

TAGS: A Test-Time Generalist–Specialist Framework with Retrieval-Augmented Reasoning and Verification

Anonymous ACL submission

Abstract

Recent advances such as Chain-of-Thought prompting have significantly improved large language models (LLMs) in zero-shot medical reasoning. However, prompting-based methods often remain shallow and unstable, while fine-tuned medical LLMs suffer from poor generalization under distribution shifts and limited adaptability to unseen clinical scenarios. To address these limitations, we present **TAGS**, a test-time framework that combines a broadly capable generalist with a domain-specific specialist to offer complementary perspectives without any model fine-tuning or parameter updates. To support this generalist–specialist reasoning process, we introduce two auxiliary modules: a hierarchical retrieval mechanism that provides multi-scale exemplars by selecting examples based on both semantic and rationale-level similarity, and a reliability scorer that evaluates reasoning consistency to guide final answer aggregation. TAGS achieves strong performance across nine MedQA benchmarks, boosting GPT-4o accuracy by 13.8%, DeepSeek-R1 by 16.8%, and improving a vanilla 7B model from 14.1% to 23.9%. These results surpass several fine-tuned medical LLMs, without any parameter updates.

1 Introduction

Large Language Models (LLMs) have shown strong potential in medical question answering (MedQA), achieving promising results on a variety of benchmarks (Singhal et al., 2025; Jin et al., 2022; Chen et al., 2023; Zhou et al., 2023). Despite these advances, recent studies reveal that even state-of-the-art models often fail on clinically challenging cases that require deep domain knowledge, multi-step reasoning, and robust generalization under distribution shifts (Xu et al., 2024; Fan et al., 2025; Shi et al., 2024).

Existing efforts to improve MedQA performance mainly follow two directions. Prompting-based approaches, including Chain-of-Thought (CoT) (Wei

et al., 2022) and multi-agent reasoning frameworks (Chen et al., 2025), aim to enhance reasoning depth through structured inference or simulated expert collaboration. However, large-scale evaluations show that such methods are frequently brittle, suffering from redundant reasoning, unstable interactions, and limited gains over single-agent baselines (Tang et al., 2025; Cemri et al., 2025). Alternatively, domain-specific fine-tuning yields specialized medical LLMs such as HUATUOGPT (Chen et al., 2024c) and MEDLLAMA (Qiu et al., 2024), but these models often overfit their training distributions and exhibit limited adaptability to unseen clinical scenarios (Yang et al., 2024b).

Overall, prior work has made progress in reasoning, retrieval, and domain adaptation, yet these components are largely developed in isolation and optimized for different assumptions. In practice, improving any single dimension alone is insufficient: stronger prompting often amplifies unstable or redundant reasoning, retrieval without role differentiation fails to induce complementary perspectives, and multi-agent collaboration without verification can exacerbate inconsistency rather than mitigate it. These limitations are particularly pronounced in medical question answering, where distribution shifts, incomplete knowledge, and evolving clinical standards are the norm, while model fine-tuning or parameter updates are often infeasible at deployment time. Consequently, robust MedQA requires a unified test-time framework that can simultaneously coordinate complementary reasoning roles, ground inference in reliable external evidence, and explicitly assess answer reliability within a single inference pipeline.

To address these challenges, we propose **TAGS** (Test-time Generalist-Specialist framework with retrieval-augmented reasoning and verification), a modular and inference-only framework for robust MedQA. TAGS introduces a structured collaboration between a generalist agent, which provides

broad clinical reasoning, and a specialist agent, which contributes domain-focused expertise. Their complementary reasoning is supported by hierarchical retrieval augmentation that supplies both semantically relevant and rationale-aligned exemplars, and finalized by an uncertainty-aware answer aggregation module that verifies reasoning consistency and selects reliable predictions.

Extensive experiments on nine MedQA benchmarks demonstrate that TAGS consistently outperforms strong prompting-based and multi-agent baselines across multiple foundation models, including GPT-4o, DeepSeek-R1, and Qwen-2.5-7B. Notably, TAGS remains effective even when excluding semantically nearest exemplars, indicating that its gains primarily stem from exposure to diverse and reliable reasoning patterns rather than answer memorization.

2 Related Work

2.1 Medical Question Answering

Medical question answering (MedQA) aims to predict accurate answers to domain-specific clinical or biomedical questions, often posed in a multiple-choice format. Existing benchmarks span a variety of formats and reasoning challenges, including clinical exam-style questions, evidence-based inference, and multi-subject distractor-rich scenarios (Pal et al., 2022). Recent approaches leverage large language models (LLMs) or chain-of-thought prompting to enhance reasoning (Singhal et al., 2025). MedCoT (Liu et al., 2024) explicitly integrates multi-step rationale generation with hierarchical expert feedback. Concurrently, biomedical LLMs such as MedLLaMA[†], HuatuoGPT (Chen et al., 2024c), and OpenBioLLM (Pal and Sankarassubbu, 2024) have achieved strong zero-shot or few-shot performance on MedQA benchmarks. However, these models typically rely on direct answer generation and lack explicit mechanisms for multi-agent reasoning or consistency verification. In contrast, our method introduces a retrieval-augmented multi-agent framework that performs staged reasoning and employs a dedicated verifier to assess the reliability of generated answers, promoting robustness and interpretability.

2.2 Retrieval-Augmented Reasoning

Retrieval-augmented reasoning enhances prediction quality, factual consistency, and interpretability

by incorporating external knowledge into the reasoning process. Early works such as RAG (Lewis et al., 2020) retrieve text passages to guide open-domain generation, while later methods extend retrieval to more structured forms, such as few-shot demonstrations (Izacard et al., 2023) or intermediate reasoning paths (Shi et al., 2023). In the context of chain-of-thought (CoT) prompting, retrieval has been explored to select relevant questions, rationales, or multi-step reasoning exemplars that better align with the target task (Xu et al., 2022; He et al., 2025). Despite these advances, most existing frameworks rely primarily on question-level similarity and often neglect deeper alignment at the reasoning level. Our framework addresses this limitation by employing a hierarchical retrieval strategy: first retrieving question-option exemplars, then refining based on CoT similarity. This enables alignment in both problem context and reasoning structure, thereby improving downstream multi-agent reasoning.

2.3 Multi-Agent Systems for Reasoning

Multi-agent systems (MAS) have emerged as a promising approach to enhance the robustness, diversity, and reliability of reasoning in complex tasks, including medical question answering. By orchestrating multiple reasoning paths or personas, MAS frameworks aim to mitigate biases and capture complementary perspectives, which are particularly critical in high-stakes medical decision-making. Recent studies have explored various MAS paradigms for medical reasoning. MedAgents (Tang et al., 2023) proposes a collaborative multi-agent framework where multiple agents independently generate answers and a majority voting scheme determines the final prediction. MDAgents (Kim et al., 2024a) further enhances this idea by introducing dynamic collaboration and adaptive feedback mechanisms among agents during the reasoning process. MedPrompt (Chen et al., 2024d) adopts a multi-round prompting strategy combined with ensemble voting to improve medical QA performance. Additionally, frameworks like Multi-Persona (Wang et al., 2023) and Self-Refine (Madaan et al., 2024) leverage self-collaboration and iterative self-feedback to strengthen individual agent reasoning capabilities. While multi-agent collaboration has demonstrated effectiveness in improving answer quality, it also introduces notable challenges. As highlighted in recent evaluations (Tang et al., 2023), excessive

