-4 to construct CoD dialogue data based on patient case data. This involves the following 8 prompts:

Prompt 1: Patient Self-report Prompt (Role: Patient) As shown in figure 11, the patient self-report prompt is used to generate the user’s initial question q_1 based on the patient’s explicit symptoms, primarily expressing the patient’s chief complaint.

Patient Self-report Prompt	Patient
<p>System Prompt: I’d like you to pretend to be a patient and describe your condition to the doctor in the voice of a patient. Please avoid using overly technical terms. Questions from the doctor will be marked with <Doctor>. The response of yourself will be marked with <Patient>. The hints I provide will be marked with <Hint>.</p> <p>Query: <Doctor> Hello, I’m a doctor. How can I help you? <Doctor> <Hint> Your symptoms are: {explicit_syms} If the symptoms include information about the patient’s age and gender, such as elderly, female, etc., please inform the doctor. Please reply in the patient’s voice, only output the patient’s words and nothing else. <Hint></p>	

Figure 11: Patient Self-report Prompt. {explicit_syms} indicates S_{exp} .

Prompt 2: Reasoning Prompt (Role: Diagnosis) When provided with the known symptoms S of a patient and the candidate diseases D' , the reasoning prompt, as illustrated in Figure 12, is utilized to generate the reasoning process T and the confidence distribution C .

Reasoning Prompt	Diagnosis
<p>You are a professional physician tasked with diagnosing a patient based on provided symptom information. You will be given a list of candidate diseases, and your role is to offer a detailed diagnostic analysis and a confidence distribution of the candidate diseases for the patient.</p> <p>You need to first analyze the patient’s condition and think about which of the candidate diseases the patient might have. Then, output the diagnostic confidence distribution of the candidate diseases in JSON format, please output a dict rather than a list.</p> <p>An output example is: {"analysis":..., "distribution": {"Animal skin disease": 0.25, "Erythema ab igne": 0.2, "Dermatitis": 0.55}}</p> <p>The patient’s explicit symptoms: {explicit_syms}, The patient’s implicit symptoms: {implicit_syms}, Candidate diseases: {candidate_diseases}</p> <p>Please first analyze the patient’s condition, then output the probability distribution of these diseases.</p>	

Figure 12: Reasoning Prompt. {explicit_syms} denotes S_{exp} . {implicit_syms} denotes the inquired symptoms $S \setminus S_{exp}$. {candidate_diseases} denotes the currently identified diseases D' .

Prompt 3: Rethinking Prompt (Role: Diagnosis) If the generated C does not meet the condition $\max C \setminus c_{d_t} > \tau$, the rethinking prompt, as shown in Figure 13, is used to have GPT4 regenerate a valid diagnosis T and C .

Rethinking Prompt	Diagnosis
<p>Your diagnostic analysis did not pass inspection because you assigned an high confidence level to a potentially incorrect disease. Please reconsider your assessment and provide a new diagnostic analysis along with the confidence distribution. Ensure that the output format remains exactly the same.</p> <p>An output example is: {"analysis":..., "distribution": {"Animal skin disease": 0.25, "Erythema ab igne": 0.2, "Dermatitis": 0.55}}</p>	

Figure 13: Rethinking Prompt.

Prompt 4: Doctor Diagnosis Prompt (Role: Doctor) If $\max C > \tau$, we prompt GPT-4 to generate a response regarding the diagnostic result. The prompt used is shown in Figure 14. The disease database information will be provided to generate more reliable suggestions. Once the diagnostic response is generated, the data generation process concludes.

