



# MEDAGENTS: Large Language Models as Collaborators for Zero-shot Medical Reasoning

Anonymous ACL submission

## Abstract

Large Language Models (LLMs), despite their remarkable progress across various general domains, encounter significant barriers in medicine and healthcare. This field faces unique challenges such as domain-specific terminologies and reasoning over specialized knowledge. To address these issues, we propose a novel Multi-disciplinary Collaboration (MC) framework for the medical domain that leverages role-playing LLM-based agents who participate in a collaborative multi-round discussion, thereby enhancing LLM proficiency and reasoning capabilities. This training-free and interpretable framework encompasses five critical steps: gathering domain experts, proposing individual analyses, summarising these analyses into a report, iterating over discussions until a consensus is reached, and ultimately making a decision. Our work focuses on the zero-shot setting, which is applicable in real-world scenarios. Experimental results on nine datasets (MedQA, MedMCQA, PubMedQA, and six subtasks from MMLU) establish that our proposed MC framework excels at mining and harnessing the medical expertise within LLMs, as well as extending its reasoning abilities.

## 1 Introduction

Large language models (LLMs) (Brown et al., 2020; Scao et al., 2022; Chowdhery et al., 2022; Touvron et al., 2023; OpenAI, 2023) have exhibited notable generalization abilities across a wide range of tasks and applications (Lu et al., 2023; Zhou et al., 2023; Park et al., 2023), with these capabilities stemming from their extensive training on vast comprehensive corpora covering diverse topics. However, in real-world scenarios, LLMs are inclined to encounter domain-specific tasks that necessitate a combination of domain expertise and complex reasoning abilities (Moor et al., 2023; Wu et al., 2023b; Singhal et al., 2023a; Yang

et al., 2023). Amidst this backdrop, a noteworthy research topic lies in the adoption of LLMs in the medical field, which has gained increasing prominence recently (Zhang et al., 2023b; Bao et al., 2023; Singhal et al., 2023a).

Currently, there are two dominant challenges that hinder LLMs from managing medical-related tasks: (i) The *volume and specificity* of training data in the medical field are limited compared with general web data used to train LLMs due to cost and privacy concerns (U.S. Department of Health and Human Services, 1996). Therefore, it remains insufficient to understand or recall the required medical expertise via simple and direct prompting (Kung et al., 2023; Singhal et al., 2023a). (ii) High performance in this field requires *extensive domain knowledge* (Schmidt and Rikers, 2007) and *sophisticated reasoning abilities* upon it (Liévin et al., 2022), thus posing heightened demands to LLMs for medical tasks.

At the same time, as opposed to the conventional single *input-output* pattern, recent research has surprisingly witnessed the success of LLM-based agents across a spectrum of tasks (Xi et al., 2023; Wang et al., 2023c). Among such work, the design of multi-agent collaboration favorably stands out by highlighting the simulation of human activities (Du et al., 2023; Liang et al., 2023; Park et al., 2023) and coordinating the potential of multiple agents (Chen et al., 2023; Li et al., 2023d; Hong et al., 2023). Through the design of multi-agent collaboration, the expertise implicitly embedded within LLMs or that the model has encountered during its training, which may not be readily accessible via traditional prompting, is effectively brought to the fore. In turn, this process enhances the model’s reasoning capabilities throughout multi-round interaction (Wang et al., 2023c,b; Du et al., 2023; Fu et al., 2023).

Inspired by the above ideas, we propose a **Multi-disciplinary Collaboration (MC)** framework in

A 66-year-old male with a history of **heart attack** and recurrent **stomach ulcers** is experiencing persistent **cough and chest pain**, and recent **CT scans** indicate a possible **lung tumor**. Designing a treatment plan that minimizes risk and maximizes outcomes is the current concern due to his deteriorating health and medical history.

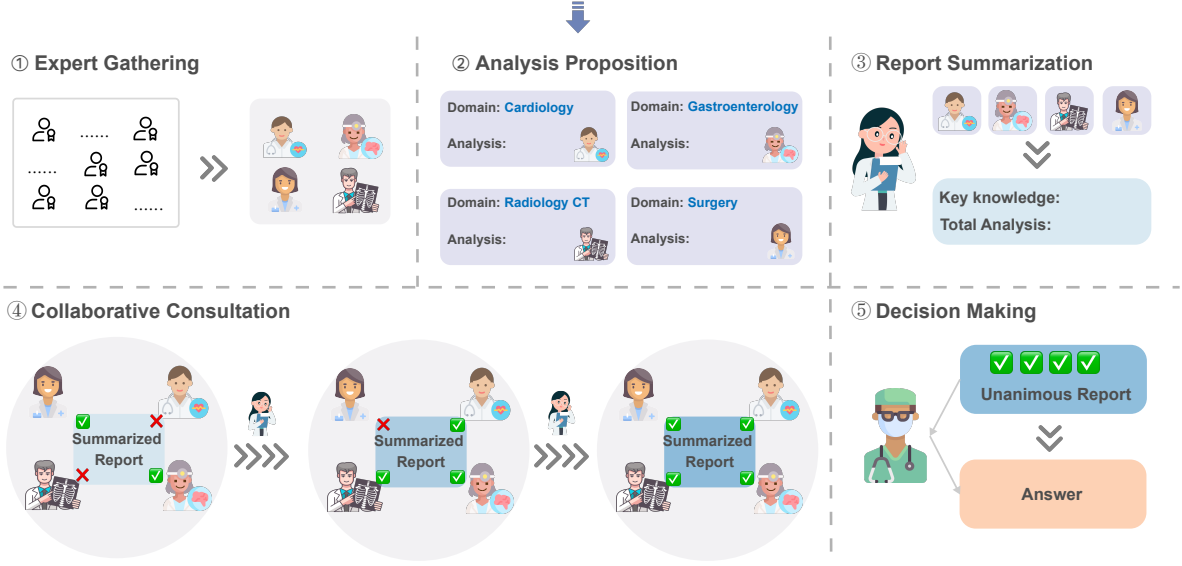


Figure 1: Diagram of our proposed multi-disciplinary collaboration framework. Given a medical question as input, the framework performs reasoning in five stages: (i) expert gathering; (ii) analysis proposition; (iii) report summarization; (iv) collaborative consultation; and (v) decision making.

the clinical domain, aiming to unveil the intrinsic medical knowledge from LLMs as well as bolster the reasoning competence in a training-free and interpretable manner. As is shown in Figure 1, the MC framework is based on five pivotal steps (i) Expert gathering: gather experts from distinct disciplines according to the clinical question. (ii) Analysis proposition: domain experts put forward their analyses with their expertise. (iii) Report summarization: compose a summarized report based on a previous series of analyses. (iv) Collaborative consultation: engage the experts in discussions over the summarized report. The report will be revised iteratively until an agreement from all experts is reached. (v) Decision making: derive a final decision from the unanimous report.

