

SFQ: A Sentence-level Framework for Medical Follow-up Questioning

Anonymous ACL submission

Abstract

Follow-up questioning is essential for safe medical decision-making but remains poorly evaluated in large language models (LLMs). In real clinical settings, patients provide incomplete information, requiring clinicians to actively elicit missing evidence before diagnosis. However, existing medical LLM benchmarks primarily evaluate diagnosis accuracy under complete or fixed case presentations, making it unclear whether models genuinely seek missing information. We propose *Sentence-level Follow-up Questioning (SFQ)*, a framework that formulates follow-up behavior as an explicit *sentence-level recall* task. Medical cases are decomposed into atomic clinical sentences, and incomplete cases are constructed by selectively hiding clinically relevant facts. Models must recover these facts through follow-up questions before committing to a diagnosis, enabling fine-grained and interpretable evaluation beyond diagnosis accuracy. Experiments on controlled sentence-abstracted benchmarks and a realistic external follow-up dataset show that recovering missing clinical information is associated with higher diagnostic accuracy. Recall-aware post-training improves both information recovery and diagnostic performance under incomplete inputs, demonstrating that explicitly modeling *follow-up questioning* can lead to safer and more reliable medical reasoning.

1 Introduction

Large language models (LLMs) have achieved strong performance in reasoning and text generation, particularly on question-answering tasks (Brown et al., 2020; Ouyang et al., 2022; Berger et al., 2025). However, many real-world applications are inherently interactive and require multi-turn conversations rather than single-shot responses (Chen et al., 2017; Biancofiore et al., 2024). In such settings, models must not only answer questions, but also identify missing information and



Figure 1: A safer and more effective medical LLM should conduct **follow-up questioning** to gather missing information from user inputs rather than making immediate judgments.

actively elicit it through follow-up questions. This information-seeking capability is critical when user inputs are partial, ambiguous, or underspecified; without it, models are forced to rely on incomplete evidence, increasing the risk of overconfident or unreliable outputs.

The limitations of follow-up questioning are particularly serious in medical decision-making (Bender et al., 2021; Yu et al., 2025). In real clinical settings, patients rarely describe their conditions completely or clearly. Initial conversations are often incomplete, messy, and informal. Clinicians must ask targeted follow-up questions to gradually gather the information needed for diagnosis. Making diagnosis decisions without sufficient inquiry is unsafe in human medicine, and this risk is amplified when AI systems are used to support clinical decisions. As illustrated in Figure 1, a competent medical assistant should first identify

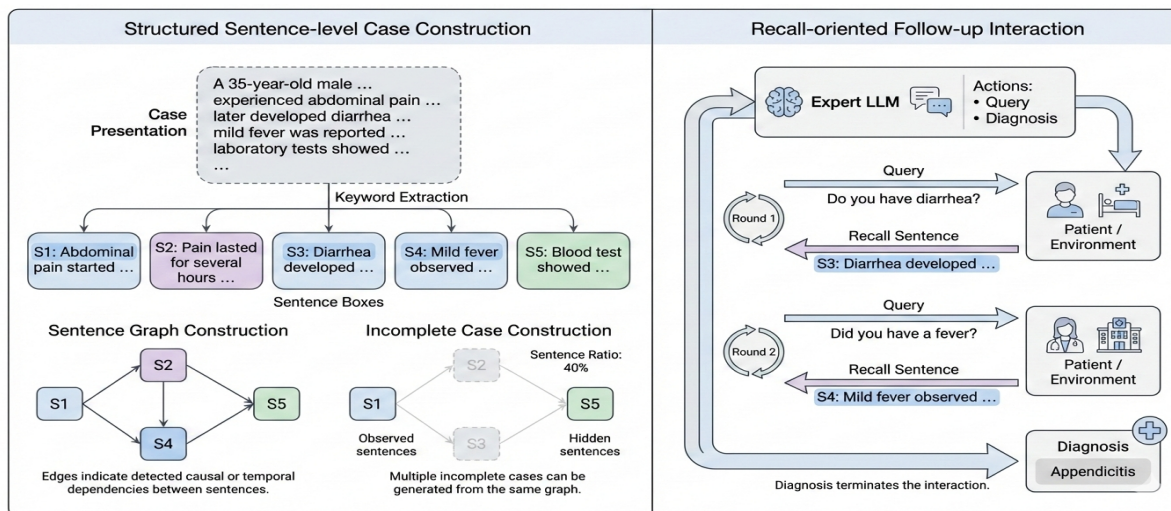


Figure 2: Our **sentence-level** formulation of medical follow-up questioning. Incomplete cases are constructed by selectively **hiding atomic clinical sentences** (left), and follow-up questioning is framed as a **recall-oriented process** (right). Diagnosis is treated as a terminal decision following multi-rounds informative interaction.

missing information and actively elicit it, rather than immediately producing a diagnosis.

Despite its importance, follow-up questioning remains a major weakness of current LLMs (Zhang et al., 2024; Meng et al., 2023; Huang et al., 2025b). Beyond architectural or inference-related limitations, a key issue lies in prevailing training paradigms. Most supervision is outcome-oriented, relying on single-turn question-answer pairs where all relevant information is already provided. In contrast, data that explicitly demonstrates how and when to ask effective follow-up questions is scarce. As a result, LLMs often fail to treat follow-up questioning as a distinct and necessary capability (Zhang et al., 2024; Huang et al., 2025b).

Although research on medical LLMs has grown rapidly, existing benchmarks and methods do not adequately evaluate follow-up questioning behavior (Liu et al., 2024; Yu et al., 2025). Most benchmarks assess diagnostic accuracy under complete or fixed case presentations (Williams et al., 2024; Luk et al., 2024). Reasoning-enhanced or dialogue-based approaches may generate follow-up questions, but these are often implicit, prompt-driven, or external to the model’s core behavior, making them difficult to control and reproduce (Wei et al., 2023; Gendron et al., 2024). Consequently, existing evaluations cannot reliably determine whether a model actively seeks missing information or simply infers a diagnosis based on incomplete evidence.

To address this gap, we propose *Sentence-level Follow-up Questioning* (SFQ), a framework for

defining and evaluating medical follow-up questioning using a sentence-level abstraction. Instead of treating medical cases as unstructured text, we decompose case presentations into *atomic clinical sentences*, each representing a single factual unit. As illustrated in Figure 2, this abstraction enables controllable construction of incomplete cases by selectively hiding subsets of sentences while preserving their causal or temporal relationships. Within this formulation, follow-up questioning is reframed as a *recall-oriented process*, requiring models to actively retrieve missing clinical facts through targeted queries before committing to a diagnosis.

Building on this formulation, we further optimize follow-up behavior through post-training. We apply supervised fine-tuning to expose models to examples of effective follow-up questions under incomplete information, and employ Group Relative Policy Optimization (GRPO, Shao et al. (2024)) to explicitly model follow-up questioning and diagnosis as distinct actions. Queries are rewarded for retrieving missing clinical facts, while redundant or uninformative questions are penalized, encouraging the model to balance continued inquiry with timely diagnosis.