[†]<https://huggingface.co/johnsnowlabs/>

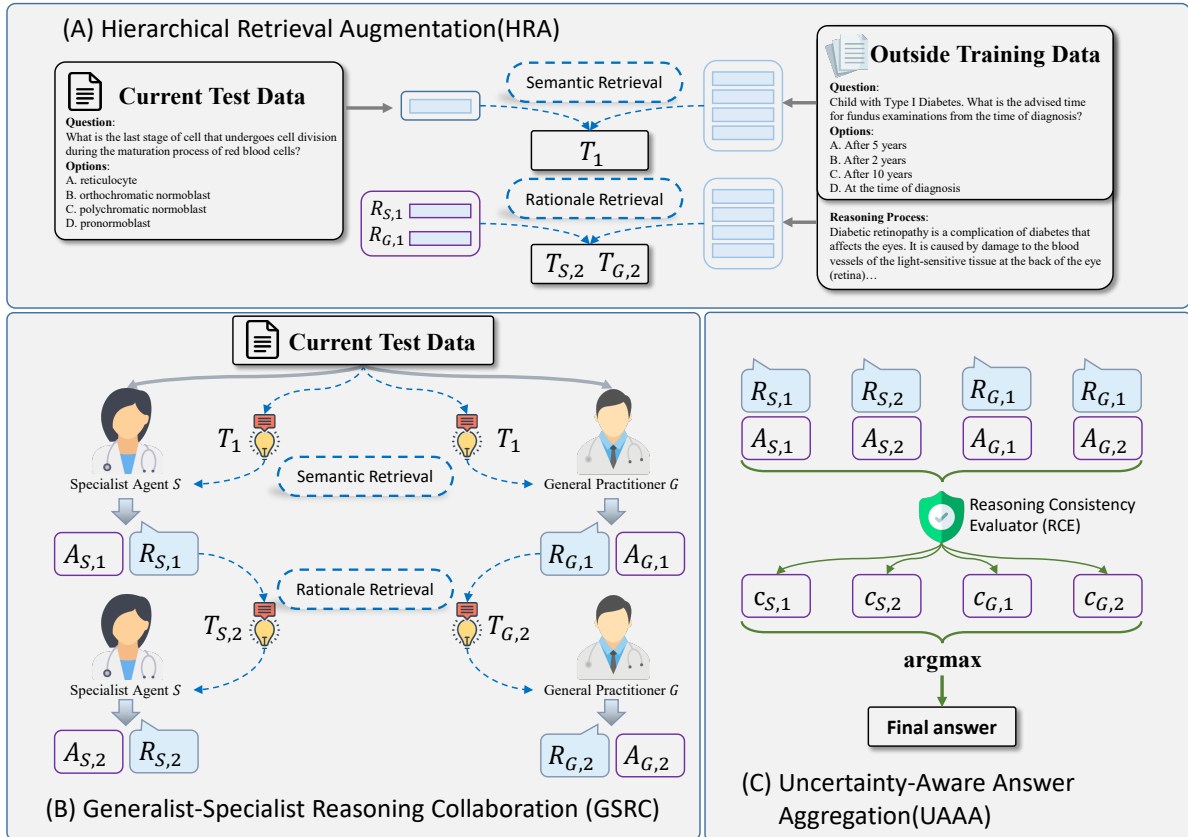


Figure 1: Overview of the proposed TAGS framework. The architecture consists of three modules: (A) **HRA (Hierarchical Retrieval Augmentation)**, a two-stage retrieval process that selects semantically relevant exemplars (T_1) and refines them based on rationale alignment ($T_{G,2}$, $T_{S,2}$). (B) **GSRC (Generalist-Specialist Reasoning Collaboration)** employs dual-agent reasoning across two rounds, generating four candidate (Rationale, Answer) pairs. (C) **UAAA (Uncertainty-Aware Answer Aggregation)** assesses rationale consistency using the RCE and aggregates reliability scores (c) to determine the final answer.

183 agent interactions may lead to reasoning conflicts, 184 unstable decision paths, and increased inference 185 costs. Recent studies (Cemri et al., 2025) further 186 reveal that over-complex MAS architectures often 187 suffer from systemic failures, such as miscommu- 188 nication, vague role specification, and weak verifi- 189 cation. Moreover, most existing MAS frameworks 190 lack explicit mechanisms to assess the internal con- 191 sistency between generated reasoning and final an- 192 swers, which can limit reliability in clinical con- 193 texts. To address these limitations, we propose a 194 lightweight General-Specialist Reasoning Collabo- 195 ration (GSRC) strategy that pairs a generalist and a 196 specialist agent in a complementary manner, pro- 197 moting stable and robust medical reasoning with 198 minimal inter-agent conflicts.

199 3 Methodology

200 We propose TAGS (Test-time General- 201 ist-Specialist Reasoning with Retrieval-

Augmentation and Uncertainty-Aware Verifi- 202 cation), a parameter-efficient framework for 203 medical question answering that operates entirely 204 during inference. At its core is the General- 205 ist-Specialist Reasoning Collaboration (GSRC), 206 a dual-agent design that promotes reasoning 207 diversity and domain alignment without requiring 208 any parameter updates. To support GSRC, we 209 introduce two auxiliary modules: Hierarchical 210 Retrieval Augmentation (HRA), which supplies 211 diverse and rationale-aligned exemplars, and 212 Uncertainty-Aware Answer Aggregation (UAAA), 213 which selects the final answer by evaluating the 214 consistency of each reasoning path. As shown in 215 Figure 1, these components form an integrated 216 pipeline that enables robust, zero-shot clinical QA 217 without model fine-tuning. 218

3.1 Hierarchical Retrieval Augmentation

Hierarchical Retrieval Augmentation (HRA) grounds reasoning in an external reference corpus while injecting diverse paths for chain-of-thought (CoT) generation through a two-stage retrieval scheme. We retrieve from a frozen medical-QA corpus \mathcal{D} whose entries are $d_i = (Q_i, O_i, A_i, R_i)$, where R_i denotes the CoT rationale. We use a frozen text encoder $\mathcal{E}(\cdot)$ based on M3-Embedding (Chen et al., 2024b), with 1024-dimensional output.

Stage 1: Initial semantic retrieval. We begin by embedding the query using a frozen encoder. Let $\mathbf{z} = \mathcal{E}(Q \oplus O)$, where the question Q and its options O are concatenated in standard order (A, B, C, D). Cosine similarity is computed against all candidate embeddings $\mathcal{E}(Q_i \oplus O_i)$ in the corpus. The top- K retrieved examples form:

$$\mathcal{T}_1 = \text{Top-K}\{d_i \in \mathcal{D} : \text{sim}(\mathbf{z}, \mathcal{E}(Q_i \oplus O_i))\}. \quad (1)$$

Stage 2: Rationale-guided retrieval. After Round-1 reasoning yields preliminary rationales $R_{G,1}$ and $R_{S,1}$, we retrieve exemplars whose stored rationales best match these CoTs:

$$\mathcal{T}_{G,2} = \text{Top-K}\{d_i \in \mathcal{D} : \text{sim}(\mathbf{r}_G, \mathcal{E}(R_i))\}, \quad (2)$$

$$\mathcal{T}_{S,2} = \text{Top-K}\{d_i \in \mathcal{D} : \text{sim}(\mathbf{r}_S, \mathcal{E}(R_i))\}. \quad (3)$$

By aligning on reasoning paths rather than surface form, Stage 2 injects complementary evidence beyond surface similarity, reducing the limitations of purely semantic matching.

3.2 Generalist–Specialist Collaboration

Given the retrieved exemplar sets from HRA, Generalist–Specialist Reasoning Collaboration (GSRC) performs dual-agent inference in two rounds by coupling broad medical knowledge with focused domain expertise. The system consists of a generalist agent \mathcal{G} and a specialist agent \mathcal{S} , both instantiated as prompted roles of the same frozen LLM without parameter updates.

An auxiliary LLM role first infers the medical specialty most relevant to the query (Q, O) , yielding a label s (e.g., cardiology). This label is then injected into the prompt for \mathcal{S} as “You are a medical specialist in the field of [s]”, guiding its reasoning toward domain-specific knowledge while preserving the core semantics of (Q, O) . Further details are provided in Appendix G. The collaboration unfolds in two rounds that iteratively refine rationales and answers.

Round 1: Initial hypothesis generation. Both agents receive the query (Q, O) together with the semantically retrieved set \mathcal{T}_1 (§3.1); the specialist additionally sees the inferred specialty s . Each agent produces an initial CoT and answer:

$$\begin{aligned} (R_{G,1}, A_{G,1}) &= \mathcal{G}(Q, O, \mathcal{T}_1), \\ (R_{S,1}, A_{S,1}) &= \mathcal{S}(Q, O, \mathcal{T}_1, s). \end{aligned} \quad (4)$$

These preliminary CoTs trigger Stage 2 of HRA, which returns the rationale-aligned exemplar sets $\mathcal{T}_{G,2}$ and $\mathcal{T}_{S,2}$.

Round 2: Refined reasoning with aligned exemplars. Using the tailored sets, the agents generate updated rationales and answers:

$$\begin{aligned} (R_{G,2}, A_{G,2}) &= \mathcal{G}(Q, O, \mathcal{T}_{G,2}), \\ (R_{S,2}, A_{S,2}) &= \mathcal{S}(Q, O, \mathcal{T}_{S,2}, s). \end{aligned} \quad (5)$$

Finally, the four (*rationale, answer*) pairs are gathered into a candidate set

$$\mathcal{C} = \{(R_{k,r}, A_{k,r}) \mid k \in \{G, S\}, r \in \{1, 2\}\}, \quad (6)$$

which is then forwarded to the Uncertainty-Aware Answer Aggregation module (§3.3) for scoring and final selection.

3.3 Uncertainty-Aware Answer Aggregation

Uncertainty-Aware Answer Aggregation (UAAA) takes as input the candidate set \mathcal{C} generated by GSRC (§3.2) and selects a single high-confidence answer through consistency-based scoring. To accomplish this, we define a *Reasoning Consistency Evaluator* (RCE), implemented as a separate zero-shot role of the same frozen LLM.

Given a candidate pair (R_k, A_k) , the RCE assesses how well the rationale supports the answer in the context of the original query (Q, O) , and assigns an integer score $c_k \in [0, 5]$, where higher values indicate stronger logical and clinical coherence. The scoring rubric is detailed in Appendix H. The final answer is selected as:

$$A_{\text{final}} = A_{k^*}, \quad k^* = \arg \max_{k \in \mathcal{C}} c_k. \quad (7)$$

In the case of ties, preference is resolved deterministically in the following order: specialist round 2, generalist round 2, specialist round 1, and generalist round 1. By explicitly verifying the internal consistency of each reasoning path, UAAA mitigates hallucination propagation and stabilizes final predictions, all without any parameter updates.