Having established the theoretical foundation of our approach, we conduct experiments on nine datasets (Singhal et al., 2023a), including MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022), PubMedQA (Jin et al., 2019) and six medical subtasks from MMLU (Hendrycks et al., 2020), similar to Flan-PaLM (Singhal et al., 2023a). To better align with real-world application scenarios, our study focuses on the zero-shot setting. Encouragingly, our proposed approach outperforms settings for both chain-of-thought (CoT) and self-consistency prompting methods.

Most notably, our approach demonstrates better performances under the zero-shot setting compared with the few-shot (5-shot) strong baselines.

Based on our results, we further investigate the influence of agent numbers and conduct human evaluations to pinpoint the limitations and issues prevalent in our approach. We find four common prevalent categories of errors: (i) lack of domain knowledge; (ii) mis-retrieval of domain knowledge; (iii) consistency errors; and (iv) CoT errors. Further refinements focused on mitigating these particular shortcomings would enhance the model’s proficiency and reliability.

Our contributions are summarized as follows:

(i) We propose a multi-disciplinary collaboration (MC) framework for zero-shot medical reasoning tasks. This novel approach endeavors to unveil the inherent clinical expertise present in LLMs and enhance their reasoning competence.

(ii) Experimental results on nine datasets demonstrate the general effectiveness of our proposed MC framework.

(iii) We identify and categorize common error types in our approach through rigorous human evaluation to shed light on future studies.

## 2 Related Work

### 2.1 LLMs in Medical Domains

Recent years have witnessed the impressive advancements brought about by LLMs across various domains (Ling et al., 2023; Wu et al., 2023b; Singhal et al., 2023a; Yang et al., 2023), among which a promising and noteworthy application lies in the medical domain (Bao et al., 2023; Nori et al., 2023; Rosol et al., 2023). Although LLMs have demonstrated their potential in distinct medical applications encompassing diagnostics (Singhal et al., 2023a; Han et al., 2023), genetics (Duong and Solomon, 2023; Jin et al., 2023), pharmacist (Liu et al., 2023), and medical evidence summarization (Tang et al., 2023b,a; Shaib et al., 2023), concerns persist when LLMs encounter clinical inquiries that demand intricate medical expertise and decent reasoning abilities (Umapathi et al., 2023; Singhal et al., 2023a). Consequently, it is of crucial importance to further tap into the medical expertise to arm LLMs with enhanced clinical reasoning capabilities. Currently, there are two major lines of research on LLMs in medical domains, namely tool-augmented methods and instruction-tuning methods.

For tool-augmented approaches, recent studies rely on external tools to acquire additional information for clinical reasoning. For instance, GeneGPT (Jin et al., 2023) guided LLMs to leverage the Web APIs of the National Center for Biotechnology Information (NCBI) to meet various biomedical information needs. Zakka et al. (2023) proposed Almanac, a framework that is augmented with retrieval capabilities for medical guidelines and treatment recommendations. Kang et al. (2023) introduced a method named KARD to improve small LMs on specific domain knowledge by fine-tuning small LMs on the rationales generated from LLMs and augmenting small LMs with external knowledge from a non-parametric memory.

For instruction tuning methods, existing research makes use of external clinical knowledge bases and self-prompted data to obtain instruction datasets (Tu et al., 2023; Zhang et al., 2023a; Singhal et al., 2023b; Tang et al., 2023c), which are then employed to tune LLMs on medical domains. For example, LLaVA-Med (Li et al., 2023a) leveraged a broad-coverage biomedical figure-caption dataset collected from PubMed Central and took advantage of GPT-4 to self-instruct open-ended instruction-following data from the captions

to fine-tune a large general-domain vision-language model. MedChatZH (Tan et al., 2023) served as a dialogue model for traditional Chinese medical QA, which was pre-trained on Chinese traditional medical books and finetuned with an elaborated medical instruction dataset. AlpaCare (Zhang et al., 2023b) benefited from its large-scale and diverse medical instruction-following data MedInstruct-52k, resulting in remarkable generality and medical proficiency. Our work, nevertheless, shifts to emphasize mining and deriving medical knowledge from within LLMs and enhancing reasoning in a training-free manner.

### 2.2 LLM-based Multi-agent Collaboration

The development of LLM-based agents has made significant progress in the community by endowing LLMs with the ability to perceive surroundings and make decisions individually (Wang et al., 2023a; Yao et al., 2022; Nakajima, 2023; Xie et al., 2023; Zhou et al., 2023). Beyond the initial single-agent mode, the multi-agent pattern has garnered increasing attention recently (Xi et al., 2023; Li et al., 2023d; Hong et al., 2023) which further explores the potential of LLM-based agents by learning from multi-turn feedback and cooperation. In essence, the key to LLM-based multi-agent collaboration is the simulation of human activities such as role-playing (Wang et al., 2023c; Hong et al., 2023) and communication (Wu et al., 2023a; Qian et al., 2023; Li et al., 2023b,c). For instance, Solo Performance Prompting (SPP) (Wang et al., 2023c) managed to combine the strengths of multiple minds to improve performance by dynamically identifying and engaging multiple personas over the course of task-solving. Camel (Li et al., 2023b) leveraged role-playing to enable chat agents to communicate with each other for task completion. Several recent works attempt to incorporate adversarial collaboration including debates (Du et al., 2023; Xiong et al., 2023) and negotiation (Fu et al., 2023) among multiple agents to further boost performance. Liang et al. (2023) proposed a multi-agent debate framework in which various agents put forward their statements in a *tit for tat* pattern. Inspired by the multi-disciplinary consultation mechanism which is common and effective in hospitals, we are thus inspired to apply this mechanism to medical reasoning tasks through LLM-based multi-agent collaboration.

**Question:** A 3-month-old infant is brought to her pediatrician because she coughs and seems to have difficulty breathing while feeding. In addition, she seems to have less energy compared to other babies and appears listless throughout the day. She was born by cesarean section to a G1P1 woman with no prior medical history and had a normal APGAR score at birth. Her parents say that she has never been observed to turn blue. Physical exam reveals a high-pitched holosystolic murmur that is best heard at the lower left sternal border. The most likely cause of this patient's symptoms is associated with which of the following abnormalities?