We evaluate the proposed framework on sentence-abstracted medical cases and existing follow-up benchmarks, demonstrating consistent improvements in sentence-level recall and diagnostic performance under incomplete information. In summary, our contributions are as follows:

- We introduce a *sentence-level abstraction* of medical case presentations that enables controllable construction of incomplete cases and reframes follow-up questioning as a *recall-oriented evaluation problem* beyond final diagnosis accuracy.
- We propose a post-training framework combining supervised fine-tuning and Group Relative Policy Optimization (GRPO) to explicitly optimize information-seeking behavior, modeling *follow-up questioning* and *diagnosis* as distinct actions. The model’s ability of follow-up questioning has been significantly enhanced after training.

2 Related Works

2.1 LLM Follow-up Questioning

Recent studies have explored follow-up questioning as a means for LLMs to acquire missing information under underspecified inputs. Meng et al. (2023) formalizes information-seeking follow-up question generation and introduces a benchmark emphasizing diverse questioning strategies, while Zhao et al. (2025b) leverages real user–LLM conversation logs to mine follow-up intents for more realistic question generation. Related work also studies structured forms of questioning, such as organizing follow-up questions according to Bloom’s taxonomy (Yadav et al., 2025). Beyond explicit question generation, recent efforts investigate proactive information acquisition under incomplete problem settings, evaluating whether models choose to ask for missing information rather than directly answering (Huang et al., 2025b), and proposing training or benchmarking frameworks that encourage evidence gathering before decision making (Huang et al., 2025a; Zhou et al., 2025). Similar ideas appear in conversational information retrieval, where intent-aware clarification and iterative clarification–rewriting are used to progressively refine user queries (Zhao et al., 2024b).

In the medical domain, follow-up questioning has been studied to support patient–provider interactions and interactive clinical reasoning. Prior work proposes benchmarks and systems for question-asking in clinical settings (Li et al., 2024), while multi-agent conversational frameworks demonstrate that structured follow-up interactions can improve diagnostic performance in complex cases (Gatto et al., 2025; Li et al., 2023).

2.2 Medical Case Diagnosis Prediction

Medical case diagnosis prediction has been extensively studied with large language models, primarily under the assumption that complete or fixed case presentations are provided as input (Singhal et al., 2022; Meyer et al., 2024). Recent benchmarks evaluate diagnostic accuracy by mapping structured or semi-structured clinical descriptions to diagnoses, ICD codes, or clinical labels (Michalopoulos et al., 2022; Nori et al., 2023). Subsequent work incorporates reasoning-oriented techniques, such as chain-of-thought prompting, retrieval-augmented generation, or tool use, to improve diagnostic performance and robustness (Chen et al., 2024; Zhao et al., 2025a; Garza et al., 2025).

More recent studies explore interactive or dialogue-based medical assistants, allowing models to engage in multi-turn conversations with patients or simulated users to gather additional information (Zeng et al., 2020; Tu et al., 2024). While these approaches move closer to real clinical interactions, evaluations typically focus on final diagnostic outcomes or conversation-level success, with limited control over what information is missing or how it is acquired. As a result, the abilities of LLMs to actively seek and recall specific missing clinical facts remain less evaluated.

3 Framework

3.1 Sentence-level Case Abstraction

As illustrated in Figure 2 (left), real-world medical case presentations are typically provided as free-text descriptions, where multiple clinical facts—such as symptoms, temporal progression, and examination results—are interwoven in an unstructured manner. To enable controllable evaluation of follow-up questioning, we introduce a sentence-level abstraction that decomposes a case presentation into atomic and interpretable clinical units. Formally, a medical case is represented as:

$$\mathcal{C} = (P, Y), \quad (1)$$

where P denotes the original case presentation and Y is the ground-truth final diagnosis. Given P , we construct a set of sentence-level clinical facts

$$\mathcal{S} = \{s_1, s_2, \dots, s_n\}, \quad (2)$$

where each sentence s_i corresponds to an atomic piece of clinically meaningful information, such as a symptom description, a temporal event, or a

test finding. These sentence units are designed to be minimal and self-contained, serving as the basic elements for information recall.

The sentence set S is generated using an LLM with domain-specific prompting, which rewrites the original presentation into a list of concise, patient-style sentences while preserving the underlying clinical content. Compared to paragraph-level or dialogue-level abstractions, sentence-level decomposition provides finer granularity and clearer attribution of information, enabling explicit control over what clinical facts are observed or hidden.

Given the full sentence set S , we construct an **incomplete case** by selecting a subset $S_{\text{init}} \subset S$, while the remaining sentences $S_{\text{hid}} = S \setminus S_{\text{init}}$ represent clinically relevant but unobserved information to be elicited through follow-up questioning. By varying the ratio $|S_{\text{init}}|/|S|$, multiple incomplete cases can be generated from the same underlying presentation.

3.2 Sentence Dependency Graph

While sentence-level abstraction enables controllable information hiding, clinical sentences are not independent. Medical case presentations often contain temporal or causal dependencies between sentences, and ignoring such structure may result in incoherent or clinically implausible incomplete cases.

For example, consider the following two sentences:

- (1) “*I had abdominal pain yesterday.*”
- (2) “*After several hours of abdominal pain, I developed diarrhea.*”

Sentence (2) is temporally and causally dependent on sentence (1). If sentence (1) is observed while sentence (2) is hidden, the incomplete case remains coherent and reflects a realistic patient description. In contrast, hiding sentence (1) while exposing sentence (2) leads to an invalid setting, as the latter presupposes the former. This motivates a structure-aware hiding strategy rather than fully random sentence masking.

To capture such dependencies, we construct a **sentence dependency graph**

$$\mathcal{G} = (S, \mathcal{E}), \quad (3)$$

where each node corresponds to a sentence $s_i \in S$, and directed edges $(s_i \rightarrow s_j) \in \mathcal{E}$ indicate detected temporal or causal dependencies between

sentences. The resulting graph is a directed acyclic graph (DAG), reflecting the progression and logical structure of the clinical narrative, as shown in Figure 2 (left).

A simplified selection procedure is shown below as Algorithm 1.

Algorithm 1 Sentences Dependency Graph

Require: Sentence graph $\mathcal{G} = (S, \mathcal{E})$, target ratio r

Ensure: Observed sentence set S_{init}

```

1: Initialize  $S_{\text{init}} \leftarrow \emptyset$ 
2: Initialize queue  $Q$  with root nodes of  $\mathcal{G}$ 
3: while  $|S_{\text{init}}|/|S| < r$  and  $Q$  not empty do
4:    $s \leftarrow Q.\text{pop}()$ 
5:   Add  $s$  to  $S_{\text{init}}$ 
6:   for each child  $c$  of  $s$  in  $\mathcal{G}$  do
7:      $Q.\text{push}(c)$ 
8:   end for
9: end while
10: return  $S_{\text{init}}$ 

```

By incorporating sentence dependencies into the hiding process, our framework generates incomplete cases that are both controllable and clinically consistent. This structure-aware construction forms the basis for meaningful evaluation of follow-up questioning behavior, which we describe next.

3.3 Follow-up Interaction Protocol

As illustrated in Figure 1 (right), we model medical follow-up as a multi-round interaction process between an expert agent and a patient (or environment) agent. Following prior work on interactive and multi-agent diagnosis systems, we adopt a follow-up interaction protocol and adapt it to our sentence-level case abstraction, enabling explicit modeling and evaluation of information recall during diagnosis. Under this protocol, the expert agent interacts with the patient agent over multiple rounds. At each round t , the expert agent selects one of two actions,

$$a_t \in \{\text{query}, \text{diagnosis}\}, \quad (4)$$

where the query action corresponds to asking a follow-up question to acquire missing clinical information, and the diagnosis action outputs a final diagnostic prediction and terminates the interaction.