Method	MedQA	PubMedQA	MedMCQA	MedBullets	MMLU	MMLU-Pro	MedExQA	MedXpert-R	MedXpert-U	Average
GPT-4o	32.0	9.0	25.0	19.1	24.7	21.0	18.0	7.0	6.0	18.0
+ few-shot	31.0	16.0	34.0	16.9	32.9	27.0	17.0	8.0	11.0	21.7
+ RAG	42.0	12.0	30.0	22.5	20.5	37.0	15.0	19.0	10.0	23.1
+ CoT	39.0	10.0	30.0	28.1	26.0	35.0	24.0	12.0	15.0	24.3
+ CoT-SC	37.0	6.0	35.0	30.3	30.1	43.0	22.0	10.0	14.0	25.3
+ Multi-Persona	45.0	15.0	25.0	29.2	37.0	42.0	21.0	10.0	16.0	26.7
+ Self-Refine	41.0	13.0	34.0	28.1	34.2	34.0	22.0	17.0	19.0	26.9
+ MedAgents	43.0	15.0	30.0	27.0	28.8	8.0	19.0	3.0	6.0	20.0
+ MDAgents	36.0	11.0	22.0	21.3	24.7	8.0	13.0	4.0	5.0	16.1
+ MedPrompt	34.0	11.0	26.0	22.5	26.0	22.0	16.0	14.0	9.0	20.1
+ Ours	54.0	13.0	32.0	33.7	45.2	47.0	17.0	22.0	22.0	31.8

Table 1: Performance heatmap by methods and datasets. All tasks are evaluated on the HARD set with Pass@1 Accuracy (%) using GPT-4o base model.

Method	MedQA	PubMedQA	MedMCQA	MedBullets	MMLU	MMLU-Pro	MedExQA	MedXpert-R	MedXpert-U	Average
DeepSeek-R1	38.0	11.0	28.0	36.0	32.9	36.0	20.0	20.0	23.0	27.2
+ few-shot	27.0	12.0	32.0	33.7	35.6	41.0	27.0	11.0	9.0	25.4
+ RAG	49.0	20.0	31.0	43.8	53.8	42.0	25.0	28.0	26.0	35.4
+ CoT	47.0	12.0	31.0	39.3	38.4	35.0	22.0	27.0	27.0	31.0
+ CoT-SC	52.0	14.0	32.0	43.8	45.2	38.0	24.0	17.0	26.0	32.4
+ Multi-Persona	52.0	18.0	37.0	42.7	42.5	38.0	26.0	23.0	26.0	33.9
+ Self-Refine	33.0	17.0	30.0	34.8	27.4	22.0	24.0	12.0	13.0	23.7
+ MedAgents	48.0	21.0	22.0	44.9	43.8	35.0	27.0	22.0	25.0	32.1
+ MedPrompt	46.0	14.0	30.0	38.2	45.2	27.0	24.0	8.0	7.0	26.6
+ Ours	55.0	28.0	35.0	52.8	61.6	53.0	26.0	36.0	49.0	44.0

Table 2: Performance heatmap by methods and datasets. All tasks are evaluated on the HARD set with Pass@1 Accuracy (%) using DeepSeek-R1 base model.

4 Experiments

4.1 Experimental Setting

Retrieval Dataset. We use the MedReason dataset (Wu et al., 2025) as our external retrieval corpus. It contains 32,682 medical QA pairs with clinically validated, step-by-step explanations generated via a knowledge graph-guided pipeline based on PrimeKG (Chandak et al., 2023). Unlike general CoT datasets, MedReason ensures factual correctness by filtering out chains that do not lead to the correct answer. We treat it as a structured knowledge base for retrieving semantically or logically relevant examples at inference. Dataset construction details and examples are provided in Appendix I.

Note that the retrieval dataset is constructed independently from all evaluation benchmarks and does not overlap with any of the test datasets used in this work. To further prevent near-duplicate leakage, we additionally filter the retrieval pool by removing candidates whose semantic similarity to the query exceeds a threshold of 0.9, ensuring that highly similar or paraphrased instances are excluded.

Test Datasets. We evaluate TAGS on a curated benchmark of nine medical QA datasets selected from the MEDAGENTS BENCH framework (Tang et al., 2025), designed to assess com-

plex clinical reasoning. The benchmark includes challenging subsets from: **MedQA** (Jin et al., 2021), a multilingual board-exam dataset (e.g., USMLE); **PubMedQA** (Jin et al., 2019), derived from biomedical literature with yes/no/maybe answers; **MedMCQA** (Pal et al., 2022), covering 21 medical subjects from Indian medical exams; **MedBullets** (Chen et al., 2024a), featuring long-context clinical questions; **MedExQA** (Kim et al., 2024b), emphasizing explainable QA across five specialties; **MedXpertQA** (Zuo et al., 2025), with subsets targeting reasoning and understanding; **MMLU** (Hendrycks et al., 2020) and **MMLU-Pro** (Wang et al., 2024), general benchmarks with medical subfields.

To better reflect real-world difficulty, we follow the hard subset construction pipeline proposed by MEDAGENTS BENCH. Questions are selected based on model failure rates (<50% accuracy across a set of strong models), medical topic coverage, and reasoning depth. Specifically, we include 100 hard questions each from MedQA, PubMedQA, MedMCQA, MedExQA, and MMLU-Pro; 100 from each MedXpertQA subset (Reasoning and Understanding); 89 from MedBullets; and 73 from MMLU. This results in a total of 862 expert-verified instances designed to stress-test the reasoning capabilities of large language models.

We additionally evaluate three representative datasets on their full test sets and observe consistent gains from TAGS (Appendix D).

Baselines. We first compare our method against several widely adopted prompting and reasoning strategies that do not involve model updates: (1) **CoT (Chain-of-Thought)**(Wei et al., 2022): A prompting technique that guides the model to articulate intermediate reasoning steps before producing a final answer. (2) **CoT-SC (Chain-of-Thought with Self-Consistency)**(Wang et al., 2022): An extension of CoT that generates multiple reasoning paths and selects the most consistent answer via majority voting. (3) **Multi-Persona**(Wang et al., 2023): A method that simulates multiple expert personas to collaboratively reason through clinical questions. (4) **Self-Refine**(Madaan et al., 2024): A self-improvement framework in which the model iteratively refines its own responses across multiple reasoning stages. (5) **MedAgents**(Tang et al., 2023): A domain-specific multi-agent framework that employs multiple specialist agents for collaborative clinical reasoning. (6) **MDAgents**(Kim et al., 2024a): A lightweight variant of MedAgents that combines minimal agent collaboration with retrieval augmentation to improve reasoning. (7) **MedPrompt** (Chen et al., 2024d): A retrieval-augmented prompting strategy that integrates semantically similar historical cases to enhance clinical inference. Note that all baselines are implemented following their original papers or publicly released prompts to ensure a fair and realistic evaluation setting.

We additionally report results under a **few-shot** baseline, where five training examples from the target dataset are retrieved and used as in-context demonstrations for single-pass inference. We also include a **RAG** baseline, which retrieves the top- K most semantically similar questions with accompanying CoTs from the MedReason dataset and feeds them directly to the model. This RAG setting shares the same retrieval setup but excludes agent collaboration and verification, highlighting the value of structured reasoning.

We further evaluate our method against several strong open-source foundation models and their medically adapted variants: (1) **Qwen2.5-7B** (Yang et al., 2024a): A 7B general-purpose model instruction-tuned for diverse tasks, evaluated both with and without CoT prompting. (2) **LLaMA-3-8B** (Grattafiori et al., 2024): Meta’s latest 8B instruction-tuned model with improved rea-

soning capabilities. (3) **HuatuoGPT-o1-7B** (Chen et al., 2024c): A 7B model fine-tuned for complex medical reasoning via reinforcement learning. (4) **HuatuoGPT-o1-8B** (Chen et al., 2024c): An enhanced 8B version of HuatuoGPT, optimized for clinical inference tasks. (5) **MedLLaMA-3-8B-v1.0** (Qiu et al., 2024): A medical-adapted variant of LLaMA-3 trained on biomedical corpora. (6) **MedLLaMA-3-8B-v2.0**: An updated release with improved performance on expert-level medical benchmarks. (7) **OpenBioLLM-8B** (Pal and Sankarasubbu, 2024): An open-source 8B biomedical language model fine-tuned for healthcare and life sciences applications.

Evaluation Metrics Following (Tang et al., 2023), we report **Pass@1 Accuracy** as the evaluation metric, which measures whether the model’s first generated answer exactly matches the ground-truth answer.

Reproducibility. All experiments were conducted using Python 3.10 and PyTorch 2.4.0 on four NVIDIA H100 GPUs, each with 80 GB of memory. For proprietary LLM baselines such as GPT-4o and DeepSeek-R1, we accessed the models through their official APIs and ensured consistent use of the same model version across all runs. For open-source models, we directly loaded checkpoint weights from their respective official Hugging Face repositories to ensure reproducibility and transparency.