**Options:** (A) 22q11 deletion (B) Deletion of genes on chromosome 7 (C) Lithium exposure in utero (D) Retinoic acid exposure in utero

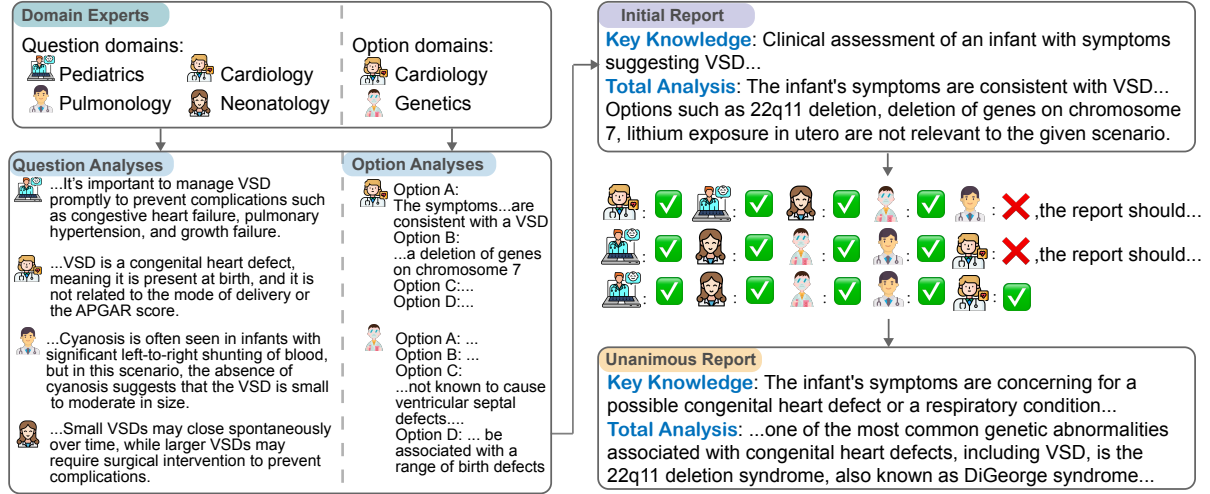


Figure 2: Illustrative example of our proposed Multi-disciplinary Collaboration (MC) framework.

### 3 Method

This section presents the details of our proposed Multi-disciplinary Collaboration (MC) framework. Figure 1 and 2 give an overview and an illustrative example of its pipeline. Our proposed MC framework works in five stages: (i) expert gathering: assemble experts from various disciplines based on the clinical question; (ii) analysis proposition: domain experts present their own analyses with their expertise; (iii) report summarization: develop a report summary on the basis of previous analyses; (iv) collaborative consultation: hold a consultation over the summarized report with the experts. The report will be revised repeatedly until every expert has given their approval. (v) decision making: derive a final decision from the unanimous report.<sup>1</sup>

#### 3.1 Expert Gathering

Given a clinical question  $q$  and a set of options  $op = \{o_1, o_2, \dots, o_k\}$ , the goal of the Expert Gathering stage is to recruit a group of question domain experts  $\mathcal{QD} = \{qd_1, qd_2, \dots, qd_m\}$  and option domain experts  $\mathcal{OD} = \{od_1, od_2, \dots, od_n\}$ . Specifically, we assign a role to the model and provide instructions to guide the model output to the corresponding domains based on the input

<sup>1</sup>Details about all guideline prompts and roles are shown in Section A for clarification.

question and options, respectively:

$$\begin{aligned} \mathcal{QD} &= \text{LLM}(q, r_{qd}, \text{prompt}_{qd}), \\ \mathcal{OD} &= \text{LLM}(q, op, r_{od}, \text{prompt}_{od}), \end{aligned} \quad (1)$$

where  $(r_{qd}, \text{prompt}_{qd})$  and  $(r_{od}, \text{prompt}_{od})$  stand for the system role and guideline prompt to gather domain experts for the question  $q$  and options  $op$ .

#### 3.2 Analysis Proposition

After gathering domain experts for the question  $q$  and options  $op$ , we aim to inquire experts to generate corresponding analyses prepared for later reasoning:  $\mathcal{QA} = \{qa_1, qa_2, \dots, qa_m\}$  and  $\mathcal{OA} = \{oa_1, oa_2, \dots, oa_n\}$ .

**Question Analyses** Given a question  $q$  and a question domain  $qd_i \in \mathcal{QD}$ , we ask LLM to serve as an expert specialized in domain  $qd_i$  and derive the analyses for the question  $q$  following the guideline prompt  $\text{prompt}_{qa}$ :

$$qa_i = \text{LLM}(q, qd_i, r_{qa}, \text{prompt}_{qa}). \quad (2)$$

**Option Analyses** Now that we have an option domain  $od_i$  and question analyses  $\mathcal{QA}$ , we can further analyze the options by taking into account both the relationship between the options and the relationship between the options and question. Concretely, we deliver the question  $q$ , the options



Table 1: Summary of the Datasets. Part of the values are from the appendix of (Singhal et al., 2023a).

Dataset	Format	Choice	Testing Size	Domain
MedQA	Question + Answer	A/B/C/D	1273	US Medical Licensing Examination
MedMCQA	Question + Answer	A/B/C/D and Explanations	6.1K	AIIMS and NEET PG entrance exams
PubMedQA	Question + Context + Answer	Yes/No/Maybe	500	PubMed paper abstracts
MMLU	Question + Answer	A/B/C/D	1089	Graduate Record Examination & US Medical Licensing Examination

### Algorithm 1: Collaborative Consultation

**Input:** Domain experts  $D = \{d_1, \dots, d_n\}$ , initial report  $R_0$ , Model  $\mathcal{M}$ , maximum attempts  $k$ , prompts  $\{p_{vote}, p_{mod}, p_{rev}\}$   
**Output:** Final report  $R_f$

// Initialize variables  
 $nocon\_flag \leftarrow True, n_{try} \leftarrow 0$   
 $R_{cur} \leftarrow R_0, Mods \leftarrow \emptyset$

// Iterative review  
**while**  $nocon\_flag$  is  $True$  and  $n_{try} < k$  **do**  
     $n_{try} \leftarrow n_{try} + 1$   
     $nocon\_flag \leftarrow False$   
    // vote for the report  
    **for**  $i$  in  $1, \dots, n$  **do**  
         $vote_i \leftarrow \mathcal{M}(R_{cur}, d_i, p_{vote})$   
        // propose modifications  
        **if**  $vote_i$  is  $no$  **then**  
             $Mod_i \leftarrow \mathcal{M}(R_{cur}, d_i, p_{mod})$   
            Update  $Mods$  with  $Mod_i$   
             $nocon\_flag \leftarrow True$   
        **end**  
    **end**  
    // modify the report  
    **if**  $nocon\_flag$  is  $True$  **then**  
         $R_{cur} \leftarrow \mathcal{M}(R_{cur}, Mods, p_{rev})$   
    **end**  
**end**  
**return**  $R_f \leftarrow R_{cur}$

$op$ , a specific option domain  $od_i \in \mathcal{OD}$ , and the question analyses  $\mathcal{QA}$  to the LLM:

$$oa_i = \text{LLM}(q, op, od_i, \mathcal{QA}, r_{oa}, \text{prompt}_{oa}). \quad (3)$$

### 3.3 Report Summarization

In the Report Summarization stage, we attempt to summarize and synthesize previous analyses from various domain experts  $\mathcal{QA} \cup \mathcal{OA}$ . Given question analyses  $\mathcal{QA}$  and option analyses  $\mathcal{OA}$ , we ask LLMs to play the role of a medical report assistant, allowing it to generate a synthesized report by extracting key knowledge and total analysis based on previous analyses:

$$Repo = \text{LLM}(\mathcal{QA}, \mathcal{OA}, r_{rs}, \text{prompt}_{rs}). \quad (4)$$

### 3.4 Collaborative Consultation

Since we have a preliminary summary report  $Repo$ , the objective of the Collaborative Consultation stage is to engage distinct domain experts in multiple rounds of discussions and ultimately render a summary report that is recognized by all experts. The overall procedure of this phase is presented in Algorithm 1. During each round of discussions, the experts give their votes (yes/no) as well as modification opinions if they vote *no* for the current report. Afterward, the report will be revised based on the modification opinions. Specifically, during the  $i$ -th round of discussion, we note the modification comments from the experts as  $Mod_i$ , then we can acquire the updated report as  $Repo_i = \text{LLM}(Repo_{i-1}, Mod_i, \text{prompt}_{mod})$ . In this way, the discussions are held iteratively until all experts vote *yes* for the final report  $Repo_f$ .