The patient agent has access to the full sentence set S and responds to each query by returning

the subset of hidden sentences that semantically match the question, if any; otherwise, it returns an empty response. The patient agent only reveals facts present in the original case presentation, ensuring that follow-up questioning corresponds to recalling existing sentence-level information.

The interaction proceeds until the expert agent outputs a diagnosis. The final diagnosis is evaluated against the ground-truth label Y , while recalled sentences are recorded throughout the interaction. This protocol grounds follow-up questioning in a sentence-level information space, enabling explicit analysis of information-seeking behavior under incomplete clinical evidence.

4 Model Optimization

4.1 Supervised Fine-tuning

We first apply supervised fine-tuning (SFT) to teach the model what to ask given an incomplete case presentation. Under the proposed sentence-level framework, SFT is formulated as a conditional generation task that maps observed sentences to follow-up questions targeting hidden clinical facts.

Formally, for each case with sentence set S , we construct an incomplete case by selecting an observed subset $S_{\text{init}} \subset S$ and a hidden subset $S_{\text{hid}} = S \setminus S_{\text{init}}$. During data preprocessing, each hidden sentence $s \in S_{\text{hid}}$ is paired with a canonical follow-up question q_s that would elicit the corresponding information. The supervision signal for a case is thus defined as the set of target questions:

$$Q_{\text{hid}} = \{q_s \mid s \in S_{\text{hid}}\}. \quad (5)$$

Given the observed sentences S_{obs} as input, the model is trained to generate the corresponding follow-up questions Q_{hid} . As illustrated in Figure 2, SFT is implemented as a single-round generation process without interaction, where the model outputs a sequence (or list) of follow-up questions conditioned on the incomplete case. This formulation avoids reliance on multi-turn dialogue during training and provides a stable learning signal aligned with sentence-level information recall.

We fine-tune a pretrained large language model using standard maximum likelihood estimation, minimizing the negative log-likelihood of the target questions given the observed sentences:

$$\mathcal{L}_{\text{SFT}} = - \sum_{q \in Q_{\text{hid}}} \log p_{\theta}(q \mid S_{\text{init}}). \quad (6)$$

During training, loss is computed only over the generated follow-up questions, while the input context is masked to prevent spurious supervision. SFT equips the model with a basic understanding of which follow-up questions are clinically relevant under incomplete information. However, it does not explicitly model when to ask questions or when to stop and make a diagnosis. We address these limitations by introducing reinforcement-based optimization in the next section.

4.2 Reinforcement Optimization with GRPO

While supervised fine-tuning teaches the model what follow-up questions are relevant, it does not explicitly model the decision process of follow-up interaction, such as when to ask questions, which questions to prioritize, and when to stop asking and commit to a diagnosis. To address these limitations, we further optimize the model using Group Relative Policy Optimization (GRPO) under the proposed follow-up interaction framework (Shao et al., 2024).

We formulate follow-up questioning as a sequential decision-making problem. At each round t , the expert agent observes the current state

$$s_t = (S_{\text{init}}^{(t)}, h_t), \quad (7)$$

where $S_{\text{init}}^{(t)}$ denotes the set of observed sentences after t rounds, and h_t represents the interaction history. The agent then samples an action

$$a_t \sim \pi_{\theta}(a_t \mid s_t), \quad a_t \in \{\text{query}, \text{diagnosis}\}, \quad (8)$$

according to the current policy π_{θ} .

If the agent selects the query action, it generates a follow-up question q_t , and the environment returns a response by revealing a subset of hidden sentences that match the query. If the agent selects the diagnosis action, the interaction terminates and the agent outputs a final diagnosis \hat{Y} .

4.2.1 Reward Design

The reward function is designed to explicitly encourage effective information-seeking behavior while discouraging inefficient or premature decisions. In particular, rewards are assigned at each step based on two principles:

1. **Sentence-level recall reward.** When a query action successfully recalls at least one previously hidden sentence, the agent receives a positive reward. Queries that fail to retrieve

any new information incur a penalty, which increases with the interaction length to discourage redundant questioning.

- Diagnosis reward.** Upon termination, the agent receives a positive reward if the predicted diagnosis \hat{Y} matches the ground-truth label Y , and a negative reward otherwise.

Formally, the total return for an interaction trajectory $\tau = (s_1, a_1, \dots, s_T, a_T)$ is defined as

$$R(\tau) = \sum_{t=1}^{T-1} r_{\text{query}}(s_t, a_t) + r_{\text{diag}}(s_T, a_T), \quad (9)$$

where r_{query} and r_{diag} denote the query-level and diagnosis-level rewards, respectively.

4.2.2 GRPO Optimization Objective

Following GRPO, policy updates are performed by comparing the relative returns of multiple sampled trajectories under the same initial state. Given a group of trajectories $\{\tau_i\}_{i=1}^N$, we compute a normalized advantage

$$\hat{A}_i = \frac{R(\tau_i) - \mu_R}{\sigma_R}, \quad (10)$$

where μ_R and σ_R are the mean and standard deviation of returns within the group. The policy is then optimized by maximizing

$$\mathcal{L}_{\text{GRPO}} = \mathbb{E}_{\tau_i \sim \pi_\theta} \left[\hat{A}_i \sum_t \log \pi_\theta(a_t | s_t) \right], \quad (11)$$

which encourages trajectories with higher relative returns while maintaining training stability.

By combining sentence-level recall rewards with diagnosis-level supervision, GRPO enables the model to jointly learn “*which questions to ask, in what order, and when to stop asking.*” This reinforcement-based optimization complements supervised fine-tuning and allows the model to better align follow-up questioning behavior with diagnostic outcomes under incomplete information.

5 Experimental Setup

5.1 Datasets

We evaluate our framework on two medical follow-up benchmarks that reflect complementary interaction settings.

Med-Sentences We construct a sentence-level follow-up benchmark by transforming cases from *MedCaseReasoning* (Wu et al., 2025) into atomic clinical sentences as described in Section 3. Given a case presentation represented as a sentence set S , we generate incomplete cases by selectively hiding a subset $S_{\text{hid}} \subset S$, while exposing the remaining sentences S_{init} . We consider two incompleteness ratios: *Quarter* (25%) and *Half* (50%). Each hidden sentence is associated with a canonical follow-up question generated during preprocessing, forming a controllable and reproducible evaluation setting for follow-up questioning.

Follow-up Text To assess external validity, we additionally evaluate on the dataset introduced by (Li et al., 2023). This dataset provides an initial case description followed by expert-defined follow-up texts (e.g., additional examinations or findings), which serve as a natural form of information supplementation. Unlike Med-Sentences, this benchmark does not explicitly define sentence-level hiding, but reflects realistic follow-up information acquisition in clinical diagnosis.

Together, these two datasets allow us to evaluate both *controlled sentence-level recall* and *realistic follow-up information utilization*.