Doctor Diagnosis Prompt	Doctor
<p>System Prompt: Please play the role of a doctor to ask the patient about their condition or diagnose the disease. The patient's responses will be marked with <Patient>. The response of yourself will be marked with <Doctor>. The hints I provide will be marked with <Hint>.</p> <p>Query: {Chat_history} <Hint>The doctor's diagnosis is marked by <diagnosis of disease>. The treatment method for this diagnosed disease is marked by <treatment method>. <diagnosis of disease> {disease_name} <diagnosis of disease> <treatment method> {treatment} <treatment method> Based on the information above and the historical conversation records, please diagnose the patient and provide detailed recommendations. Reply in the tone of a doctor, and do not start with the word 'doctor'.<Hint></p>	

Figure 14: Doctor Diagnosis Prompt. {Chat_history} represents the previously generated conversation history. {disease_name} represents d_t . {treatment} represents the "treatment" of d_t from the database.

Prompt 5: Symptom Generation Prompt (Doctor) If $\max C \leq \tau$, we will have the LLM generate the symptom s_{gen} it wants to inquire about, using the prompt shown in Figure 15. Then, we will select the inquired symptom s_t from $S_{imp} \cup \{s_{gen}\}$ based on $H(C|s)$.

Symptom Generation Prompt	Doctor
<p>You are now a professional physician, and you need to infer the next symptom to ask the patient based on the following information.</p> <p>The patient's explicit symptoms: {explicit_syms}</p> <p>The patient's implicit symptoms: {implicit_syms}</p> <p>The patient may currently have {predicted_disease}.</p> <p>Please infer the next symptom to ask the patient, asking only one symptom that has not been previously inquired.</p> <p>The output format should be json, for example: {"symptom": "headache"}</p>	

Figure 15: Symptom Generation Prompt. {explicit_syms} denotes S_{exp} . {implicit_syms} denotes the inquired symptoms $S \setminus S_{exp}$. {predicted_disease} represents the currently most likely disease $\arg\max_{d \in \mathcal{D}'} \{c_d\}$.

Prompt 6: Doctor Inquiry Prompt (Role: Doctor) After confirming the symptom s_t , the Doctor Inquiry Prompt, shown in Figure 16, generates questions regarding the symptom.

Doctor Inquiry Prompt	Doctor
<p>System Prompt: Please play the role of a doctor to ask the patient about their condition or diagnose the disease. The patient's responses will be marked with <Patient>. The response of yourself will be marked with <Doctor>. The hints I provide will be marked with <Hint>.</p> <p>Query: {Chat_history} <Hint>Please inquire about the patient's condition based on these symptoms: {current_sym}</p> <p>Please reply in the tone of a doctor, asking only one question in a conversational manner that the patient can understand. Do not start with the word 'doctor'.<Hint></p>	

Figure 16: Doctor Inquiry Prompt. {Chat_history} represents the previously generated conversation history. {current_sym} represents s_t .

Prompt 7: Symptom Assessment Prompt (Role: Patient) As shown in Figure 17, the symptom assessment prompt is used to determine whether the patient exhibits the symptoms inquired about by the doctor.

Symptom Assessment Prompt	Role: Patient
<p>You are now a professional physician. Please judge whether the patient has the symptom based on the patient's information.</p> <p>It is known the patient's main symptoms are {explicit_syms}, and the implicit symptoms are {implicit_syms}.</p> <p>Please determine whether the patient has {choose_sym}.</p> <p>Search in the patient's existing symptoms, paying attention to synonyms. If found, output true; if not found, output false.</p> <p>The output format is json, for example: {"headache": true}</p>	

Figure 17: Symptom Assessment Prompt. {explicit_syms} and {implicit_syms} represent S_{exp} and S_{imp} respectively. {choose_sym} represents s_t , the symptom being inquired about.

Prompt 8: Patient Response Prompt (Role: Patient) The Patient Response Prompt, as shown in Figure 18, is used to generate verbal responses from patients regarding the symptom of inquiry. Then, s_t will be added to the doctor's known symptoms S , initiating the next doctor's response generation.