4.2 Compared with Prompting and MAS

We evaluate the effectiveness of TAGS by comparing it with a diverse set of prompting-based and multi-agent reasoning baselines across nine challenging MedQA benchmarks. Tables 1 and 2 summarize the results in terms of **Pass@1 Accuracy**, evaluated on the HARD split using two foundational LLMs: GPT-4o and DeepSeek-R1. Under the GPT-4o setting, TAGS achieves the highest average accuracy of 31.8%, outperforming all baselines including Self-Refine (26.9%), MedAgents (20.0%), and CoT-SC (25.3%). The most notable improvements appear on MedQA (+9.0 over Multi-Persona), MMLU (+8.2 over Multi-Persona), and MedXpert-R (+3.0 over RAG), highlighting the impact of verifier-guided aggregation and structured multi-agent reasoning. TAGS also surpasses standard few-shot and RAG baselines by margins of +10.1 and +8.7 respectively. With the DeepSeek-R1 base model, TAGS achieves an average accuracy of 44.0%, outperforming CoT-SC (32.4%), Multi-

Method	MedQA	PubMedQA	MedMCQA	MedBullets	MMLU	MMLU-Pro	MedExQA	MedXpert-R	MedXpert-U	Average
Qwen2.5-7B	16.0	16.0	24.0	4.5	13.7	26.0	9.0	10.0	8.0	14.1
Llama-3-8B	18.0	13.0	23.0	16.9	11.0	23.0	11.0	10.0	4.0	14.4
HuatuoGPT-o1-7B	22.0	21.0	26.0	12.4	13.7	29.0	9.0	11.0	7.0	16.8
HuatuoGPT-o1-8B	29.0	20.0	33.0	20.2	21.9	17.0	18.0	16.0	7.0	20.2
MedLlama-3-8B-v1.0	24.0	20.0	22.0	14.6	16.4	12.0	12.0	11.0	11.0	15.9
MedLlama-3-8B-v2.0	28.0	25.0	22.0	22.5	32.9	12.0	22.0	9.0	8.0	20.2
OpenBioLLM-8B	19.0	29.0	20.0	19.1	21.9	2.0	17.0	7.0	3.0	15.3
Qwen2.5-7B + Ours	28.0	25.0	24.0	14.6	35.6	25.0	16.0	18.0	29.0	23.9

Table 3: Comparison with fine-tuned medical LLMs on nine MedQA benchmarks.

Persona (33.9%), and MedAgents (32.1%). TAGS also surpasses the few-shot and RAG baselines by margins of +18.6 and +8.5, respectively, demonstrating the scalability of our framework across both general and domain-specific tasks.

4.3 Compared with Fine-Tuned LLMs

To further contextualize the performance of our TAGS framework, we evaluate its effectiveness when integrated with the Qwen2.5-7B base model and compare its performance against a series of prominent open-source and medically fine-tuned large language models across the same nine challenging MedQA datasets. The results of this comparison are presented in Table 3. As shown in the table, our TAGS framework substantially boosts the zero-shot question answering capability of the base Qwen2.5-7B model, improving its average accuracy from 14.1% to 23.9%. In particular, TAGS demonstrates robust performance gains on difficult benchmarks such as MedQA (+12.0 percentage points), MMLU (+21.9 percentage points), and MedXpert-U (+21.0 percentage points). Notably, our inference-only strategy even outperforms several models that have been fine-tuned with domain-specific medical corpora or expert feedback, such as MedLLaMA-3-8B and the HuatuoGPT-o1 variants, on the majority of the evaluated datasets. These results strongly highlight the significant potential of structured retrieval and multi-agent reasoning, combined with uncertainty-aware verification, to effectively close the performance gap with models requiring extensive fine-tuning, while retaining the inherent flexibility and adaptability of a zero-shot approach.

4.4 Ablation Study

We conduct ablation studies on Qwen2.5-7B to assess the contribution of each component in TAGS, focusing on MMLU and MedXpert-U.

Overall impact of individual components. As shown in Table 4, the base Qwen2.5-7B model achieves 13.7% and 8.0% accuracy on MMLU and MedXpert-U, respectively. Adding retrieval augmentation (RAG) improves performance to 20.5% and 10.0%, demonstrating the benefit of external knowledge access. Introducing generalist (G) and specialist (S) agents with majority voting further boosts MMLU to 30.1%, while yielding limited gains on MedXpert-U, suggesting that agent collaboration alone is insufficient for highly specialized domains.

Effect of hierarchical retrieval augmentation (HRA). Incorporating HRA leads to substantial improvements for both agents, achieving 34.2% for the generalist and 32.9% for the specialist. This highlights the importance of reasoning-guided retrieval beyond purely semantic matching. The complete TAGS framework, which further integrates uncertainty-aware answer aggregation (UAAA), attains the best results of 35.6% on MMLU and 29.0% on MedXpert-U, confirming the synergistic effect of structured retrieval, dual-agent reasoning, and verification.

Robustness to retrieval overlap and specialist mismatch. To examine whether TAGS relies on retrieving semantically closest examples, we evaluate **RAG-w/o-topk**, which excludes the top-10 most similar candidates during retrieval. Performance decreases only marginally, indicating that TAGS benefits primarily from exposure to valid reasoning patterns rather than copying specific answers. We further consider **S-w-3rd**, which assigns the third-ranked specialist instead of the top choice. The modest performance drop under this perturbation suggests that TAGS remains effective even with imperfect specialist selection.

Robustness to retrieval corpus choice. We additionally evaluate TAGS using m1k (Huang et al., 2025), a substantially smaller CoT reference set

RAG	\mathcal{G}	\mathcal{S}	HRA	UAAA	MMLU	MedXpert-U
					13.7	8.0
✓					20.5	10.0
✓	✓	✓			30.1	10.0
✓	✓		✓	✓	34.2	16.0
✓		✓	✓	✓	32.9	18.0
✓	✓	✓	✓		31.5	22.0
✓	✓	✓	✓	✓	35.6	29.0
w/o-top10	✓	✓	✓	✓	37.0	24.0
✓	✓	w-3rd	✓	✓	34.2	23.0

Table 4: Ablation Study on TAGS using Qwen2.5-7B. RAG: Retrieval-Augmented Generation; \mathcal{G} : Generalist Agent; \mathcal{S} : Specialist Agent; HRA: Hierarchical Retrieval Augmentation; UAAA: Uncertainty-Aware Answer Aggregation.

containing 1,000 exemplars. Results reported in Appendix C show that TAGS maintains strong performance, indicating that it is not tightly coupled to a specific retrieval corpus.

Finally, we observe that the specialist agent consistently produces more domain-specific rationales than the generalist, supported by both qualitative examples and quantitative analyses in Appendix E.

4.5 Hyperparameter Analysis

Figure 2 shows TAGS’ sensitivity to two key hyperparameters: the number of specialist agents and the retrieval size K . In Figure 2(a), adding one specialist to the generalist improves accuracy from 34.2% to 45.2% on MMLU and from 16.0% to 22.0% on MedXpert-U. However, adding more specialists brings limited or no further gains, likely due to redundancy or conflicts in reasoning paths. Figure 2(b) shows that accuracy peaks at $K = 2$ and declines with larger K , as additional exemplars may introduce noise or irrelevant content that misleads the model. These results support our choice of using one specialist and $K = 2$ as the default configuration, balancing diversity and robustness.

4.6 Inference Efficiency

On the MedQA dataset with GPT-4o, TAGS takes 72 seconds per question on average, which is longer than CoT-SC (27.7s) but shorter than Multi-Persona (109.6s). Although slower than simple prompting, TAGS achieves substantially higher accuracy. Both reasoning and verification are parallelizable, enabling efficient deployment in real-world clinical settings. This moderate inference cost represents a favorable trade-off for improved robustness and reliability. Additionally, these stages can be parallelized across GPU streams or executed via asyn-

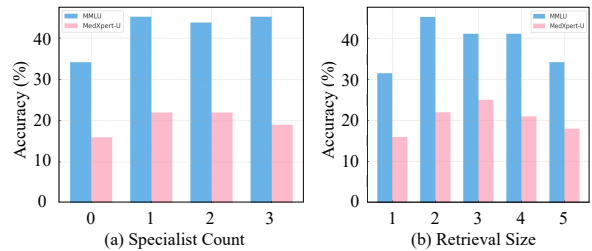


Figure 2: Hyper-parameter sensitivity analysis of specialist count and retrieval size in relation to accuracy.

chronous API calls to further speed up inference.

5 Conclusion

We presented TAGS, a parameter-efficient and test-time-only framework designed to enhance reliability in medical question answering without model fine-tuning. TAGS introduces a structured reasoning paradigm through generalist–specialist reasoning collaboration, which combines the breadth of a generalist with the depth of a specialist to generate complementary inference paths. This collaboration is guided by hierarchical retrieval augmentation, which retrieves exemplars at both semantic and rationale levels to enrich reasoning diversity, and finalized by uncertainty-aware answer aggregation to select robust answers. Extensive experiments on nine challenging MedQA benchmarks, spanning general-purpose and fine-tuned LLMs, demonstrate TAGS’ consistent superiority over prompting-based, retrieval-augmented, and multi-agent baselines. Notably, our method delivers substantial improvements even for compact 7B models, highlighting its adaptability across model scales. TAGS offers a practical, inference-only alternative for trustworthy medical AI and opens promising directions for adaptive retrieval, dynamic agent collaboration, and scaling to multimodal or real-world clinical QA workflows.