5.2 Models

We evaluate a range of representative large language models, including Qwen2.5-7B-Instruct, GPT-5, GPT-5 mini, Llama3.1-8b and DeepSeek-V3.2 (Qwen et al., 2025; OpenAI et al., 2024; Grattafori et al., 2024; DeepSeek-AI et al., 2025). We report results for our fine-tuned models (Ours (SFT) and Ours (SFT+GRPO)). All models follow the same interaction protocol and prompt format to ensure fair comparison.

5.3 Interaction Protocol

All evaluations follow the recall-oriented follow-up interaction protocol described in Section 3.3. At each turn, the model selects one of two actions: (1) issuing a follow-up query to request additional information, or (2) producing a final diagnosis, which terminates the interaction. For Med-Sentences, the environment returns hidden sentences whose content matches the model’s query; for Follow-up Text, relevant follow-up descriptions are returned. A maximum number of interaction rounds is enforced to prevent unbounded querying.

Model	Med-Sentences (Quarter)		Med-Sentences (Half)		Follow-up Text	
	IRR	HKR	IRR	HKR	IRR	HKR
GPT-5 mini	28.89%	23.1%	53.3%	37.9%	73.1%	61.2%
DeepSeek-V3.2	27.59%	22.5%	54.2%	38.1%	73.3%	63.8%
Llama3.1-8b	27.38%	22.0%	51.9%	37.4%	72.5%	60.6%
Qwen-2.5-7b-Instruct	27.42%	21.6%	53.1%	37.6%	72.3%	60.5%
Ours (SFT)	44.30%	27.8%	57.2%	40.3%	76.8%	67.7%
Ours (SFT+GRPO)	42.80%	27.0%	57.7%	40.8%	77.8%	67.9%

Table 1: Sentence-level information recall rate (IRR) and keyword hit rate (HKR) on sentence-abstracted medical cases with different observation ratios and on follow-up text. Our models (SFT and SFT+GRPO) are post-trained on Qwen2.5-7B-Instruct. The results clearly demonstrate that although current LLMs generally do not possess the ability to actively conduct follow-up questioning, with the appropriate training data for post-training, this capability can be **significantly enhanced**.

5.4 Evaluation Metrics

We evaluate models from both information acquisition and diagnostic performance perspectives.

Information Recall Rate (IRR). IRR measures the proportion of case information recovered during interaction:

$$\text{IRR} = \frac{|S_{\text{init}} \cup \hat{S}_{\text{rec}}|}{|S|}, \quad (12)$$

where \hat{S}_{rec} denotes the set of sentences recalled through follow-up questioning.

Hit Keyword Rate (HKR). To assess the quality of recalled information, we define HKR as the fraction of diagnosis-relevant keywords covered by recalled sentences. Keywords are generated during preprocessing by LLMs with access to the ground-truth diagnosis:

$$\text{HKR} = \frac{|\text{KW}(S_{\text{init}} \cup \hat{S}_{\text{rec}})|}{|\text{KW}(S)|}. \quad (13)$$

Diagnosis Accuracy (Acc). We report Top-1 and Top-3 diagnostic accuracy. Due to the inherent variability in medical diagnosis phrasing, we adopt a judge-based evaluation scheme using Qwen2.5-7B and GPT-5 as diagnosis judges to determine semantic equivalence between predicted and ground-truth diagnoses.

6 Results

6.1 Overall Performance

Table 1 shows the main results on Med-Sentences (Quarter / Half) and Follow-up Text. Across all settings, our approach consistently improves sentence-level information acquisition while maintaining or improving diagnostic performance.

On *Med-Sentences (Quarter)*, our SFT model achieves an IRR of 44.30%, substantially outperforming all base models, indicating a stronger ability to actively recover hidden clinical information under highly incomplete inputs. Incorporating GRPO further stabilizes performance, yielding comparable IRR (42.80%) while maintaining competitive HKR. Similar trends are observed under the Half setting, where our models achieve the highest IRR (57.7%) and HKR (40.8%), demonstrating robustness as case completeness increases.

On the *Follow-up Text benchmark*, which reflects more realistic follow-up information acquisition, our GRPO-optimized model achieves the strongest overall performance, with IRR reaching 77.8% and HKR 67.9%, outperforming both general-purpose and domain-specific baselines. These results suggest that explicit optimization for follow-up questioning generalizes beyond controlled sentence-level settings.

6.2 Follow-up Recall on Diagnostic Accuracy

Tables 2-4 present diagnosis accuracy evaluated under different recall agents and judge models. Across all datasets, introducing an explicit recall agent consistently improves Top-1 and Top-3 diagnostic accuracy.

On *Med-Sentences (Half)* (Table 2), when evaluated by a GPT-5 judge, diagnosis accuracy improves from 22.2% (no recall agent) to 26.1% with our GRPO-trained recall agent, with Top-3 accuracy increasing to 41.2%. Similar gains are observed when using Qwen2.5-7B as the judge, indicating that improvements are not judge-specific.

Results on *Med-Sentences (Quarter)* (Table 3) further highlight the importance of follow-up under more severe information sparsity. Our model

Med-Sentences (Half)			
Judge	Recall_agent	Top-1 Acc	Top-3 Acc
Qwen2.5-7b	None	9.96%	16.8%
	Qwen2.5-7b	10.7%	17.2%
	gpt-5-mini	10.8%	17.2%
	Ours (SFT+GRPO)	12.3%	21.4%
gpt-5	None	22.2%	37.1%
	Qwen2.5-7b	24.1%	38.2%
	gpt-5-mini	24.7%	38.7%
	Ours (SFT+GRPO)	26.1%	41.2%

Table 2: Top-1 and Top-3 diagnostic accuracy with different recall agents on Med-Sentences (Half), evaluated using Qwen2.5-7B and GPT-5 as diagnosis judges.

Med-Sentences (Quarter)			
Judger	Recall_agent	Top-1 Acc	Top-3 Acc
Qwen2.5-7b	None	8.45%	14.1%
	Qwen2.5-7b	8.83%	14.4%
	gpt-5-mini	8.75%	15.2%
	Ours (SFT+GRPO)	9.84%	17.2%
gpt-5	None	13.3%	22.9%
	Qwen2.5-7b	15.1%	25.7%
	gpt-5-mini	15.4%	26.3%
	Ours (SFT+GRPO)	19.6%	31.7%

Table 3: Top-1 and Top-3 diagnostic accuracy with different recall agents on Med-Sentences (Quarter), evaluated using Qwen2.5-7B and GPT-5 as diagnosis judges.

achieves the highest Top-1 accuracy (9.84% under Qwen2.5-7B and 19.6% under GPT-5), outperforming baselines by a clear margin, despite the increased difficulty of the setting.

Follow-up Text			
Judger	Recall_agent	Top-1 Acc	Top-3 Acc
Qwen2.5-7b	None	16.1%	22.5%
	Qwen2.5-7b	16.3%	22.9%
	gpt-5-mini	18.1%	25.1%
	Ours (SFT+GRPO)	18.9%	26.8%
gpt-5	None	31.5%	40.6%
	Qwen2.5-7b	32.8%	40.9%
	gpt-5-mini	33.3%	42.1%
	Ours (SFT+GRPO)	36.1%	45.8%

Table 4: Top-1 and Top-3 diagnostic accuracy with different recall agents on Follow-up Text, evaluated using Qwen2.5-7B and GPT-5 as diagnosis judges.