Patient Response Prompt	Patient
<p>System Prompt: I'd like you to pretend to be a patient and describe your condition to the doctor in the voice of a patient. Please avoid using overly technical terms. Questions from the doctor will be marked with <Doctor>. The response of yourself will be marked with <Patient>. The hints I provide will be marked with <Hint>.</p> <p>Query: {Chat_history} <Hint>Please answer the doctor's questions based on the information, note that you {do_or_do_not} have this symptom: {choose_sym} Please respond in the tone of a patient, and avoid using overly technical terms. Do not begin with the word 'patient'.<Hint></p>	

Figure 18: Patient Response Prompt. {Chat_history} represents the previously generated conversation history. {current_sym} represents s_t . {do_or_do_not} indicates whether the patient has the symptom determined by the Symptom Assessment Prompt.

G Building of DxBench

We used MedDialog [28], a publicly available dataset with real doctor-patient dialogues from various departments, to create a realistic evaluation dataset. Despite its chaotic nature, we filtered 3,121 cases with clear dialogue and definite disease diagnoses. Then, using prompts designed to extract structured case information, as shown in Figure 19, we extracted structured case information from these cases. We further filtered out cases where the diseases were not included in the database or where there was insufficient symptom information, resulting in 1,429 cases. Finally, we manually filtered 1,148 high-quality cases, verifying the accuracy and completeness of symptom extraction, to construct DxBench. Additionally, for each case, we identified the two most relevant diseases using

the disease retriever of DiagnosisGPT to form three candidate diseases. The construction process is detailed in Appendix G.

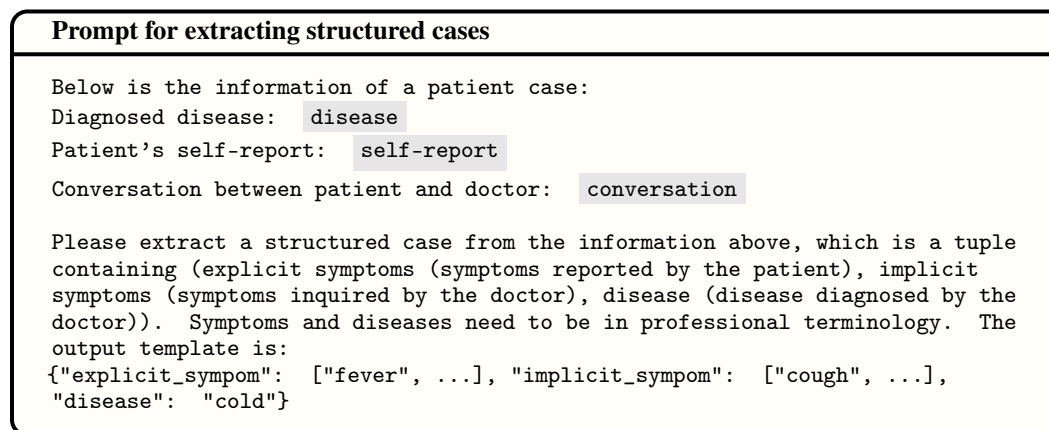


Figure 19: Prompt for extracting structured cases. `{disease}` indicates the doctor's diagnosed disease. `self-report` indicates the patient's self-reported symptoms. `conversation` represents the content of the dialogue between the doctor and the patient.

H Case Study

Below, two detailed diagnostic cases will be provided to better understand the response of our model.

User

I recently found a lump in my abdomen, and the pain is severe. Can you help me?

DiagnosisGPT

Analyzing patient symptoms:

"abdominal lump", "severe pain"

Based on the information provided, the likely diagnoses include:

- **Greater omentum cyst**, common symptoms include an abdominal mass or progressive abdominal enlargement, a feeling of heaviness in the abdomen when lying on the back, and severe abdominal pain may occur when there is concurrent intestinal torsion or obstruction.
- **Intestinal tumor**, common symptoms include an abdominal lump, severe pain, anemia, rectal bleeding, frequent and malformed stools, and fatigue with weight loss.
- **Recurrent retroperitoneal tumor**, common symptoms include an abdominal lump, abdominal distension, abdominal pain, dull pain, and severe pain.