Limitations

While TAGS offers a robust, inference-only approach to medical QA, it carries several limitations. First, it depends heavily on the coverage and quality of the external retrieval corpus: gaps or biases in the QA database may lead to missing or misleading exemplars, particularly for rare diseases or newly emerging clinical scenarios. Second, the Reasoning Consistency Evaluator (RCE) is itself a zero-shot LLM prompt and may inherit the same hallucination tendencies or biases as the genera-

tor agents, potentially mis-scoring perfectly valid but unconventional reasoning chains. Third, the two-round retrieval and dual-agent design, while effective, substantially increases inference latency and API cost compared to single-pass prompting; this may limit real-time deployment in resource-constrained clinical settings.

Additionally, our current evaluation focuses solely on answer accuracy (Pass@1), without assessing the interpretability or faithfulness of reasoning paths. Future work may benefit from human evaluation or rationale consistency metrics to further assess clinical applicability.

Moreover, our specialty inference component can occasionally misclassify the most relevant domain, which, although gracefully handled, may still introduce suboptimal reasoning contexts. Finally, our evaluation is confined to English-language, multiple-choice benchmarks and does not cover open-ended clinical dialogs, multimodal data (e.g., images, lab reports), or non-English patient populations. Addressing these limitations will require enriching and updating the retrieval corpus, developing more calibrated or human-in-the-loop verifier mechanisms, optimizing retrieval budgets and round counts, and extending evaluation to diverse, real-world clinical workflows.

References

Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. 2025. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657*.

Payal Chandak, Kexin Huang, and Marinka Zitnik. 2023. Building a knowledge graph to enable precision medicine. *Scientific Data*, 10(1):67.

Hanjie Chen, Zhouxiang Fang, Yash Singla, and Mark Dredze. 2024a. Benchmarking large language models on answering and explaining challenging medical questions. *arXiv preprint arXiv:2402.18060*.

Jianlyu Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024b. M3-embedding: Multi-linguality, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2318–2335.

Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. 2024c. Huatuogpt-o1, towards medical complex reasoning with llms. *arXiv preprint arXiv:2412.18925*.

Qingyu Chen, Jingcheng Du, Yan Hu, Vipina Kut-tichi Keloth, Xueqing Peng, Kalpana Raja, Rui Zhang, Zhiyong Lu, and Hua Xu. 2023. Large language models in biomedical natural language processing: benchmarks, baselines, and recommendations. *arXiv e-prints*, pages arXiv–2305.

Xi Chen, Huahui Yi, Mingke You, WeiZhi Liu, Li Wang, Hairui Li, Xue Zhang, Yingman Guo, Lei Fan, Gang Chen, et al. 2025. Enhancing diagnostic capability with multi-agents conversational large language models. *NPJ digital medicine*, 8(1):159.

Xuhang Chen, Shenghong Luo, Chi-Man Pun, and Shuqiang Wang. 2024d. MedPrompt: Cross-modal prompting for multi-task medical image translation. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, pages 61–75. Springer.

Zhihao Fan, Lai Wei, Jialong Tang, Wei Chen, Wang Siyuan, Zhongyu Wei, and Fei Huang. 2025. Ai hospital: Benchmarking large language models in a multi-agent medical interaction simulator. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10183–10213.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Qiangqiang He, Shuwei Qian, Jie Zhang, and Chongjun Wang. 2025. Inference retrieval-augmented multimodal chain-of-thoughts reasoning for language models. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.

Xiaoke Huang, Juncheng Wu, Hui Liu, Xianfeng Tang, and Yuyin Zhou. 2025. m1: Unleash the potential of test-time scaling for medical reasoning with large language models. *arXiv preprint arXiv:2504.00869*.

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, 24(251):1–43.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? A large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. PubMedQA: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*.

727	Qiao Jin, Zheng Yuan, Guangzhi Xiong, Qianlan Yu,	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres,	783
728	Huaiyuan Ying, Chuanqi Tan, Mosha Chen, Song-	Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin	784
729	fang Huang, Xiaozhong Liu, and Sheng Yu. 2022.	Clark, Stephen R Pfohl, Heather Cole-Lewis, et al.	785
730	Biomedical question answering: a survey of ap-	2025. Toward expert-level medical question answer-	786
731	proaches and challenges. <i>ACM Computing Surveys</i>	ing with large language models. <i>Nature Medicine</i> ,	787
732	(<i>CSUR</i>), 55(2):1–36.	pages 1–8.	788
733	Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu	Xiangru Tang, Daniel Shao, Jiwoong Sohn, Jiapeng	789
734	Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee,	Chen, Jiayi Zhang, Jinyu Xiang, Fang Wu, Yilun	790
735	Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won	Zhao, Chenglin Wu, Wenqi Shi, et al. 2025. Meda-	791
736	Park. 2024a. MDAgents: An adaptive collaboration	gentsbench: Benchmarking thinking models and	792
737	of llms in medical decision making. <i>arXiv preprint</i>	agent frameworks for complex medical reasoning.	793
738	<i>arXiv:2404.15155</i> .	<i>arXiv preprint arXiv:2503.07459</i> .	794
739	Yunsoo Kim, Jinge Wu, Yusuf Abdulle, and Honghan	Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming	795
740	Wu. 2024b. MedExQA: Medical question answer-	Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and	796
741	ing benchmark with multiple explanations. <i>arXiv</i>	Mark Gerstein. 2023. MedAgents: Large language	797
742	<i>preprint arXiv:2406.06331</i> .	models as collaborators for zero-shot medical reason-	798
743	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	ing. <i>arXiv preprint arXiv:2311.10537</i> .	799
744	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	800
745	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	Ed H Chi, Sharan Narang, Aakanksha Chowdhery,	801
746	täschel, et al. 2020. Retrieval-augmented generation	and Denny Zhou. 2022. Self-consistency improves	802
747	for knowledge-intensive nlp tasks. <i>Advances in neu-</i>	chain of thought reasoning in language models. In	803
748	<i>ral information processing systems</i> , 33:9459–9474.	<i>The Eleventh International Conference on Learning</i>	804
749	Jiaxiang Liu, Yuan Wang, Jiawei Du, Joey Tianyi Zhou,	<i>Representations</i> .	805
750	and Zuozhu Liu. 2024. Medcot: Medical chain	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng	806
751	of thought via hierarchical expert. <i>arXiv preprint</i>	Ni, Abhranil Chandra, Shiguang Guo, Weiming	807
752	<i>arXiv:2412.13736</i> .	Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al.	808
753	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	2024. MMLU-Pro: A more robust and challenging	809
754	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	multi-task language understanding benchmark. <i>arXiv</i>	810
755	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	<i>preprint arXiv:2406.01574</i> .	811
756	et al. 2024. Self-Refine: Iterative refinement with	Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao	812
757	self-feedback. <i>Advances in Neural Information Pro-</i>	Ge, Furu Wei, and Heng Ji. 2023. Unleashing the	813
758	<i>cessing Systems</i> , 36.	emergent cognitive synergy in large language mod-	814
759	Ankit Pal, Logesh Kumar Umapathi, and Malaikan-	els: A task-solving agent through multi-persona self-	815
760	nan Sankarasubbu. 2022. MedMCQA: A large-scale	collaboration. <i>arXiv preprint arXiv:2307.05300</i> .	816
761	multi-subject multi-choice dataset for medical do-	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	817
762	main question answering. In <i>Conference on health,</i>	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	818
763	<i>inference, and learning</i> , pages 248–260. PMLR.	et al. 2022. Chain-of-thought prompting elicits rea-	819
764	Malaikannan Sankarasubbu Ankit Pal and Malaikannan	soning in large language models. <i>Advances in Neural</i>	820
765	Sankarasubbu. 2024. Openbiollms: Advancing open-	<i>Information Processing Systems</i> , 35:24824–24837.	821
766	source large language models for healthcare and life	Juncheng Wu, Wenlong Deng, Xingxuan Li, Sheng Liu,	822
767	sciences. <i>Hugging Face repository</i> .	Taomian Mi, Yifan Peng, Ziyang Xu, Yi Liu, Hyunjin	823
768	Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong	Cho, Chang-In Choi, et al. 2025. Medreason: Elicit-	824
769	Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and	ing factual medical reasoning steps in llms via knowl-	825
770	Weidi Xie. 2024. Towards building multilingual lan-	edge graphs. <i>arXiv preprint arXiv:2504.00993</i> .	826
771	guage model for medicine. <i>Nature Communications</i> ,	Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu,	827
772	15(1):8384.	Ding Zhao, Joshua Tenenbaum, and Chuang Gan.	828
773	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	2022. Prompting decision transformer for few-shot	829
774	joon Seo, Rich James, Mike Lewis, Luke Zettle-	policy generalization. In <i>international conference on</i>	830
775	moyer, and Wen-tau Yih. 2023. Replug: Retrieval-	<i>machine learning</i> , pages 24631–24645. PMLR.	831
776	augmented black-box language models. <i>arXiv</i>	Shaochen Xu, Yifan Zhou, Zhengliang Liu, Zihao Wu,	832
777	<i>preprint arXiv:2301.12652</i> .	Tianyang Zhong, Huaqin Zhao, Yiwei Li, Hanqi	833
778	Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Haotian	Jiang, Yi Pan, Junhao Chen, et al. 2024. Towards	834
779	Sun, Hang Wu, Carl Yang, and May D Wang. 2024.	next-generation medical agent: How o1 is reshaping	835
780	Medadapter: Efficient test-time adaptation of large	decision-making in medical scenarios. <i>arXiv preprint</i>	836
781	language models towards medical reasoning. <i>arXiv</i>	<i>arXiv:2411.14461</i> .	837
782	<i>preprint arXiv:2405.03000</i> .		