On Follow-up Text (Table 4), the benefits of recall-oriented interaction are most pronounced. Under GPT-5 evaluation, our GRPO-trained model reaches 36.1% Top-1 accuracy and 45.8% Top-3 accuracy, compared to 31.5% / 40.6% without follow-up reasoning. This confirms that improved information recall directly translates into more reliable final diagnoses in realistic follow-up scenarios.

6.3 Case Study and Error Analysis

To better understand model behavior beyond aggregate metrics, we conduct a qualitative case study based on a rare disease case and our SFT model. Overall, we observe that effective sentence-level recall often leads to more informative follow-up interactions, but correct recall alone does not always guarantee correct diagnosis, particularly for rare diseases with overlapping clinical features. A detailed analysis of a representative case is provided in Appendix 7.

6.4 Summary

Overall, these results demonstrate that:

- Explicitly modeling follow-up questioning as a recall-oriented process significantly improves information acquisition (IRR, HKR).
- Improvements in sentence-level recall consistently lead to higher diagnostic accuracy across datasets and judges.
- Reinforcement optimization with GRPO further enhances robustness and generalization, particularly in realistic follow-up settings.

Together, these findings validate sentence-level recall as a meaningful and actionable evaluation signal for medical follow-up questioning.

7 Discussion

Our results highlight the importance of explicitly modeling follow-up questioning as an information acquisition problem rather than treating it as a by-product of diagnosis-oriented reasoning. Across both controlled sentence-level benchmarks and realistic follow-up settings, we observe a consistent correlation between improved information recall and diagnostic accuracy, suggesting that diagnostic errors often stem from missing evidence rather than insufficient inference.

The sentence-level abstraction adopted in this work provides a controllable and interpretable perspective for analyzing follow-up behavior. By decomposing free-text case presentations into atomic clinical facts, our framework enables fine-grained evaluation of what information is missing, what is actively recovered, and how such recovery influences downstream diagnostic decisions.

603 Limitations

604 This work has several limitations that point to di-
605 rections for future research.

- 606 • First, while our framework provides a con-
607 trolled and interpretable abstraction for eval-
608 uating follow-up questioning, it focuses on
609 sentence-level factual information and does
610 not aim to fully capture all aspects of real-
611 world clinical decision-making.
- 612 • Second, our evaluation is conducted in a sim-
613 ulated interaction setting, which allows for
614 systematic comparison but may differ from
615 open-ended patient–clinician conversations in
616 practice.
- 617 • Finally, we primarily study a limited set
618 of model backbones and optimization strate-
619 gies, and exploring broader model families
620 and training paradigms may further improve
621 follow-up reasoning capabilities.

622 References

623 Emily Bender, Timnit Gebru, Angelina McMillan-
624 Major, and Shmargaret Shmitchell. 2021. [On the
625 dangers of stochastic parrots: Can language models
626 be too big?](#) pages 610–623.

627 Armin Berger, Sarthak Khanna, David Berghaus, and
628 Rafet Sifa. 2025. [Reasoning llms in the medi-
629 cal domain: A literature survey.](#) *arXiv preprint*
630 *arXiv:2508.19097*.

631 Giovanni Maria Biancofiore, Yashar Deldjoo, Tom-
632 maso Di Noia, Eugenio Di Sciascio, and Fedelucio
633 Narducci. 2024. [Interactive question answering sys-
634 tems: Literature review.](#) *ACM Computing Surveys*,
635 56(9):1–38.

636 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
637 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
638 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
639 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
640 Gretchen Krueger, Tom Henighan, Rewon Child,
641 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
642 Clemens Winter, and 12 others. 2020. [Lan-
643 guage models are few-shot learners.](#) *Preprint*,
644 *arXiv:2005.14165*.

645 Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang
646 Tang. 2017. [A survey on dialogue systems: Recent
647 advances and new frontiers.](#) *ACM SIGKDD Explor-
648 ations Newsletter*, 19(2):25–35.

649 Junying Chen, Chi Gui, Anningzhe Gao, Ke Ji, Xidong
650 Wang, Xiang Wan, and Benyou Wang. 2024. [Cod,
651 towards an interpretable medical agent using chain
652 of diagnosis.](#) *Preprint*, *arXiv:2407.13301*.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-
uan Wang, Bochao Wu, Chengda Lu, Chenggang
Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
Damai Dai, Daya Guo, Dejian Yang, Deli Chen,
Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,
and 181 others. 2025. [Deepseek-v3 technical report.](#)
Preprint, *arXiv:2412.19437*.

Leon Garza, Anantaa Kotal, Michael A. Grasso, and
Emre Umucu. 2025. [Retrieval-augmented framework
for llm-based clinical decision support.](#) *Preprint*,
arXiv:2510.01363.

Joseph Gatto, Parker Seegmiller, Timothy Burdick,
Inas S. Khayal, Sarah DeLozier, and Sarah M.
Preum. 2025. [Follow-up question generation for
enhanced patient-provider conversations.](#) *Preprint*,
arXiv:2503.17509.

Gaël Gendron, Qiming Bao, Michael Witbrock, and
Gillian Dobbie. 2024. [Large language mod-
els are not strong abstract reasoners.](#) *Preprint*,
arXiv:2305.19555.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
Abhinav Pandey, Abhishek Kadian, Ahmad Al-
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
tra, Archie Sravankumar, Artem Korenev, Arthur
Hinsvark, and 542 others. 2024. [The llama 3 herd of
models.](#) *Preprint*, *arXiv:2407.21783*.

Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,
Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen,
Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun
Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni,
and Jian Guo. 2025. [A survey on llm-as-a-judge.](#)
Preprint, *arXiv:2411.15594*.

Tenghao Huang, Sihao Chen, Muhao Chen, Jonathan
May, Longqi Yang, Mengting Wan, and Pei Zhou.
2025a. [Teaching language models to gather infor-
mation proactively.](#) *Preprint*, *arXiv:2507.21389*.

Youcheng Huang, Bowen Qin, Chen Huang, Duan-
yu Feng, Xi Yang, and Wenqiang Lei. 2025b. [Beyond
solving math quiz: Evaluating the ability of large
reasoning models to ask for information.](#) *Preprint*,
arXiv:2508.11252.

Jian Li, Xi Chen, Weizhi Liu, Li Wang, Yingman Guo,
Mingke You, Gang Chen, and Kang Li. 2023. [One
is not enough: Multi-agent conversation framework
enhances rare disease diagnostic capabilities of large
language models.](#)

Shuyue Stella Li, Vidhisha Balachandran, Shangbin
Feng, Jonathan S. Ilgen, Emma Pierson, Pang Wei
Koh, and Yulia Tsvetkov. 2024. [Mediq: Question-
asking llms and a benchmark for reliable interactive
clinical reasoning.](#) *Preprint*, *arXiv:2406.00922*.