Diagnostic reasoning:

The patient presents symptoms of abdominal mass, severe pain, and rectal bleeding. These symptoms may be associated with various intra-abdominal conditions. The abdominal mass and severe pain may be direct manifestations of some abdominal tumor-related disease, while rectal bleeding could indicate involvement of the in- testines. When considering candidate diseases, the likelihood of colorectal tumors is high as they can account for the occurrence of all these symptoms. Mesenteric cysts typically present as abdominal masses but may not necessarily cause significant pain or rectal bleeding if they are small or located in specific positions. Recurrent retroperitoneal tumors, depending on their origin and type, could also result in sim- ilar symptoms, especially if they affect the intestines or other abdominal structures. Considering all factors, colorectal tumors are more closely associated with symptoms such as rectal bleeding due to their direct involvement with the intestines. Recur- rent retroperitoneal tumors are less likely than colorectal tumors but more likely than mesenteric cysts, as the latter typically show a weaker association with pain and rectal bleeding.

Diagnostic confidence:

- "Greater Omental cyst": 0.15,

- "Intestinal tumor": 0.4,
- "Recurrent retroperitoneal tumor": 0.45

Insufficient information for diagnosis. Additional symptom query:

Have you noticed any abnormalities or blood streaks in your stool color recently?

User

I've noticed some bloody stools recently, and the color doesn't seem right.

DiagnosisGPT

Analyzing patient symptoms:

"abdominal lump", "severe pain", "bloody stool"

Based on the information provided, the likely diagnoses include:

- **Intestinal tumor**, ...
- **Greater omentum cyst**, ...
- **Recurrent retroperitoneal tumor**, ...

(The disease retriever identified the same three diseases as last time, but their priority has changed.)

Diagnostic reasoning:

The patient presents symptoms of abdominal mass, severe pain, and rectal bleeding. These symptoms may be associated with various intra-abdominal conditions. The abdominal mass and severe pain may be direct manifestations of some abdominal tumor-related disease, while rectal bleeding could indicate involvement of the intestines. When considering candidate diseases, the likelihood of colorectal tumors is high as they can account for the occurrence of all these symptoms. Mesenteric cysts typically present as abdominal masses but may not necessarily cause significant pain or rectal bleeding if they are small or located in specific positions. Recurrent retroperitoneal tumors, depending on their origin and type, could also result in similar symptoms, especially if they affect the intestines or other abdominal structures. Considering all factors, colorectal tumors are more closely associated with symptoms such as rectal bleeding due to their direct involvement with the intestines. Recurrent retroperitoneal tumors are less likely than colorectal tumors but more likely than mesenteric cysts, as the latter typically show a weaker association with pain and rectal bleeding.

Diagnostic confidence:

- "Intestinal tumor": 0.6,
- "Greater Omental cyst": 0.15,
- "Recurrent retroperitoneal tumor": 0.25

Diagnosis

Based on the symptoms and rectal bleeding you described, there is a possibility of a **Intestinal tumor**. I recommend you to visit a hospital for further examinations, including endoscopy and imaging tests, to confirm the diagnosis. Early-stage Intestinal tumors can be treated with surgical removal, while malignant tumors may require additional treatments such as radiation therapy or chemotherapy. Wishing you a speedy recovery.

Case 1: Interaction between an intestinal tumor patient and the DiagnosisGPT model.

User

I recently discovered that I have developed an inflammatory lump, and I have also been running a fever, feeling generally unwell.

DiagnosisGPT

Analyzing patient symptoms:

"inflammatory mass", "fever"

Based on the information provided, the likely diagnoses include:

- **Granulomatous lymphadenitis**, common symptoms include inflammatory mass, lymph node enlargement and pain, fever.
- **Geotrichosis**, common symptoms include fever, chronic cough, sputum production and hemoptysis, stomach pain, diarrhea, cheesy white pseudomembrane on the oral mucosa, tongue, and palate, erythematous base with a burning sensation, involvement of the pharynx, larynx, and their junctions, skin erythema and exudation,

cutaneous candidiasis in skin folds, tumor-like nodules, granulomas, itching and pain, edema, erythema, epithelial desquamation, and crust formation.