### 3.5 Decision Making

In the end, we demand LLM act as a medical decision maker to derive the final answer to the clinical question  $q$  referring to the unanimous report  $Repo_f$ :

$$ans = \text{LLM}(q, op, Repo_f, r_{dm}, \text{prompt}_{dm}). \quad (5)$$

## 4 Experiments

### 4.1 Setup

**Tasks and Datasets.** We evaluate our MC framework on three benchmark datasets MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022), and PubMedQA (Jin et al., 2019), as well as six subtasks most relevant to the medical domain from MMLU datasets (Hendrycks et al., 2020) including anatomy, clinical knowledge, college medicine, medical genetics, professional medicine, and college biology. Table 1 summarizes the data statistics. MedQA consists of USMLE-style questions with four or five possible answers.

Table 2: Main results on MedQA, MedMCQA, PubMedQA, and six subtasks from MMLU including anatomy, clinical knowledge, college medicine, medical genetics, professional medicine, and college biology (Acc). SC denotes the self-consistency prompting method. Results in **bold** are the best performances.

Method	MedQA	MedMCQA	PubMedQA	Anatomy	Clinical knowledge	College medicine	Medical genetics	Professional medicine	College biology	Avg.
<b>Flan-Palm</b>										
Few-shot CoT	60.3	53.6	77.2	66.7	77.0	83.3	75.0	76.5	71.1	71.2
Few-shot CoT + SC	67.6	57.6	75.2	71.9	80.4	88.9	74.0	83.5	76.3	75.0
<b>GPT-3.5</b>										
<i>*few-shot setting</i>										
Few-shot	54.7	56.7	67.6	65.9	71.3	59.0	72.0	75.7	73.6	66.3
Few-shot CoT	55.3	54.7	71.4	48.1	65.7	55.5	57.0	69.5	61.1	59.8
Few-shot CoT + SC	62.1	58.3	73.4	70.4	76.2	69.8	78.0	79.0	77.2	71.6
<i>*zero-shot setting</i>										
Zero-shot	54.3	56.3	73.7	61.5	76.2	63.6	74.0	75.4	75.0	67.8
Zero-shot CoT	44.3	47.3	61.3	63.7	61.9	53.2	66.0	62.1	65.3	58.3
Zero-shot CoT + SC	61.3	52.5	<b>75.7</b>	<b>71.1</b>	75.1	68.8	76.0	<b>82.3</b>	75.7	70.9
MC framework ( <b>Ours</b> )	<b>64.1</b>	<b>59.3</b>	72.9	65.2	<b>77.7</b>	<b>69.8</b>	<b>79.0</b>	82.1	<b>78.5</b>	<b>72.1</b>
<b>GPT-4</b>										
<i>*few-shot setting</i>										
Few-shot	76.6	70.1	73.4	79.3	89.5	75.6	<b>93.0</b>	91.5	91.7	82.3
Few-shot CoT	73.3	63.2	74.9	75.6	89.9	61.0	79.0	79.8	63.2	73.3
Few-shot CoT + SC	82.9	73.1	75.6	80.7	90.0	<b>88.2</b>	90.0	95.2	93.0	85.4
<i>*zero-shot setting</i>										
Zero-shot	73.0	69.0	76.2	78.5	83.3	75.6	90.0	90.0	90.0	80.6
Zero-shot CoT	61.8	69.0	71.0	82.1	85.2	80.8	92.0	93.5	91.7	80.8
Zero-shot CoT + SC	74.5	70.1	75.3	80.0	86.3	81.2	<b>93.0</b>	94.8	91.7	83.0
MC framework ( <b>Ours</b> )	<b>83.7</b>	<b>74.8</b>	<b>76.8</b>	<b>83.5</b>	<b>91.0</b>	87.6	<b>93.0</b>	<b>96.0</b>	<b>94.3</b>	<b>86.7</b>

MedMCQA encompasses four-option multiple-choice questions from Indian medical entrance examinations (AIIMS/NEET). MMLU (Massive Multitask Language Understanding) covers 57 subjects across various disciplines, including STEM, humanities, social sciences, and many others. The scope of its assessment stretches from elementary to advanced professional levels, evaluating both world knowledge and problem-solving capabilities. While the subject areas tested are diverse, encompassing traditional fields like mathematics and history, as well as more specialized areas like law and ethics, we deliberately limit our selection to the sub-subjects within the medical domain for this exercise, following (Singhal et al., 2023a).

**Implementation.** We utilize the popular and publicly available GPT-3.5-Turbo and GPT-4 (OpenAI, 2023) from Azure OpenAI Service.<sup>2</sup> All experiments are conducted in the **zero-shot** setting. The temperature is set to 1.0 and *top\_p* to 1.0 for all generations. The number *k* of options is 4 except for PubMedQA (3). The numbers of

domain experts for the question and options are set as:  $m = 5, n = 2$  except for PubMedQA ( $m = 4, n = 2$ ). Considering the costly API expenses, we randomly sample 300 examples for each dataset and conduct experiments on them.

## 4.2 Main Results

Table 2 presents the main results on the nine datasets, including MedQA, MedMCQA, PubMedQA, and six subtasks from MMLU. We compare our method with several baselines including CoT and self-consistency prompting in both zero-shot and few-shot settings. Notably, our proposed MC framework outperforms the zero-shot baseline methods by a large margin, indicating the effectiveness of our MC framework in real-world application scenarios. Furthermore, our approach surprisingly demonstrates comparable performance under the zero-shot setting compared with the strong baseline *Few-shot CoT+SC*.

## 5 Analysis

### 5.1 Ablation Study

Since our MC framework simulates a multi-disciplinary collaboration process that contains

<sup>2</sup><https://learn.microsoft.com/en-us/azure/ai-services/openai/>

Table 3: Ablation study for different processes in our MC framework. Anal: Analysis proposition, Summ: Report summarization, Cons: Collaborative consultation.