Mianxin Liu, Jinru Ding, Jie Xu, Weiguo Hu, Xiaoyang
Li, Lifeng Zhu, Zhian Bai, Xiaoming Shi, Benyou
Wang, Haitao Song, Pengfei Liu, Xiaofan Zhang,

709	Shanshan Wang, Kang Li, Haofen Wang, Tong Ruan, Xuanjing Huang, Xin Sun, and Shaoting Zhang. 2024. Medbench: A comprehensive, standardized, and reliable benchmarking system for evaluating chinese medical large language models . <i>Preprint</i> , arXiv:2407.10990.	765
710		766
711		767
712		768
713		769
714		770
715	D. W. A. Luk, W. C. T. Ip, and Y. F. Shea. 2024. Performance of GPT-4 and GPT-3.5 in generating accurate and comprehensive diagnoses across medical subspecialties . <i>Journal of the Chinese Medical Association</i> , 87(3):259–260. Epub 2024 Feb 2.	771
716		772
717		773
718		774
719		775
720	Yan Meng, Liangming Pan, Yixin Cao, and Min-Yen Kan. 2023. FollowupQG: Towards information-seeking follow-up question generation . In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 252–271, Nusa Dua, Bali. Association for Computational Linguistics.	776
721		777
722		778
723		779
724		780
725		781
726		782
727		783
728		784
729	Annika Meyer, Janik Riese, and Thomas Streichert. 2024. Comparison of the performance of gpt-3.5 and gpt-4 with that of medical students on the written german medical licensing examination: Observational study . <i>JMIR Medical Education</i> , 10.	785
730		786
731		787
732		788
733		789
734	George Michalopoulos, Michal Malyska, Nicola Sahar, Alexander Wong, and Helen Chen. 2022. Icddbig-bird: A contextual embedding model for icd code classification . <i>Preprint</i> , arXiv:2204.10408.	790
735		791
736		792
737		793
738	Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems . <i>Preprint</i> , arXiv:2303.13375.	794
739		795
740		796
741		797
742	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	798
743		799
744		800
745		801
746		802
747		803
748		804
749		805
750	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . <i>Preprint</i> , arXiv:2203.02155.	806
751		807
752		808
753		809
754		810
755		811
756		812
757		813
758	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. Qwen2.5 technical report . <i>Preprint</i> , arXiv:2412.15115.	814
759		815
760		816
761		817
762		818
763		819
764		820
		821
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models . <i>Preprint</i> , arXiv:2402.03300.	821
		822
	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguera y Arcas, Dale Webster, and 11 others. 2022. Large language models encode clinical knowledge . <i>Preprint</i> , arXiv:2212.13138.	822
		823
	Tao Tu, Anil Palepu, Mike Schaeckermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, Shekoofeh Azizi, Karan Singhal, Yong Cheng, Le Hou, Albert Webson, Kavita Kulkarni, S Sara Mahdavi, Christopher Sementurs, Juraj Gottweis, and 6 others. 2024. Towards conversational diagnostic ai . <i>Preprint</i> , arXiv:2401.05654.	823
		824
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . <i>Preprint</i> , arXiv:2201.11903.	824
		825
	C. Y. K. Williams, B. Y. Miao, A. E. Kornblith, and 1 others. 2024. Evaluating the use of large language models to provide clinical recommendations in the emergency department . <i>Nature Communications</i> , 15:8236.	825
		826
	Kevin Wu, Eric Wu, Rahul Thapa, Kevin Wei, Angela Zhang, Arvind Suresh, Jacqueline J. Tao, Min Woo Sun, Alejandro Lozano, and James Zou. 2025. Medcasereasoning: Evaluating and learning diagnostic reasoning from clinical case reports . <i>Preprint</i> , arXiv:2505.11733.	826
		827
	Archana Yadav, Harshvivek Kashid, Medchalimi Sruthi, B JayaPrakash, Chintalapalli Raja Kullayappa, Mandala Jagadeesh Reddy, and Pushpak Bhattacharyya. 2025. From recall to creation: Generating follow-up questions using bloom’s taxonomy and Grice’s maxims . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 6: Industry Track)</i> , pages 1322–1338, Vienna, Austria. Association for Computational Linguistics.	827
		828
	E. Yu, X. Chu, W. Zhang, X. Meng, Y. Yang, X. Ji, and C. Wu. 2025. Large language models in medicine: Applications, challenges, and future directions . <i>International Journal of Medical Sciences</i> , 22(11):2792–2801.	828
		829
	Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, Hongchao Fang, Penghui Zhu, Shu Chen, and Pengtao Xie. 2020.	829

822 [MedDialog: Large-scale medical dialogue datasets.](#)
823 In *Proceedings of the 2020 Conference on Empirical*
824 *Methods in Natural Language Processing (EMNLP)*,
825 pages 9241–9250, Online. Association for Computa-
826 tional Linguistics.

827 Tong Zhang, Peixin Qin, Yang Deng, Chen Huang,
828 Wenqiang Lei, Junhong Liu, Dingnan Jin, Hongru
829 Liang, and Tat-Seng Chua. 2024. [Clamber: A bench-](#)
830 [mark of identifying and clarifying ambiguous infor-](#)
831 [mation needs in large language models.](#) *Preprint*,
832 arXiv:2405.12063.

833 Xuejiao Zhao, Siyan Liu, Su-Yin Yang, and Chun-
834 yan Miao. 2025a. [Medrag: Enhancing retrieval-](#)
835 [augmented generation with knowledge graph-](#)
836 [elicited reasoning for healthcare copilot.](#) *Preprint*,
837 arXiv:2502.04413.

838 Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang,
839 Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu,
840 Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda
841 Chen. 2024a. [Swift:a scalable lightweight infrastruc-](#)
842 [ture for fine-tuning.](#) *Preprint*, arXiv:2408.05517.

843 Ziliang Zhao, Zhicheng Dou, and Yujia Zhou. 2024b.
844 [Generating intent-aware clarifying questions in con-](#)
845 [versational information retrieval systems.](#) In *Pro-*
846 *ceedings of the 33rd ACM International Conference*
847 *on Information and Knowledge Management, CIKM*
848 *'24*, page 3384–3394, New York, NY, USA. Associa-
849 tion for Computing Machinery.

850 Ziliang Zhao, Shiren Song, and Zhicheng Dou. 2025b.
851 [Followgpt: A framework of follow-up question gen-](#)
852 [eration for large language models via conversation](#)
853 [log mining.](#) In *Proceedings of the 34th ACM In-*
854 *ternational Conference on Information and Knowl-*
855 *edge Management, CIKM '25*, page 4412–4422, New
856 York, NY, USA. Association for Computing Machin-
857 ery.

858 Zhanke Zhou, Xiao Feng, Zhaocheng Zhu, Jiangchao
859 Yao, Sanmi Koyejo, and Bo Han. 2025. [From passive](#)
860 [to active reasoning: Can large language models ask](#)
861 [the right questions under incomplete information?](#)
862 *Preprint*, arXiv:2506.08295.

Sentence Construction Prompt

You are a clinical text processing assistant. Your task is to transform a medical case presentation into two outputs:

1. A list of extracted clinical keywords.
2. A list of short, patient-style sentences ("I have...", "I am...", "My ...") where each sentence contains exactly ONE clinical fact.

TASK INSTRUCTIONS

Given an input case presentation, perform the following steps:

(1) KEYWORD EXTRACTION

- Extract minimal clinical units.
- Include demographics, duration, symptoms, characteristics, radiation, severity, aggravating/relieving factors, relevant history.
- Do NOT include function words or unnecessary phrases.
- Keep keywords short (1–5 words each).