- **Suppurative mastitis**, common symptoms include fever, localized breast pain, lump, breast fistula with abscess formation that may rupture internally or externally.

Diagnostic reasoning:

The primary symptoms of the patient are inflammatory masses and fever. Without more detailed symptom descriptions, lymphadenitis granulomatous, geotrichosis, and suppurative mastitis can all present with these symptoms. Lymphadenitis granulomatous often causes lymph node enlargement and may be accompanied by inflammatory masses. Geotrichosis, an infection caused by *Geotrichum*, can also cause systemic symptoms such as fever and may lead to granuloma formation in immunocompromised individuals. Suppurative mastitis primarily occurs in breastfeeding women, with symptoms including pain in the breast area, masses, and fever. Given the presence of inflammatory masses, if the patient is a breastfeeding woman, the likelihood of suppurative mastitis increases. Due to the lack of explicit latent symptoms, negating symptoms information, and detailed background information of the patient (such as gender, age, occupation, etc.), the specific diagnosis requires further clinical examination and auxiliary examination data.

Diagnostic confidence:

- **Granulomatous lymphadenitis**: 0.3,
- **Geotrichosis**: 0.3,
- **Suppurative mastitis**: 0.4

Insufficient information for diagnosis. Additional symptom query:

Do you have a headache?

User

I have a bit of a headache, but it's not very serious.

DiagnosisGPT

Analyzing patient symptoms:

"inflammatory mass", "fever", "headache"

Based on the information provided, the likely diagnoses include:

- **Granulomatous lymphadenitis**, common symptoms include inflammatory mass, lymph node enlargement and pain, fever.
- **Occipital indurative folliculitis**, common symptoms include fever, hard rashes, swollen hair follicles, pain, and poor sleep.
- **Hair cysts**, common symptoms include cystic masses on the scalp, often solitary, more common in women.

Diagnostic reasoning:

The patient's primary symptoms are inflammatory masses accompanied by fever, which typically indicate an inflammatory response in a certain area of the body. Headaches may be a systemic symptom caused by the body's overall inflammatory response. Common symptoms of lymphadenitis granulomatosa include fever and localized lymph node enlargement, and the local inflammatory masses may be related to the swelling of the lymph nodes. Suboccipital kerion celsi is characterized by inflammation and induration of hair follicles in a specific area of the head and typically does not cause systemic symptoms such as fever. Pilar cysts are cystic formations within hair follicles, usually accompanied by localized inflammation, but are unlikely to cause fever and systemic symptoms. Considering these candidate conditions and the patient's symptoms, a preliminary judgment can be made accordingly.

Diagnostic confidence:

- **Granulomatous lymphadenitis**: 0.6,
- **Occipital indurative folliculitis**: 0.1,
- **Hair cysts**: 0.3

Diagnosis

It seems you may have **granulomatous inflammation**. I recommend that you see a doctor as soon as possible, as you may need some blood tests and imaging exams to confirm the diagnosis. The doctor will develop a detailed treatment plan based on the examination results. Meanwhile, make sure to rest, drink plenty of water, and avoid overexertion.

Case 2: Interaction between an granulomatous inflammation patient and the DiagnosisGPT model.

I Performance of Disease Retriever

We allocated 10% of the data as a validation set to evaluate retrieval performance. Table 7 shows the retrieval performance of diseases under Diagnosis on the validation set. It can be seen that the top 3 diseases achieve a recall rate of 73%, indicating that most diseases can be effectively excluded.

Metric	Value
MRR@100	0.6926
Recall@3	0.7357
Recall@5	0.7807
Recall@10	0.8430
Recall@30	0.9047
Recall@50	0.9313
Recall@100	0.9573

Table 7: Performance of the Disease Retriever on the Validation Set

J Standard Errors of Result

We report the standard errors of the results from our model in Table 8. The standard errors were obtained by conducting five random experiments.