Method	Accuracy(%)
Direct Prompting	49.0
CoT Prompting	55.0
<b>w/ MedAgents</b>	
+ Anal	62.0(↑ 7.0)
+ Anal & Summ	65.0(↑ 10.0)
+ Anal & Summ & Cons	67.0(↑ 12.0)

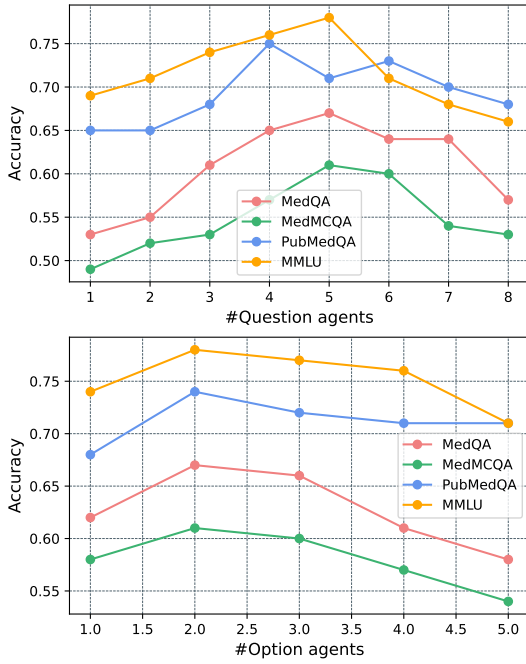


Figure 3: Influence of the number of question agents and option agents on MedQA, MedMCQA, PubMedQA, and MMLU.

multiple intermediate steps, a natural question is whether each intermediate step contributes to the ultimate result. To investigate this, we ablate three major processes, namely *analysis proposition*, *report summarization* and *collaborative consultation*. Results in Table 3 show that all of these processes are non-trivial. Notably, the proposition of MEDAGENTS substantially boosts the performance (i.e., 55.0%→62.0%), whereas the subsequent processes achieve relatively slight improvements over the previous one (i.e., 62.0%→65.0/67.0%). This suggests that the initial role-playing agents are responsible for exploring medical knowledge of various levels and aspects within LLMs, while the following processes play a role in further verification and revision.

Table 4: Optimal number of agents on MedQA, MedMCQA, PubMedQA, and MMLU.

Dataset	MedQA	MedMCQA	PubMedQA	MMLU
#Question agents	5	5	4	5
#Option agents	2	2	2	2

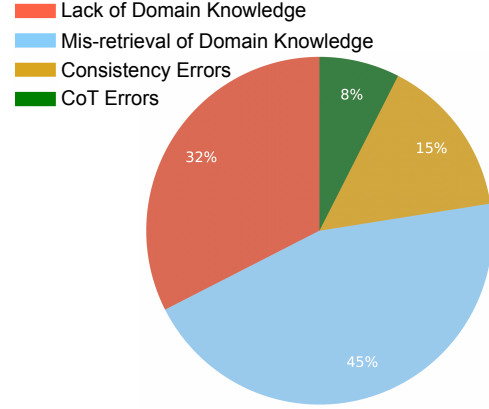


Figure 4: Ratio of different categories in error cases.

## 5.2 Number of agents

As our proposed MC framework involves multiple agents that play certain roles to acquire the ultimate answer, we explore how the number of collaborating agents influences the overall performance. We vary the number of question agents and option agents while fixing the other variable to observe the performance trends on the MedQA dataset. As shown in Figure 3, the number of question agents and option agents peaks at 5 and 2, respectively. Table 4 summarizes that the optimal number of question agents is 5 for MedQA, MedMCQA, and 4 for PubMedQA, beyond which there may be diminishing returns or potential confusion caused by information overload.

## 5.3 Error Analysis

Based on our results, we conduct a human evaluation to pinpoint the limitations and issues prevalent in our model. We distill these errors into four major categories: (i) **Lack of Domain Knowledge**: these errors occur when the model demonstrates an inadequate understanding of the specific medical knowledge necessary to provide an accurate response; (ii) **Mis-retrieval of Domain Knowledge**: the model has the necessary domain knowledge but fails to retrieve or apply it correctly in the given context; (iii) **Consistency Errors**: such errors arise when the model provides differing responses to the same statement. The inconsistency

Category	Example	Interpretation
Lack of Domain Knowledge	...The hypopigmented rash <span style="color: red;">✗</span> is a classic symptom of <span style="color: orange;">cutaneous larva migrans</span> . To confirm the diagnosis, a skin biopsy <span style="color: red;">✗</span> would be the most appropriate test.	About <span style="color: orange;">cutaneous larva migrans</span> : 1. symptoms: <span style="color: red;">✗</span> not simply hypopigmented rash 2. diagnostic method: <span style="color: red;">✗</span> skin biopsy is not preferred
Mis-retrieval of Domain Knowledge	...The physician instructs the patient to stand from a supine position while still wearing the stethoscope. It is known as the " <span style="color: orange;">Valsalva maneuver</span> " <span style="color: red;">✗</span> . During the Valsalva maneuver, ...	The patient is asked to merely stand from a supine position. It does not involve the <span style="color: orange;">Valsalva maneuver</span> . <span style="color: red;">✗</span>
Consistency Errors	...Option A states that there is a decrease in systolic blood pressure of 20 mmHg within 6 minutes. This is a correct statement, as a drop in systolic blood pressure of at least <span style="color: orange;">20 mmHg within 3 minutes</span> of standing up is a diagnostic criterion for postural hypotension...	Correct statement: <span style="color: orange;">20mmHg within 3 minutes</span> Option A: <span style="color: orange;">20mmHg within 6 minutes</span> <span style="color: red;">✗</span>
CoT Errors	Q: Deciduous teeth do not show fluorosis because: ... (A) Placenta acts as a barrier: While it's true that placenta can act as a barrier for certain substances, this option is <span style="color: red;">not relevant</span> <span style="color: red;">✗</span> to the question...	placenta can as a barrier for certain substances such as fluoride, which is <span style="color: orange;">part of the reason</span> why deciduous teeth do not show fluorosis...

Figure 5: Examples of error cases from MedQA and MedMCQA datasets in four major categories including: lack of domain knowledge, mis-retrieval of domain knowledge, consistency errors, and CoT errors.

suggests confusion in the model’s understanding or application of the underlying knowledge; (iv) **CoT Errors**: errors under this category pertain to flawed reasoning sequences or lapses in logical cohesion. The model may form and follow inaccurate rationales, leading to incorrect conclusions.

We randomly select 40 error cases in MedQA and MedMCQA datasets and analyze the percentage of different categories in these error cases. As is shown in Figure 4, the majority (77%) of the error examples are due to confusion about the domain knowledge (including the lack and mis-retrieval of domain knowledge), which illustrates that although our method further mines medical knowledge concealed within LLMs via multi-disciplinary consultation, there still exists a portion of domain knowledge that is explicitly beyond the intrinsic knowledge of LLMs, leading to a bottleneck of our proposed MC framework. As a result, our analysis sheds light on future directions to mitigate the aforementioned drawbacks and further strengthen the model’s proficiency and reliability. One potential solution is incorporating credible medical knowledge sources to complement the existing shortcomings.