(2) SENTENCE ATOMIZATION

- Convert the case into patient-style first-person statements.
- Each sentence must contain exactly ONE clinical fact.
- Use simple structures: "I am...", "I have...", "My symptoms..."
- No merging of multiple facts into one sentence, unless they share a core keyword.
- No extra explanations or new clinical information.
- If two sentences have a very clear sequence relationship, place them in one sentence.
- Minimize the number of sentences as much as possible

(3) Question Generation

- Generate a sample question for each sentence, focusing on the clinical fact it contains.
- For example, sentence "I have had a fever for 3 days." could generate the question "Do you have a fever and how long it takes?"
- Each question should be concise and directly related to the sentence.
- Each question should be in the same index with its corresponding sentence.

(4) JSON OUTPUT FORMAT

Return ONLY a JSON object with fields:

```
{
  "keywords": ["...", "..."],
  "sentences": ["...", "..."],
  "questions": ["...", "..."]
}
```

INPUT CASE PRESENTATION

{CASE-TEXT}

A. Sentence Construction and Details

A.1 Sentence Extraction Prompt

To enable controllable and fine-grained evaluation of follow-up questioning, we first transform free-text medical case presentations into a structured sentence-level representation. This process decomposes each case into atomic clinical sentences, associated keywords, and sentence-level follow-up questions, which together serve as the foundation for subsequent incomplete case construction and evaluation.

Given a raw case presentation C , we use a large language model to generate:

1. a set of extracted clinical keywords \mathcal{K} ;
2. a list of atomic, patient-style sentences $\mathcal{S} = \{s_1, \dots, s_n\}$, where each sentence contains exactly one clinical fact; and
3. a corresponding list of follow-up questions $\mathcal{Q} = \{q_1, \dots, q_n\}$, each aligned with one sentence.

The transformation is performed using the following prompt:

A.2 Keyword Generation

For each medical case, clinical keywords are generated jointly with sentence construction as described in Appendix A.1. Keywords are intended to represent minimal semantic units corresponding to the core clinical facts expressed in each atomic sentence, including symptoms, temporal descriptors, and relevant contextual attributes.

During evaluation, keyword matching is used to compute the keyword hit rate (HKR), which measures whether information elicited through follow-up questioning captures clinically relevant concepts. Keyword matching is performed by comparing recalled sentences against the keyword set associated with the original case. Minor lexical variations are allowed through normalization and relaxed matching rules, while function words and non-clinical terms are ignored.

This keyword-based signal complements sentence-level recall by providing a coarse-grained measure of semantic coverage, without relying on exact sentence matching.

A.3 Dataset Statistics and Splits

We report basic statistics of the sentence-level abstraction in Table 5. The table summarizes the average number of observed sentences in the initial

Type	Med-Sentences		
	%	Avg. Initial	Avg. Full
Quarter	25%	5.89	25.01
Half	50%	12.27	25.01

Table 5: The average length of sentences-list for different division methods.

presentation and the total number of atomic sentences per case under different observation ratios (25% and 50%). These statistics reflect the granularity of sentence decomposition and ensure that incomplete cases are constructed with comparable levels of information across experimental settings.

For the sentence-abstracted version of the MedCase-Reasoning dataset, we randomly split the data into training, validation, and test sets with a ratio of 8:1:1. All models are trained on the training set, and final results are reported on the held-out test set.

B. Recall-oriented Follow-up Interaction

B.1 Expert Prompt for Follow-up Questioning

In our recall-oriented interaction, the expert agent acts as a clinician who iteratively decides whether to ask a targeted follow-up question or to stop and output a final diagnosis. At each turn, the expert observes the currently known patient statements, the dialogue history, and the remaining turn budget, and must output exactly one action in a strict JSON format.

B.2 Patient / Environment Prompt

The patient (environment) agent serves as a non-generative information source. It holds the complete set of atomic clinical sentences and responds to the expert’s follow-up questions by returning only matched sentences, without introducing new facts or performing additional reasoning.

B.3 Multi-round Interaction Protocol

Following prior work on interactive medical reasoning and multi-turn diagnosis, we adopt a structured multi-round interaction protocol within our sentence-level recall framework. Each episode proceeds as a sequence of turns between the expert agent and the patient (environment) agent.

At each turn t , the expert observes the currently known patient statements and the dialogue history, and outputs exactly one action: either a follow-up query or a final diagnosis. If a query is issued, the patient agent returns all matched sentences from the

Patient / Environment Prompt

You are a patient who knows a set of true statements about yourself.

You will receive:

- A list of ALL-SENTENCES, each describing one true fact about you.
- A doctor’s question.

Your task is to answer ONLY based on ALL-SENTENCES.

Rules:

- If the question is clearly asking about one or more facts that are present in ALL-SENTENCES, you should return those sentences.
- If multiple sentences are relevant, you may return more than one.
- If there is no relevant information, answer that you don’t have this information.
- Do NOT invent new facts.

You MUST output ONLY a JSON object with the following format:

If one or more sentences match the doctor’s question:

```
{ "matched": true,
  "sentences": [ "<one matching sentence>", "<another matching sentence>" ] }
```

If no sentence matches:

```
{ "matched": false,
  "sentences": [] }
```

full sentence set, or an empty response if no relevant information is available. To ensure bounded interaction, the maximum number of rounds is fixed to $T = 10$ for all experiments. If no diagnosis is produced within T rounds, the expert is forced to output a diagnosis at the final step.

This protocol ensures that follow-up questioning is explicit, bounded, and fully observable, enabling precise measurement of information recall and its impact on downstream diagnostic decisions. For experiments involving local model inference, including Qwen2.5-7B-Instruct, Llama3.1-8B, and our post-trained models, all interactions are executed on a single NVIDIA A40 GPU.

C. Optimization with GRPO

To further improve follow-up questioning behavior beyond supervised fine-tuning, we apply reinforcement optimization using Group Relative Policy Optimization (GRPO). This section provides implementation details of the reward design and training procedure, which are omitted from the main text for clarity.

C.1 Action and Episode Definition

Each training sample corresponds to a static incomplete case constructed at the sentence level. An episode consists of a bounded sequence of expert

Expert Prompt

You are an experienced clinician evaluating a patient case.

You will see:

- A list of patient statements that are currently KNOWN.
- A log of previous questions you asked and the patient’s answers.
- The maximum number of turns you are allowed to use.

Your task is to decide, at each step, whether to:

- 1) ASK a focused follow-up question, or
- 2) DIAGNOSE the most likely primary diagnosis.

Rules:

- You MAY ONLY ask questions about potential symptoms, history, or risk factors that could plausibly appear as short patient-style sentences (e.g., "I have chest pain", "I have a fever", "I am a smoker", etc.).
- When you ASK, your question must have a single answer, such as: "Do you have chest tightness?", "What is your age?", "Do your symptoms radiate to your left arm?"
- When you DIAGNOSE, provide your best single main diagnosis.
- Be concise and clinically focused.
- Do NOT repeat the same question.
- Stop asking questions once you believe you have enough information to make diagnosis.

You MUST respond ONLY in valid JSON with the following format:

If you choose to ask a question:

```
{ "action": "query", "question": "Do you have ... ?" }
```

If you choose to provide a final diagnosis:

```
{ "action": "diagnosis", "diagnosis": "<your final diagnosis>" }
```

Do not include any other fields or free text.

actions

$$a_t \in \{\text{query}, \text{diagnosis}\}, \quad t = 1, \dots, T, \quad (14)$$

where $T = 10$ is the maximum number of allowed rounds. The episode terminates when a diagnosis action is produced or when the round limit is reached, in which case a diagnosis is forced.