838 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,
839 Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,
840 Fei Huang, Haoran Wei, et al. 2024a. Qwen2. 5
841 technical report. *arXiv preprint arXiv:2412.15115*.

842 Zhaorui Yang, Tianyu Pang, Haozhe Feng, Han Wang,
843 Wei Chen, Minfeng Zhu, and Qian Liu. 2024b.
844 Self-distillation bridges distribution gap in language
845 model fine-tuning. In *Proceedings of the 62nd Annual
846 Meeting of the Association for Computational
847 Linguistics*, pages 1028–1043.

848 Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu Zou,
849 Jinfa Huang, Jing Wu, Yiru Li, Sam S Chen, Peilin
850 Zhou, Junling Liu, et al. 2023. A survey of large
851 language models in medicine: Progress, application,
852 and challenge. *arXiv preprint arXiv:2311.05112*.

853 Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai
854 Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and
855 Bowen Zhou. 2025. MedXpertQA: Benchmarking
856 expert-level medical reasoning and understanding.
857 *arXiv preprint arXiv:2501.18362*.

A Ethics Statement 858

859 This work relies solely on publicly available medi-
860 cal question answering datasets, including MedQA,
861 PubMedQA, MedMCQA, and others curated in
862 the MEDAGENTS BENCH framework (Tang et al.,
863 2025). These datasets are de-identified and col-
864 lected from open educational or biomedical sources
865 such as medical board exams and peer-reviewed
866 literature. No private health records or patient-
867 identifiable information were used.

868 Our proposed framework operates entirely at test
869 time and does not require any model fine-tuning or
870 user data collection. All evaluations are conducted
871 offline on benchmark datasets, and no deployment
872 in real clinical settings has been performed.

873 While our method is designed to improve the ro-
874 bustness and reliability of medical LLMs, it is not
875 intended for use in high-stakes clinical decision-
876 making without appropriate human oversight. We
877 emphasize that the generated answers should not
878 be interpreted as medical advice. Future work may
879 involve incorporating human-in-the-loop mecha-
880 nisms and broader impact assessments before real-
881 world deployment.

B Ablation Study on DeepSeek-R1 Backbone 882

883 Table 5 reports ablation results on the DeepSeek-R1
884 backbone. We observe consistent and monotonic
885 improvements as GSRC, HRA, and UAAA are pro-
886 gressively added, confirming that each component
887 contributes positively and that the overall TAGS
888 framework generalizes beyond the backbone used
889 in the main paper. 890

Table 5: Ablation Study on DeepSeek-R1 Backbone

Model Variant	MMLU	MedXpert-U
Base (DeepSeek-R1)	32.9%	23.0%
+ GSRC	42.5%	29.0%
+ HRA	52.1%	40.0%
+ UAAA (Full TAGS)	61.6%	49.0%

C Alternative 1K Reference Corpus 891

892 As shown in Table 4 and §4.4, we conduct a per-
893 turbation experiment (“RAG-w/o-top10”) where
894 the 10 most semantically similar exemplars are re-
895 moved from the MedReason retrieval pool. To fur-
896 ther demonstrate generalizability beyond MedRea-
897 son, we evaluate TAGS using an independent and

much smaller CoT corpus: the m1k subset (1K examples), Unleash the Potential of Test-Time Scaling for Medical Reasoning with LLMs (Huang et al., 2025). This corpus is constructed via stratified sampling and expert filtering from MedQA, PubMedQA, HeadQA, and MedMCQA, and contains only 1,000 QA–rationale exemplars, compared to 32,000 in MedReason. Even under this limited-resource setting, TAGS achieves 47.9% on MMLU and 30.0% on MedXpert-U with DeepSeek-R1, and 46.6% / 20.0% with GPT-4o. These results exceed several strong prompting-based and multi-agent baselines, demonstrating that TAGS is not dependent on the MedReason corpus and generalizes well across different domains and retrieval sources.

D Generalizability Beyond Hard Subsets (Full-set Results)

Our main paper follows MEDAGENTS-BENCH (Tang et al., 2025) and reports HARD-subset results for fair comparison with prior prompting and multi-agent baselines. To validate that TAGS generalizes beyond hard cases, we additionally evaluate on the *full* test sets of three representative datasets: PubMedQA (500 cases), MedBullets (308 cases), and MedXpert-U (589 cases), under both GPT-4o and DeepSeek-R1 backbones (Table 6). Despite strong full-set performance of the base models, TAGS consistently yields additional improvements over competitive baselines across all three datasets. For GPT-4o, TAGS achieves 79.8% on PubMedQA, 73.3% on MedBullets, and 47.2% on MedXpert-U, outperforming RAG, CoT-SC, and Self-Refine. For DeepSeek-R1, TAGS further reaches 82.0%, 79.5%, and 53.6% on the same datasets, again providing the best results among compared methods. These full-set results corroborate that the advantages of TAGS are not limited to hard subsets and extend to broader evaluation settings.

E Specialist Validation

To examine whether the specialist agent exhibits domain-specific expertise, we conduct both qualitative and quantitative analyses. In a representative carpal tunnel syndrome case, the specialist produces rationales containing precise clinical terminology (e.g., *flexor tendon*, *operative report compliance*, *physical therapy*), whereas the generalist relies on more generic descriptions. Quantitatively,

Method	PubMedQA	MedBullets	MedXpert-U
Backbone: GPT-4o (Full set)			
GPT-4o	73.2	68.5	25.3
+ RAG	76.6	69.8	32.9
+ CoT-SC	74.2	71.1	36.3
+ Self-Refine	76.2	70.1	32.4
+ TAGS (ours)	79.8	73.3	47.2
Backbone: DeepSeek-R1 (Full set)			
DeepSeek-R1	71.4	75.0	34.4
+ RAG	75.6	77.6	39.5
+ CoT-SC	74.6	77.9	40.4
+ Self-Refine	76.2	75.9	37.3
+ TAGS (ours)	82.0	79.5	53.6

Table 6: Full-set Pass@1 accuracy (%) on three representative datasets. Main paper uses HARD subsets for baseline comparability; these results confirm TAGS generalizes beyond hard cases.

across a random sample of 20 questions, specialist rationales contain 1,548 MeSH tokens (Medical Subject Headings from the U.S. National Library of Medicine), compared to 1,176 tokens for the generalist (+31%), indicating richer domain-specific content. Finally, ablation results on MedXpert-U further demonstrate their complementarity: the generalist alone achieves 16%, the specialist alone 18%, while their combination under GSRC reaches 29%, yielding a +21 percentage point improvement over the base model.

F Comparison with rStar

We provide an additional comparison with rStar (Qi et al., 2024), a recent test-time reasoning method based on mutual consistency and Monte Carlo Tree Search (MCTS), which has demonstrated strong performance on difficult reasoning benchmarks.

As described in the original rStar paper, the method is primarily designed to enhance small language models (SLMs) through extensive stochastic rollouts and internal consistency scoring. In particular, rStar relies on (i) a large number of model invocations per question and (ii) access to token-level log probabilities and stochastic decoding mechanisms (e.g., dropout-based sampling) to compute mutual reasoning scores. These requirements make rStar computationally expensive and limit its applicability to closed-source API-based models.

To provide a direct empirical comparison, we conducted a preliminary evaluation of rStar on the MedXpert-U (HARD) benchmark using the Qwen2.5-7B backbone, following the implementation details described in the original paper. Table 7 reports the accuracy, inference time, and token us-