To illustrate the error examples more intuitively, we select four typical samples from the four error categories, which can be shown in Figure 4: (i) The first error is due to a lack of domain knowledge regarding *cutaneous larva migrans*, whose symptoms are not purely *hypopigmented rash*, as well as the fact that *skin biopsy* is not

an appropriate test method, which results in the hallucination phenomenon. (ii) The second error is caused by mis-retrieval of domain knowledge, wherein the fact in green is not relevant to *Valsalva maneuver*. (iii) The third error is attributed to consistency errors, where the model incorrectly regards *20 mmHg within 6 minutes* and *20 mmHg within 3 minutes* as the same meaning. (iv) The fourth error is provoked by incorrect inference about the relevance of a fact and option A in CoT.

## 6 Conclusion

This paper presents a novel multi-disciplinary collaboration framework for the medical domain that leverages role-playing LLM-based agents who participate in a collaborative multi-round discussion. The framework is training-free and interpretable, encompassing five critical steps: gathering domain experts, proposing individual analyses, summarising these analyses into a report, iterating over discussions until a consensus is reached, and ultimately making a decision. Experimental results on nine datasets show that our proposed framework outperforms all the zero-shot baselines by a large margin and demonstrates comparable performance with the strong few-shot baseline with self-consistency. According to our human evaluations on error cases, future studies may further improve the framework by mitigating the mistakes due to the lack of domain knowledge, mis-retrieval of domain knowledge, and addressing consistency errors and CoT errors.



## Limitation

The proposed Multi-disciplinary Collaboration (MC) framework for medical reasoning tasks has shown promising results, but it is important to address certain limitations.

Firstly, the time-varying nature of medical consensus poses a challenge as the parameterized knowledge in LLMs may need to be updated over time. Changes in medical knowledge can significantly impact the accuracy of medical decision-making, and continuous efforts are required to keep the framework up-to-date.

Also, there is a lack of explicit access to domain knowledge. While the MC framework leverages the inherent medical knowledge embedded within LLMs, there is still a portion of domain knowledge that remains beyond the intrinsic understanding of LLMs. This limitation can result in errors in reasoning and decision-making, highlighting the need for complementary sources of domain knowledge to enhance the framework’s accuracy.

Additionally, the applicability of the MC framework may be limited in low-resource languages. These languages often lack sufficient training data and resources, making it difficult to perform effective medical evaluations. Addressing this limitation requires additional research and development to adapt the framework to low-resource languages and their specific medical needs.

Finally, model biases and ethical considerations must be carefully addressed. Just like any AI system, the MC framework may inadvertently inherit biases present in the training data, which can lead to unfair or discriminatory outcomes. It is essential to actively mitigate these biases and consider the ethical implications in order to ensure fairness, accuracy, and equity in medical decision-making.

## References

- Zhijie Bao, Wei Chen, Shengze Xiao, Kuang Ren, Jiaao Wu, Cheng Zhong, Jiajie Peng, Xuanjing Huang, and Zhongyu Wei. 2023. [Disc-medllm: Bridging general large language models and real-world medical consultation](#).
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child,

- Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, et al. 2023. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [Palm: Scaling language modeling with pathways](#). *ArXiv preprint, abs/2204.02311*.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. [Improving factuality and reasoning in language models through multiagent debate](#).
- Dat Duong and Benjamin D Solomon. 2023. Analysis of large-language model versus human performance for genetics questions. *European Journal of Human Genetics*, pages 1–3.
- Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023. [Improving language model negotiation with self-play and in-context learning from ai feedback](#).
- Tianyu Han, Lisa C. Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander L  user, Daniel Truhn, and Keno K. Bressen. 2023. [Medalpaca – an open-source collection of medical conversational ai models and training data](#).
- 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*.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, and Chenglin Wu. 2023. [Metagpt: Meta programming for multi-agent collaborative framework](#).
- 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*.