C.2 Recall-aware Reward Design

We design a recall-aware reward function that explicitly balances information acquisition efficiency and diagnostic correctness.

1. Query reward. For a follow-up query issued at round t , the reward depends on whether the query successfully recalls at least one hidden sentence:

$$r_{\text{query}}(t) = \begin{cases} \frac{\alpha}{1 + \lambda t}, & \text{if the query matches,} \\ -\beta(1 + \lambda t), & \text{otherwise,} \end{cases} \quad (15)$$

where $\alpha > 0$ and $\beta > 0$ control the relative reward and penalty, and λ governs time-dependent decay. This formulation encourages early, informative questions while discouraging prolonged or unproductive querying.

2. Diagnosis reward. When a diagnosis action is produced, a terminal reward is assigned:

$$r_{\text{diag}} = \begin{cases} +\gamma, & \text{if the diagnosis is correct,} \\ -\gamma, & \text{otherwise,} \end{cases} \quad (16)$$

where γ is set to dominate individual query rewards, incentivizing the model to stop querying once sufficient information has been gathered.

C.3 GRPO Training Details

GRPO Training Configuration

```
train_type = lora
lora_rank = 8
lora_alpha = 16
learning_rate = 2e-5
max_length = 2048
num_train_epochs = 5
per_device_train_batch_size = 1
gradient_accumulation_steps = 8
eval_strategy = steps
eval_steps = 50
logging_steps = 20
num_generations = 4
```

GRPO training is initialized from the supervised fine-tuned (SFT) model described in Section 4.1, with Qwen2.5-7B-Instruct serving as the base model.

All GRPO training is conducted using the **ms-swift** framework (Zhao et al., 2024a) with LoRA-based parameter-efficient updates. We use 1,400 sentence-abstracted Med-Sentences training cases for reinforcement learning. Training is executed on a single NVIDIA A100 GPU. Unless otherwise specified, other optimization hyperparameters follow the default GRPO configuration provided by



Figure 3: The training loss curve and the evaluation loss curve in the supervised fine-tuning (SFT) method

ms-swift.

D. Supervised Fine-tuning Details

D.1 Training Configuration

All supervised fine-tuning experiments are conducted on top of the Qwen2.5-7B-Instruct base model. We fine-tune the model using a standard causal language modeling objective on the constructed sentence-level follow-up question data. Training is performed on a single NVIDIA A40 GPU.

We adopt a small per-device batch size combined with gradient accumulation to balance memory efficiency and optimization stability. Mixed-precision training with bfloat16 is enabled, together with gradient checkpointing, to reduce memory consumption. Model checkpoints are saved periodically, and only a limited number of recent checkpoints are retained.

The key training hyperparameters are summarized in the configuration box below.

```

Supervised Fine-tuning Configuration

CONFIG =
per_device_train_batch_size = 2
per_device_eval_batch_size = 2
gradient_accumulation_steps = 8
num_train_epochs = 3
learning_rate = 5e-5
logging_steps = 50
eval_strategy = "steps"
eval_steps = 100
bf16 = True
gradient_checkpointing = True

```

D.2 Training Performance

Figure 3 illustrates the training and evaluation loss curves during supervised fine-tuning. The model is trained for three epochs, and both training and evaluation losses decrease steadily over the course

of optimization.

The training loss shows a rapid decline in the early stage, followed by a more gradual convergence, indicating stable optimization behavior. The evaluation loss closely tracks the training loss without divergence, suggesting limited overfitting under the chosen training configuration. Overall, the loss trends demonstrate that the supervised fine-tuning process converges smoothly and yields a well-trained initialization for subsequent reinforcement optimization.

E. LLM-based Diagnosis Evaluation

Final diagnosis correctness is evaluated using an *LLM-as-a-Judge* (Gu et al., 2025) protocol to account for synonymous expressions and minor lexical variations in free-form diagnostic outputs. Specifically, we use *GPT-5-mini* as a lightweight judge model that receives the predicted diagnosis and the ground-truth label and outputs a binary decision (yes or no) indicating whether they refer to the same primary diagnosis. The judge is instructed to be concise and does not generate explanations. This evaluation protocol is applied consistently across all models and experimental settings and is used solely for evaluation, not for training or inference.

```

LLM-as-a-Judge

Is this predicted final diagnosis correct (yes/no)?
Predicted final diagnosis: {prediction}
True final diagnosis: {truth}
Answer [yes/no].

```

We emphasize that the judge model is only used for evaluation and does not participate in training or inference of the proposed models. Since the same judge and prompt are applied consistently, relative comparisons across models remain fair and controlled.

F. Detailed Case Study

```

F.1 Case Overview

• Case ID: PMC8561223
• Ground-truth diagnosis: Pituitary adenoma
• Predicted diagnosis: Cystic pituitary adenoma
• Recall rate: 23.1% → 61.5%
• Category: Correct follow-up, incorrect diagnosis

```

F.2 Initial Observed and Hidden Sentences

Initially Observed Sentences (S_{init}):

- MRI of my brain showed a sellar cystic lesion ...
- The cyst was only partially removed ...
- My headaches worsened after a thunderclap episode ...

Hidden Sentences (S_{hid}):

- I am a 33-year-old soldier ...
- I have had intermittent headaches for the ...
- I had no focal neurologic deficits ...
- The cyst wall enhanced with gadolinium ...
- I had low thyroid-stimulating hormone ...
- I have low testosterone ...
- I did not have diabetes insipidus ...
- An endoscopic transsphenoidal approach was ...
- A yellow capsule released clear, straw-colored ...
- The pulsating cyst capsule was tightly adherent ...

evant information through follow-up questioning. However, despite high-quality recall, the final diagnosis remains incorrect due to the rarity of the true condition and its substantial overlap with more common pituitary lesions. This suggests that while sentence-level recall provides a meaningful signal for follow-up ability, diagnosis accuracy remains constrained by *prior biases* and *disease prevalence*.

1087
1088
1089
1090
1091
1092
1093
1094

F.3 Follow-up Interaction (Dialog)

Doctor Query: What type of surgery did you undergo and under what anesthesia?

Patient Answer: I don't have this information.

Doctor Query: What symptoms do you have related to vision?

Patient Answer: I have had intermittent headaches for the past year ...

Patient Answer: My headaches worsened after a thunderclap episode ...

Patient Answer: MRI of my brain showed a sellar cystic lesion ...

Doctor Query: What was found during the surgery regarding the cystic lesion?

Patient Answer: An endoscopic transsphenoidal approach was undertaken ...

Patient Answer: A yellow capsule released clear, straw-colored fluid ...

Patient Answer: The pulsating cyst capsule was tightly adherent ...

Patient Answer: The cyst was only partially removed ...

Doctor Query: What were the findings from the MRI of your brain?

Patient Answer: MRI of my brain showed a sellar cystic lesion ...

Patient Answer: The cyst wall enhanced with gadolinium ...

Final Diagnosis: *Cystic pituitary adenoma*

F.4 Analysis and Discussion

This case illustrates that the model can effectively recover a substantial portion of diagnostically rel-