981	age per question.		
982	While rStar improves over the base model, it		
983	incurs substantially higher inference cost, requiring		
984	nearly 20× more inference time and an order of		
985	magnitude more tokens per question. In contrast,		
986	TAGS achieves significantly higher accuracy with		
987	much lower computational overhead.		
988	These results highlight a fundamental design dif-		
989	ference between the two approaches: rStar prior-		
990	itizes exhaustive search and internal consistency		
991	estimation, whereas TAGS focuses on lightweight,		
992	retrieval-augmented, and role-specialized reason-		
993	ing that remains compatible with both open-source		
994	and API-based LLMs. We emphasize that the goal		
995	of this comparison is not to diminish the contribu-		
996	tion of rStar, but to clarify the accuracy–efficiency		
997	trade-off and the complementary application sce-		
998	narios of the two methods.		
999			
1000	G Generalist–Specialist System Roles and		
	Prompt Templates		
1001	G.1 System Prompt for Specialist and		
1002	Generalist		
1003	To ensure consistency and clarity across different		
1004	model roles, we define structured system prompts		
1005	tailored to each classifier in our multi-agent frame-		
1006	work. These prompts specify role-specific reason-		
1007	ing strategies and output formats, enabling the mod-		
1008	els to adopt appropriate clinical reasoning behav-		
1009	iors under zero-shot test-time conditions.		
1010	The system prompt for the specialist categoriza-		
1011	tion agent is presented in Table 8, while the diag-		
1012	nostic prompt for the specialist agent is shown in		
1013	Table 9. The prompt for the generalist agent is		
1014	provided in Table 10.		
1015	G.2 Prompt Organization and Structure		
1016	To ensure faithful and consistent model behavior		
1017	across different roles and stages of inference, we		
1018	design modular and task-specific prompt templates.		
1019	These templates guide the models in both few-shot		
1020	reasoning and auxiliary classification tasks.		
1021	Specifically, the specialist classification prompt		
1022	(Table 11) is used to determine the relevant sub-		
1023	fields of medicine needed to solve a given question,		
1024	-serving as a basis for downstream role assignment		
1025	and retrieval. Meanwhile, the few-shot prompt tem-		
1026	plate (Table 12) provides structured instructions		
1027	and reference examples to facilitate reasoning trans-		
1028	fer for clinical question answering.		
	H The Reasoning Consistency Evaluator		1029
	Rubric and Prompt		1030
	To robustly aggregate multi-agent responses, we		1031
	introduce a reliability scoring mechanism that eval-		1032
	uates the consistency between an agent’s reasoning		1033
	and its final answer. In scenarios where a question		1034
	has many answer options (e.g., N), simple majority		1035
	voting becomes inefficient — achieving a reliable		1036
	consensus typically requires at least $N + 1$ agreeing		1037
	agents.		1038
	To address this, we employ a scoring-based ver-		1039
	ification strategy: each agent’s reasoning is eval-		1040
	uated by a separate verifier agent that assigns a		1041
	reliability score ranging from 1 to 5. This enables		1042
	us to treat scores as soft confidence signals and		1043
	aggregate responses more efficiently, even when		1044
	only a few answers are available. The resulting		1045
	per-sample reliability sum lies in the range of 4–20		1046
	(with 4 verifiers), providing fine-grained guidance		1047
	for final answer selection. The full scoring prompt		1048
	is shown in Table 13.		1049
	I Reference CoT Dataset Examples		1050
	We adopt the MedReason dataset (Wu et al.,		1051
	2025) as our external reference corpus to support		1052
	retrieval-augmented reasoning. MedReason com-		1053
	prises 32,682 high-quality question–answer pairs,		1054
	each accompanied by detailed, clinically grounded		1055
	chain-of-thought (CoT) explanations. The dataset		1056
	is constructed through a knowledge graph–guided		1057
	pipeline that ensures both logical consistency and		1058
	medical factuality.		1059
	Specifically, the authors first collect QA pairs		1060
	from seven public medical benchmarks, includ-		1061
	ing MedQA, MedMCQA, PubMedQA, MMLU,		1062
	MedXpert, Huatuo, and HLE. For each QA pair,		1063
	relevant medical entities are extracted from both the		1064
	question and the answer using GPT-4o and are then		1065
	mapped to a structured medical knowledge graph,		1066
	PrimeKG. Next, the shortest reasoning paths con-		1067
	necting the question and answer entities within the		1068
	graph are retrieved and pruned using LLM-based		1069
	selection to retain only clinically relevant paths.		1070
	These paths serve as scaffolds for guiding step-by-		1071
	step CoT generation.		1072
	To guarantee data quality, each generated rea-		1073
	soning trace is verified by prompting the LLM to		1074
	reproduce the original answer solely based on the		1075
	CoT explanation. If the answer cannot be recov-		1076
	ered, the CoT is discarded. This quality filtering		1077
	process reduces 45K generated samples to a final		1078

Table 7: Comparison with rStar on MedXpert-U (HARD) using Qwen2.5-7B.

Method	Accuracy (%)	Inference Time (s / case)	Token Usage (tokens / case)
Qwen2.5-7B + rStar	14.0	571.4	142.4K
Qwen2.5-7B + TAGS (ours)	29.0	31.6	5.2K

Specialist Categorization — System Prompt

You are a senior medical expert tasked with classifying clinical multiple-choice problems into the most relevant areas of medical science.

Your role is strictly to determine and output the classification.

Important: Do not provide any explanation, reasoning, or commentary. Only output the final classification strictly following the format.

Table 8: System prompt for the specialist categorization.

Specialist Agent — System Prompt

You are an experienced specialist in {domain}. Your role is to carefully analyze clinical multiple-choice questions from the standpoint of a {domain.lower()} expert. You should reason by focusing on the interpretation of symptoms, underlying pathophysiology, and domain-specific diagnostic principles.

First, review the provided reference examples and understand their reasoning patterns.

Then, based on your specialist knowledge, perform *structured, step-by-step* reasoning for the new question.

Required output format

Thought: [your detailed step-by-step reasoning]
 Answer: [one of A, B, C, . . .]

Table 9: System prompt for the specialist agent.

Generalist Agent — System Prompt

You are a general practitioner trained to manage a wide range of clinical conditions. Your task is to evaluate clinical multiple-choice questions using broad, cross-disciplinary medical knowledge. Focus on extracting key clinical findings, ruling out unlikely diagnoses, and applying general reasoning principles.

First, analyze the reference examples to understand their diagnostic thought process.

Then, produce a *step-by-step* analysis for the new question.

Required output format

Thought: [your detailed step-by-step reasoning]
 Answer: [one of A, B, C, . . .]

Table 10: System prompt for the generalist agent.

Specialist Classifier — Prompt

Task Instructions

- Carefully analyze the following medical question:

'''{question}'''

- The corresponding options are:

'''{options}'''

- Based on both the question and the options, determine the top **3 most relevant subfields of medicine** that are required to solve this question.

- You must only output in the exact format:

Medical Field: Field1 | Field2 | Field3

Table 11: Prompt used for classifying medical questions into relevant specialist subfields.

dataset of 32,682 verified CoTs, ensuring that each retained example is both logically sound and clinically accurate. 1079
1080
1081

The 7,168th indexed sample from MedReason is shown below as an example: 1082
1083

Question: 1084

A young male patient presents with complete rectal prolapse. The surgery of choice is? 1085
1086

Answer Choices: 1087

A. Abdominal rectopexy 1088

B. Delerom’s procedure 1089

C. Anterior resection 1090

D. Goodsall’s procedure 1091

Answer: 1092

Abdominal rectopexy. Explanation: Surgery is required, and the operation can be performed the 1093
1094

perineal or the abdominal approaches. An abdominal 1095
1096

rectopexy has a lower rate of recurrence, as 1097
1098

an abdominal procedure risks damage to the pelvic 1099
1100

autonomic nerves, resulting in possible sexual dys- 1101
1102

function, a perineal approach is also usually preferred 1103
1104

in young men. Ref: Bailey & Love’s Short 1105
1106

Practice of Surgery, 25 ed., p. 1225 1107
1108

Reasoning: 1109

Finding reasoning paths: 1. Rectal prolapse (dis- 1110
1111

Few-shot Prompt Template
<p>Header “Your task is to solve the following clinical multiple-choice question.”</p>
<p>Question Block Present the target question text, then list answer options (A/B/C/...).</p>
<p>Instruction Block</p> <ul style="list-style-type: none"> • The prompt shows N solved <i>reference examples</i>. • Each example contains: <ul style="list-style-type: none"> – Finding Reasoning Paths: brainstorming approaches – Reasoning Process: a structured, step-by-step solution • Focus <i>only</i> on learning the reasoning patterns. • For the <i>new</i> question, generate your own reasoning and answer. • Constraints: <ul style="list-style-type: none"> – Always choose one of the provided options—never “unknown” – Follow the exact output format shown below
<p>Output Format Hint Thought: [your detailed step-by-step reasoning] Answer: [one of A, B, C, . . .]</p>
<p>Reference Examples For each retrieved example i insert: Reference Example i: <question text> <options> Thought: <reference rationale></p>

Table 12: Prompt template used for few-shot reasoning with retrieved reference examples.

Reliability Scoring — System Prompt
<p>A clinical AI agent has answered the following multiple-choice question:</p>
<p>Question: “{question}”</p>
<p>Options: “{options_str\$}”</p>
<p>The agent provided the following reasoning: Thought: “{thought}”</p>
<p>Final Answer: “{answer}”</p>
<p>Your Role: You are a critical-thinking medical reviewer. Your task is to assign a reliability score from 1 to 5 based on how well the reasoning supports the answer.</p>
<p>Scoring Scale:</p> <ul style="list-style-type: none"> • 5 — Reasoning is complete, medically accurate, and fully supports the answer. • 4 — Mostly correct with minor issues, but the answer is still justified. • 3 — Reasoning has some issues or omissions, but partially supports the answer. • 2 — Reasoning is flawed or incomplete; answer is weakly supported. • 1 — Reasoning is incorrect or misleading; answer is not justified.
<p>Output Format: Score: [1-5]</p>

Table 13: System prompt used by the reasoning consistency evaluator to assign a soft reliability score.