597	Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu.	Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-	651
598	2023. Genegpt: Augmenting large language models	Wei Chang, Ying Nian Wu, Song-Chun Zhu, and	652
599	with domain tools for improved access to biomedical	Jianfeng Gao. 2023. Chameleon: Plug-and-play	653
600	information. <i>ArXiv</i> .	compositional reasoning with large language models.	654
		<i>arXiv preprint arXiv:2304.09842</i> .	655
601	Minki Kang, Seanie Lee, Jinheon Baek, Kenji	Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein	656
602	Kawaguchi, and Sung Ju Hwang. 2023. Knowledge-	Abad, Harlan M Krumholz, Jure Leskovec, Eric J	657
603	augmented reasoning distillation for small language	Topol, and Pranav Rajpurkar. 2023. Foundation	658
604	models in knowledge-intensive tasks. <i>arXiv preprint</i>	models for generalist medical artificial intelligence.	659
605	<i>arXiv:2305.18395</i> .	<i>Nature</i> , 616(7956):259–265.	660
606	Tiffany H Kung, Morgan Cheatham, Arielle Medenilla,	Y Nakajima. 2023. Task-driven autonomous agent	661
607	Czarina Sillos, Lorie De Leon, Camille Elepaño,	utilizing gpt-4, pinecone, and langchain for	662
608	Maria Madriaga, Rimel Aggabao, Giezel Diaz-	diverse applications. See <a href="https://yoheinakajima.com/task-driven-autonomous-agent-utilizing-gpt-4-pinecone-and-langchain-for-diverse-applications">https://yoheinakajima.</a>	663
609	Candido, James Maningo, et al. 2023. Performance	<i>com/task-driven-autonomous-agent-utilizing-gpt-4-</i>	664
610	of chatgpt on usmle: Potential for ai-assisted medical	<i>pinecone-and-langchain-for-diverse-applications</i>	665
611	education using large language models. <i>PLoS digital</i>	(accessed 18 April 2023).	666
612	<i>health</i> , 2(2):e0000198.		
613	Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto	Harsha Nori, Nicholas King, Scott Mayer McKinney,	667
614	Usuyama, Haotian Liu, Jianwei Yang, Tristan	Dean Carignan, and Eric Horvitz. 2023. Capabilities	668
615	Naumann, Hoifung Poon, and Jianfeng Gao. 2023a.	of gpt-4 on medical challenge problems. <i>arXiv</i>	669
616	Llava-med: Training a large language-and-vision	<i>preprint arXiv:2303.13375</i> .	670
617	assistant for biomedicine in one day. <i>arXiv preprint</i>		
618	<i>arXiv:2306.00890</i> .	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> . <i>ArXiv preprint</i> ,	671
		<a href="#">abs/2303.08774</a> .	672
619	Guohao Li, Hasan Abed Al Kader Hammoud, Hani	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan	673
620	Itani, Dmitrii Khizbullin, and Bernard Ghanem.	Sankarasubbu. 2022. Medmcqa: A large-scale multi-	674
621	2023b. Camel: Communicative agents for" mind"	subject multi-choice dataset for medical domain	675
622	exploration of large scale language model society.	question answering. In <i>Conference on Health,</i>	676
623	<i>arXiv preprint arXiv:2303.17760</i> .	<i>Inference, and Learning</i> , pages 248–260. PMLR.	677
624	Huaoli, Yu Quan Chong, Simon Stepputtis, Joseph	Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai,	678
625	Campbell, Dana Hughes, Michael Lewis, and Katia	Meredith Ringel Morris, Percy Liang, and Michael S.	679
626	Sycara. 2023c. Theory of mind for multi-agent	Bernstein. 2023. Generative agents: Interactive	680
627	collaboration via large language models. <i>arXiv</i>	simulacra of human behavior. In <i>In the 36th Annual</i>	681
628	<i>preprint arXiv:2310.10701</i> .	<i>ACM Symposium on User Interface Software and</i>	682
629	Yuan Li, Yixuan Zhang, and Lichao Sun. 2023d.	<i>Technology (UIST ’23)</i> , UIST ’23, New York, NY,	683
630	<a href="#">Metaagents: Simulating interactions of human</a>	USA. Association for Computing Machinery.	684
631	<a href="#">behaviors for llm-based task-oriented coordination</a>		
632	<a href="#">via collaborative generative agents</a> .	Chen Qian, Xin Cong, Cheng Yang, Weize Chen,	685
633	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong	686
634	Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and	Sun. 2023. Communicative agents for software	687
635	Shuming Shi. 2023. <a href="#">Encouraging divergent thinking</a>	development. <i>arXiv preprint arXiv:2307.07924</i> .	688
636	<a href="#">in large language models through multi-agent debate</a> .		
637	Valentin Liévin, Christoffer Egeberg Hother, and	Maciej Rosoł, Jakub S Gąsior, Jonasz Łaba,	689
638	Ole Winther. 2022. Can large language models	Kacper Korzeniewski, and Marcel Młyńczak. 2023.	690
639	reason about medical questions? <i>arXiv preprint</i>	Evaluation of the performance of gpt-3.5 and gpt-4	691
640	<i>arXiv:2207.08143</i> .	on the medical final examination. <i>medRxiv</i> , pages	692
		2023–06.	693
641	Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng,	Teven Le Scao, Angela Fan, Christopher Akiki,	694
642	Can Zheng, Junxiang Wang, Tanmoy Chowdhury,	Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman	695
643	Yun Li, Hejie Cui, Tianjiao Zhao, et al. 2023.	Castagné, Alexandra Sasha Luccioni, François Yvon,	696
644	Beyond one-model-fits-all: A survey of domain	Matthias Gallé, et al. 2022. <a href="#">Bloom: A 176b-</a>	697
645	specialization for large language models. <i>arXiv</i>	<a href="#">parameter open-access multilingual language model</a> .	698
646	<i>preprint arXiv:2305.18703</i> .	<i>ArXiv preprint</i> , <a href="#">abs/2211.05100</a> .	699
647	Zhengliang Liu, Zihao Wu, Mengxuan Hu, Bokai Zhao,	Henk G Schmidt and Remy MJP Rikers. 2007.	700
648	Lin Zhao, Tianyi Zhang, Haixing Dai, Xianyan Chen,	How expertise develops in medicine: knowledge	701
649	Ye Shen, Sheng Li, et al. 2023. Pharmacygpt: The ai	encapsulation and illness script formation. <i>Medical</i>	702
650	pharmacist. <i>arXiv preprint arXiv:2307.10432</i> .	<i>education</i> , 41(12):1133–1139.	703

704	Chantal Shaib, Millicent L Li, Sebastian Joseph, Iain J	U.S. Department of Health and Human Services.	759
705	Marshall, Junyi Jessy Li, and Byron C Wallace. 2023.	1996. The hipaa privacy rule. <a href="https://www.hhs.gov/hipaa/for-professionals/privacy/index.html">https://www.hhs.gov/hipaa/for-professionals/</a>	760
706	Summarizing, simplifying, and synthesizing medical	<a href="https://www.hhs.gov/hipaa/for-professionals/privacy/index.html">privacy/index.html</a> .	761
707	evidence using gpt-3 (with varying success). <i>arXiv</i>		762
708	<i>preprint arXiv:2305.06299</i> .		
709	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Mahdavi,	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	763
710	Jason Wei, Hyung Chung, Nathan Scales, Ajay	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	764
711	Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry	Xu Chen, Yankai Lin, et al. 2023a. A survey on large	765
712	Payne, Martin Seneviratne, Paul Gamble, Chris	language model based autonomous agents. <i>arXiv</i>	766
713	Kelly, Abubakr Babiker, Nathanael SchÅd'rli,	<i>preprint arXiv:2308.11432</i> .	767
714	Aakanksha Chowdhery, Philip Mansfield, Dina	Zekun Moore Wang, Zhongyuan Peng, Haoran Que,	768
715	Demner-Fushman, and Vivek Natarajan. 2023a.	Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu,	769
716	<a href="#">Large language models encode clinical knowledge</a> .	Hongcheng Guo, Ruitong Gan, Zehao Ni, Man	770
717	<i>Nature</i> , 620:1–9.	Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu,	771
718	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres,	Wenhu Chen, Jie Fu, and Junran Peng. 2023b.	772
719	Ellery Wulczyn, Le Hou, Kevin Clark, Stephen	Rolellm: Benchmarking, eliciting, and enhancing	773
720	Pfohl, Heather Cole-Lewis, Darlene Neal, et al.	role-playing abilities of large language models. <i>arXiv</i>	774
721	2023b. Towards expert-level medical question	<i>preprint arXiv: 2310.00746</i> .	775
722	answering with large language models. <i>arXiv</i>	Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao	776
723	<i>preprint arXiv:2305.09617</i> .	Ge, Furu Wei, and Heng Ji. 2023c. <a href="#">Unleashing</a>	777
724	Yang Tan, Mingchen Li, Zijie Huang, Huiqun Yu, and	<a href="#">cognitive synergy in large language models: A</a>	778
725	Guisheng Fan. 2023. Medchatz: a better medical	<a href="#">task-solving agent through multi-persona self-</a>	779
726	adviser learns from better instructions. <i>arXiv preprint</i>	<a href="#">collaboration</a> .	780
727	<i>arXiv:2309.01114</i> .		
728	Liyan Tang, Zhaoyi Sun, Betina Idnay, Jordan G	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran	781
729	Nestor, Ali Soroush, Pierre A Elias, Ziyang Xu,	Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,	782
730	Ying Ding, Greg Durrett, Justin F Rousseau, et al.	Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023a.	783
731	2023a. Evaluating large language models on medical	Autogen: Enabling next-gen llm applications via	784
732	evidence summarization. <i>npj Digital Medicine</i> ,	multi-agent conversation framework. <i>arXiv preprint</i>	785
733	6(1):158.	<i>arXiv:2308.08155</i> .	786
734	Xiangru Tang, Arman Cohan, and Mark Gerstein. 2023b.	Yiquan Wu, Siying Zhou, Yifei Liu, Weiming Lu,	787
735	Aligning factual consistency for clinical studies	Xiaozhong Liu, Yating Zhang, Changlong Sun, Fei	788
736	summarization through reinforcement learning. In	Wu, and Kun Kuang. 2023b. <a href="#">Precedent-enhanced</a>	789
737	<i>Proceedings of the 5th Clinical Natural Language</i>	<a href="#">legal judgment prediction with llm and domain-</a>	790
738	<i>Processing Workshop</i> , pages 48–58.	<a href="#">model collaboration</a> .	791
739	Xiangru Tang, Andrew Tran, Jeffrey Tan, and Mark	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen	792
740	Gerstein. 2023c. Gersteinlab at medqa-chat 2023:	Ding, Boyang Hong, Ming Zhang, Junzhe Wang,	793
741	Clinical note summarization from doctor-patient	Senjie Jin, Enyu Zhou, et al. 2023. The rise and	794
742	conversations through fine-tuning and in-context	potential of large language model based agents: A	795
743	learning. <i>arXiv preprint arXiv:2305.05001</i> .	survey. <i>arXiv preprint arXiv:2309.07864</i> .	796
744	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	Tianbao Xie, Fan Zhou, Zhoujun Cheng, Peng Shi,	797
745	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	Luoxuan Weng, Yitao Liu, Toh Jing Hua, Junning	798
746	Baptiste Rozière, Naman Goyal, Eric Hambro,	Zhao, Qian Liu, Che Liu, et al. 2023. Openagents:	799
747	Faisal Azhar, et al. 2023. <a href="#">Llama: Open and</a>	An open platform for language agents in the wild.	800
748	<a href="#">efficient foundation language models</a> . <i>ArXiv preprint</i> ,	<i>arXiv preprint arXiv:2310.10634</i> .	801
749	<a href="#">abs/2302.13971</a> .	Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing	802
750	Tao Tu, Shekoofeh Azizi, Danny Driess, Mike	Qin. 2023. Examining the inter-consistency of large	803
751	Schaeckermann, Mohamed Amin, Pi-Chuan Chang,	language models: An in-depth analysis via debate.	804
752	Andrew Carroll, Chuck Lau, Ryutaro Tanno, Ira	<i>arXiv e-prints</i> , pages arXiv–2305.	805
753	Ktena, et al. 2023. Towards generalist biomedical ai.	Yi Yang, Yixuan Tang, and Kar Yan Tam. 2023.	806
754	<i>arXiv preprint arXiv:2307.14334</i> .	<a href="#">Investlm: A large language model for investment</a>	807
755	Logesh Kumar Umapathi, Ankit Pal, and Malaikannan	<a href="#">using financial domain instruction tuning</a> .	808
756	Sankarasubbu. 2023. Med-halt: Medical domain	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	809
757	hallucination test for large language models. <i>arXiv</i>	Shafraan, Karthik Narasimhan, and Yuan Cao. 2022.	810
758	<i>preprint arXiv:2307.15343</i> .	React: Synergizing reasoning and acting in language	811
		models. <i>arXiv preprint arXiv:2210.03629</i> .	812