Model	DxBench	Muzhi Dataset	Dxy Dataset
GPT4-0125-preview	62.5 \pm 2.1	56.3 \pm 2.4	65.4 \pm 0.5
DiagnosisGPT ($\tau = 0.6$)	64.5 \pm 1.8	65.5 \pm 2.1	75.4 \pm 0.3

Table 8: Statistical Results of Standard Errors.

K Review of Synthetic Cases by Medical Experts

To verify the quality of the synthetic cases, we had two licensed physicians review the data. Each physician was given 50 randomly sampled synthetic cases and asked to assess whether any cases posed a risk of errors. Based on their feedback, they identified that out of the 100 cases, only 6 might be incorrect, as the symptom information was less likely to be associated with the respective diseases. This suggests that synthesizing cases from a medical encyclopedia is a fairly reliable method.

L Limitations

Despite its promising performance in diagnostic tasks, DiagnosisGPT has several limitations that must be considered:

- **Limited Disease Coverage:** DiagnosisGPT is trained to identify only a specific set of diseases. This constraint means that the model’s diagnostic capabilities are confined to this predefined list, and it may not recognize or provide accurate diagnoses for conditions that fall outside its training parameters. Consequently, this limitation could hinder the model’s applicability in a real-world medical setting where a wide range of diseases, including rare and emerging conditions, need to be diagnosed.
- **Synthetic Data Annotation:** The dataset used to train DiagnosisGPT relies on annotations created by Large Language Models (LLMs). While utilizing LLMs for annotation is a cost-effective approach, it raises concerns about the quality and reliability of the data. LLMs can sometimes generate plausible but incorrect information—often referred to as "hallucinations"—which can introduce biases or errors into the training data. This could potentially lead to the model making incorrect or misleading diagnoses.
- **Reliance on Synthetic Cases:** DiagnosisGPT’s training is based on synthetic medical cases, which are constructed to avoid the privacy concerns associated with using real patient

data. However, these synthetic cases may not always accurately reflect the complexity and variability of actual patient presentations. The nuances of real-life medical conditions, including co-morbidities and patient-specific factors, are difficult to replicate in artificial scenarios. This gap between the training data and real-world contexts may impact the model’s diagnostic accuracy and its generalizability to real patient populations.

M Impact

M.1 Positive Impact

- **Promotes medical AI development:** DiagnosisGPT promotes the development of medical AI, as diagnostics are crucial in healthcare AI. Accurate diagnostic capabilities enhance patient outcomes and streamline clinical processes.
- **Improves interpretability in healthcare:** DiagnosisGPT improves the interpretability of medical AI by utilizing a disease retriever function and knowledge base integration. This increased interpretability builds trust in AI systems among healthcare providers and patients. By making the diagnostic process more transparent, DiagnosisGPT helps users understand the reasoning behind AI-generated suggestions, fostering greater confidence in AI-assisted medical practices.
- **Addresses privacy concerns in medical cases:** DiagnosisGPT offers a solution to privacy issues prevalent in medical case handling by constructing cases using a knowledge base, thereby eliminating patient privacy concerns. This approach also alleviates the problem of data scarcity.
- **Assists healthcare professionals:** DiagnosisGPT assists healthcare professionals by rapidly collecting patient symptom information and providing preliminary diagnoses. This capability enables medical practitioners to save time and focus on more complex aspects of patient care.

M.2 Potential Negative Impact

The development of DiagnosisGPT raises several potential risks.

- **Risk of Misdiagnosis:** Despite the promising results shown by DiagnosisGPT in diagnosis, it is crucial to underscore that at this stage, it should not be used to provide any medical advice. There is a possibility that it could provide incorrect interpretations or inaccurate diagnoses. Considering the nature of this field, our model and data will only be available for download by researchers. Our model will not be available for public use.
- **Data Privacy and Ethics:** The diagnostic field may involve ethical issues related to patient privacy. To address this, we use synthetic data. The training data for CoD is entirely generated by GPT-4, ensuring that there are no privacy or ethical concerns. As for DxBench, we constructed it using open-source licensed datasets, ensuring compliance with ethical standards.