1104 ease) -> Surgical treatment options -> Perineal ap-
1105 proach vs. Abdominal approach 2. Abdominal
1106 approach -> Abdominal rectopexy -> Risk of re-
1107 currence vs. risk of complications 3. Abdominal
1108 rectopexy -> Potential complications -> Pelvic au-
1109 tonomic nerve damage -> Male sexual dysfunction
1110 Reasoning Process: 1.**Understanding Rectal**
1111 **Prolapse and Treatment Options:** Rectal pro-
1112 lapse is a condition where the rectum protrudes
1113 through the anus. Surgical intervention is often
1114 required, and the two main approaches are the pe-
1115 rineal and abdominal approaches.
1116 2.**Evaluating Surgical Approaches: -Perineal**
1117 **Approach:** This is generally less invasive and may
1118 be preferred in older or frail patients due to lower
1119 surgical risk. However, it is associated with a
1120 higher rate of recurrence. -**Abdominal Approach:**
1121 This includes procedures like abdominal rectopexy,
1122 which tend to have a lower recurrence rate. How-
1123 ever, it is more invasive and carries a risk of com-
1124 plications.

3.**Abdominal Rectopexy and Complications:** - 1125
Recurrence Rate: Abdominal rectopexy is favored 1126
for its lower recurrence rate compared to perineal 1127
procedures. -**Complications:** A significant risk 1128
associated with the abdominal approach is poten- 1129
tial damage to the pelvic autonomic nerves. This 1130
can lead to complications such as male sexual dys- 1131
function, which is an important consideration in a 1132
young male patient. 1133
Conclusion: Abdominal rectopexy is the 1134
preferred surgical option for a young male patient 1135
with complete rectal prolapse due to its lower rate 1136
of recurrence. However, it is important to consider 1137
the risk of pelvic autonomic nerve damage, which 1138
could lead to male sexual dysfunction. This risk 1139
should be weighed against the benefits of a lower 1140

recurrence rate when deciding on the surgical approach.

J Test Dataset Examples

Test Datasets. We conduct experiments on a curated benchmark comprising challenging subsets of nine medical QA datasets, selected from the MEDAGENTS BENCH framework (Tang et al., 2025). This benchmark is specifically designed to evaluate complex medical reasoning and includes questions that require multi-step clinical inference, diagnosis formulation, and treatment planning. The original datasets include: **MedQA** (Jin et al., 2021), a multilingual dataset with English questions derived from professional medical board exams (e.g., USMLE), focusing on diagnosis and treatment knowledge. Each question includes four answer choices, with an average length of 167 tokens. **PubMedQA** (Jin et al., 2019), constructed from biomedical research articles in PubMed, requires reasoning over scientific evidence to determine yes/no/maybe answers. It contains long passages with an average length of 316 tokens. **MedMCQA** (Pal et al., 2022), sourced from Indian medical entrance exams (AIIMS and NEET-PG), comprises short multiple-choice questions across 21 medical subjects, averaging 18.7 tokens per question. **MedBullets** (Chen et al., 2024a), a high-yield medical education platform, provides detailed clinical questions with up to five answer options and longer contexts (avg. 213 tokens). **MedExQA** (Kim et al., 2024b), covers five medical specialties and emphasizes both correctness and explainability. It includes 935 questions, each with four options. **MedXpertQA** (Zuo et al., 2025), consists of two subsets: *Reasoning* and *Understanding*, specifically annotated to distinguish complex clinical inference tasks from knowledge recall. **MMLU** (Hendrycks et al., 2020) and **MMLU-Pro** (Wang et al., 2024), general-purpose multitask benchmarks that include medical subsets testing both factual recall and abstract reasoning.

K Fail case of existing method

To illustrate the limitations of existing methods, we present a representative failure case from DeepSeek-R1 on a clinically complex MedQA question. **Scenario:** A junior orthopaedic surgery resident is completing a carpal tunnel repair with the department chairman as the attending

physician. During the procedure, the resident inadvertently cuts a flexor tendon, which is then repaired without complication. The attending advises the resident to exclude this incident from the operative report to avoid patient concern.

Question: *What is the correct next action for the resident to take?*

Answer Options:

- (A) Disclose the error to the patient and put it in the operative report
- (B) Tell the attending that he cannot fail to disclose this mistake
- (C) Report the physician to the ethics committee
- (D) Refuse to dictate the operative report

Predicted Answer: (A) - Disclose the error and document it.

Model’s Reasoning: <think> Okay, let’s see. The question is about a resident who accidentally cuts a flexor tendon during carpal tunnel repair and the attending tells them not to report it. The options are A to D, and the generated solutions all say A.

First, the key points here are medical ethics, specifically honesty and patient disclosure. The resident made an error, even if it was repaired. The attending is advising to omit this from the operative report and not tell the patient to avoid worry.

In medical practice, complications should be documented and disclosed to the patient. The principle of informed consent requires patients to know about any complications, even if resolved. Failing to document could lead to legal issues and violates professional standards.

Option A says to disclose and document, which aligns with ethical guidelines. The other options either don’t directly address informing the patient (like B, which is telling the attending but not necessarily documenting) or escalate without first handling it properly (C, reporting; D, refusing to dictate).

Since all solutions suggest A, the answer is likely A. They must inform the patient and include it in the report despite the attending’s advice. </think> A

Dataset	Question & Options & Answer
MedQA	Question A 39-year-old woman presents with headaches, constipation, increased thirst, and increased urinary frequency. Her blood pressure is 160/100 mm Hg and pulse is 67/min. Blood tests reveal high blood glucose and low potassium level. Abdominal CT scan shows a well-circumscribed homogeneous mass in the left adrenal gland. The hormone excess responsible for her symptoms uses which of the following pathways to exert its action? Options (A): "Intracellular receptors", (B): "cAMP pathway", (C): "cGMP pathway", (D): "JAK/STAT pathway". Answer (A)
PubMedQA	Question Sternal fractures in childhood are rare. The aim of the study was to investigate the accident mechanism, the detection of radiological and sonographical criteria and consideration of associated injuries. In the period from January 2010 to December 2012 all inpatients and outpatients with sternal fractures were recorded according to the documentation. A total of 4 children aged 5-14 years with a sternal fracture were treated in 2000 years, 2000 children were hospitalized for pain management and 2 remained in outpatient care. Sternal fracture in growing children : A rare and often overlooked fracture? Options (A): "yes", (B): "no", (C): "maybe". Answer (C)
MedMCQA	Question Minimum number of lobes require to form a tooth? Options (A): "1", (B): "2", (C): "3", (D): "4". Answer (C)
MedBullets	Question A 22-year-old woman presents to the emergency department with shortness of breath. She was hiking when she suddenly felt unable to breathe and had to take slow deep breaths to improve her symptoms. The patient is a Swedish foreign exchange student and does not speak any English. Her medical history and current medications are unknown. Her temperature is 99.500b0F (37.500b0C), blood pressure is 127/68 mmHg, pulse is 120/min, respirations are 22/min, and oxygen saturation is 90% on room air. Physical exam is notable for poor air movement bilaterally and tachycardia. The patient is started on treatment. Which of the following parameters including forced expiratory volume in 1 second (FEV1), forced vital capacity (FVC), and diffusing capacity of carbon monoxide (DLCO) most appropriately describes this patient's underlying pathology? Options (A): "Decreased airway tone", (B): "Increased FEV1", (C): "Increased FEV1/FVC", (D): "Increased FVC", (E): "Normal DLCO". Answer (E)
MMLU	Question How many different types of microorganisms may colonize the mouth? Options (A): "35", (B): "100", (C): "350", (D): "500". Answer (C)
MMLU-Pro	Question How are new polyomaviruses detailed? Options (A): "Shot gun sequencing", (B): "Cultivation in human neural cells", (C): "Deep pyro sequencing (NGS)", (D): "Monoclonal antibody techniques". Answer (A)
MedExQA	Question Which biological tissue has the highest viscosity? Options (A): "Blood", (B): "Bone", (C): "Soft tissue", (D): "Water". Answer (B)
MedXpert-R	Question A 52-year-old paralegal comes to the clinic reporting chronic low back pain that has gradually worsened over two years. She describes an achy pain rated at 7/10 and spends considerable time at her computer desk. She experiences numbness, tingling, and radiating pain down her right leg during prolonged sitting or standing. Her examination reveals:- No urinary/bowel incontinence or perineal numbness- Right lower extremity strength of 4/5- Decreased sensation in right leg- Right patellar reflex grade 1/4- Positive straight leg raise test on right side Which nerve root levels are most likely involved in this presentation? Options (A) L1-L2 (B) L5-S1 (C) L2-L3 (D) S2-S3 (E) L4-L5 (F) L4-S1 (G) L3-L4 (H) S1-S2 (I) T12-L1 (J) L5-S2. Answer G
MedXpert-U	Which hypoxic cell radiosensitizer, known for undergoing redox recycling or decomposing into a toxic product, has demonstrated effectiveness in the treatment of head and neck cancer? Options (A): Doranidazole, (B): Tirapazamine, (C): Camptothecin, (D): Misonidazole, (E): Pimonidazole, (F): Nimorazole, (G): Sanazole, (H): Cetuximab, (I): Etanidazole, (J): Methotrexate. Answer (F)

Table 14: Dataset examples with corresponding questions and answer options from the nine test datasets.

Issue:

The model selected (A) as the answer, which is inconsistent with the reference solution (B). The error occurred because the model focused solely on disclosure and documentation, ignoring the specific instruction to address the attending's unethical guidance.