Cyril Zakka, Akash Chaurasia, Rohan Shad, Alex R Dalal, Jennifer L Kim, Michael Moor, Kevin Alexander, Euan Ashley, Jack Boyd, Kathleen Boyd, et al. 2023. Almanac: Retrieval-augmented language models for clinical medicine. *Research Square*.

Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023a. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*.

Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold. 2023b. [Alpacare:instruction-tuned large language models for medical application](#).

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. 2023. [Webarena: A realistic web environment for building autonomous agents](#).

## A Prompt Templates

Prompt templates involved in the experiments are presented in Table 5:



Table 5: Prompt templates and role descriptions employed in our MC framework.

---

<b>r<sub>qd</sub></b> :	You are a medical expert who specializes in categorizing a specific medical scenario into specific areas of medicine.
<b>prompt<sub>qd</sub></b> :	You need to complete the following steps: 1. Carefully read the medical scenario presented in the question: question. 2. Based on the medical scenario in it, classify the question into five different subfields of medicine. 3. You should output in the same format as: Medical Field:   .
<b>r<sub>od</sub></b> :	As a medical expert, you possess the ability to discern the two most relevant fields of expertise needed to address a multiple-choice question encapsulating a specific medical context.
<b>prompt<sub>od</sub></b> :	You need to complete the following steps: 1. 1. Carefully read the medical scenario presented in the question: question. 2. The available options are: options. Strive to understand the fundamental connections between the question and the options. 3. Your core aim should be to categorize the options into two distinct subfields of medicine. You should output in the same format as: Medical Field:   .
<b>r<sub>qa</sub></b> :	You are a medical expert in the domain of question_domain. From your area of specialization, you will scrutinize and diagnose the symptoms presented by patients in specific medical scenarios.
<b>prompt<sub>qa</sub></b> :	Please meticulously examine the medical scenario outlined in this question: question. Drawing upon your medical expertise, interpret the condition being depicted. Subsequently, identify and highlight the aspects of the issue that you find most alarming or noteworthy.
<b>r<sub>oa</sub></b> :	You are a medical expert specialized in the op_domain domain. You are adept at comprehending the nexus between questions and choices in multiple-choice exams and determining their validity. Your task, in particular, is to analyze individual options with your expert medical knowledge and evaluate their relevancy and correctness.
<b>prompt<sub>oa</sub></b> :	Regarding the question: question, we procured the analysis of five experts from diverse domains. The evaluation from the question_domain expert suggests: question_analysis. The following are the options available: options. Reviewing the question’s analysis from the expert team, you’re required to fathom the connection between the options and the question from the perspective of your respective domain and scrutinize each option individually to assess whether it is plausible or should be eliminated based on reason and logic. Pay close attention to discerning the disparities among the different options and rationalize their existence. A handful of these options might seem right at first glance but could potentially be misleading in reality.
<b>r<sub>rs</sub></b> :	You are a medical assistant who excels at summarizing and synthesizing based on multiple experts from various domain experts.
<b>prompt<sub>rs</sub></b> :	Here are some reports from different medical domain experts. You need to complete the following steps: 1. Take careful and comprehensive consideration of the following reports. 2. Extract key knowledge from the following reports. 3. Derive the comprehensive and summarized analysis based on the knowledge. 4. Your ultimate goal is to derive a refined and synthesized report based on the following reports. You should output in exactly the same format as: Key Knowledge;; Total Analysis:
<b>prompt<sub>mod</sub></b> :	Here is advice from a medical expert specialized in domain: advice. Based on the above advice, output the revised analysis in the same format as: Key Knowledge;; Total Analysis:
<b>prompt<sub>dm</sub></b> :	Here is a synthesized report: syn_report. Based on the above report, select the optimal choice to answer the question. Points to note: 1. The analyses provided should guide you towards the correct response. 2. Any option containing incorrect information inherently cannot be the correct choice. 3. Please respond only with the selected option’s letter, like A, B, C, D, or E, using the following format: ’’Option: [Selected Option’s Letter]’’. Remember, it’s the letter we need, not the full content of the option.

---