

# MED-MAD: Breaking Medical Mental Set with Mindset-Diversified Multi-Agent Debate

Anonymous ACL submission

## Abstract

Large language models exhibit strong performance in medical question answering yet remain vulnerable to persistence bias, wherein early diagnostic hypotheses anchor subsequent reasoning and induce correlated, overconfident errors. We conceptualize this failure mode as the medical mental set and quantify its impact by measuring trajectory collapse using answer agreement and rationale similarity across repeated inferences. Prior work inadequately addresses this bias, lacking quantitative diagnostics, failing to promote hypothesis diversity, and relying on structurally homogeneous multi-agent frameworks that are insufficient to mitigate correlated diagnostic errors. To address these challenges, we introduce MED-MAD, an anonymous multi-agent framework designed to enhance hypothesis-level diversity while ensuring auditability. The framework incorporates mindset-specialized clinical roles, structured anonymous debate emphasizing safety-oriented critique and concession, optional backbone heterogeneity, and confidence-aware aggregation reinforced by counterfactual consistency checks. These components work together to promote rigorous and diverse diagnostic reasoning while minimizing correlated errors. Experimental evaluations demonstrate that MED-MAD consistently improves diagnostic accuracy and reduces correlated errors across benchmark datasets, under matched inference budgets. These findings highlight its potential to support safer and more reliable clinical reasoning, advancing the role of LLMs in medical decision support.

## 1 Introduction

Large language models (LLMs) have shown strong performance on medical exam-style benchmarks such as MedQA (Jin et al., 2021) and other comprehensive medical QA suites

(Singhal et al., 2023; Saab et al., 2024), indicating their potential for medical reasoning in controlled settings. However, these LLM-based systems can suffer from persistence bias; when an initial hypothesis is incorrect, the models often remain fixated on the original diagnostic path, producing variations of the same flawed reasoning (Wang et al., 2025; Liu et al., 2025b; Huang et al., 2024; Valmeekam et al., 2023). This premature closure in differential diagnosis is particularly problematic, as safety-oriented decision-making requires the consideration and evaluation of plausible alternative hypotheses (Vázquez-Costa and Costa-Alcaraz, 2013).

To address the challenges of persistence bias in LLMs, common strategies involve iterative prompting and decoding techniques, such as chain-of-thought and self-consistency (Wei et al., 2022; Wang et al., 2023), along with iterative self-correction methods (Madaan et al., 2023; Shinn et al., 2023; Huang et al., 2023). However, without incorporating new information or external feedback, these repeated attempts often fail to alter the initial hypothesis and may even reinforce previous errors (Wang et al., 2025; Liu et al., 2025b; Huang et al., 2024; Stechly et al., 2023; Valmeekam et al., 2023). Multi-agent debate (MAD) offers a potential solution by enabling agents to critique and revise each other’s reasoning (Du et al., 2024; Liang et al., 2024; Chan et al., 2024; Liu et al., 2025b). However, many frameworks still rely on agents derived from the same model with similar prompts, resulting in homogeneous reasoning and correlated errors. This limitation is particularly significant in the medical field, where differential diagnosis is most effective when it incorporates diverse clinical perspectives rather than minor variations in agent personas (Tang et al., 2024; Kim et al., 2024; Chen et al., 2024; Hong and Page, 2004).

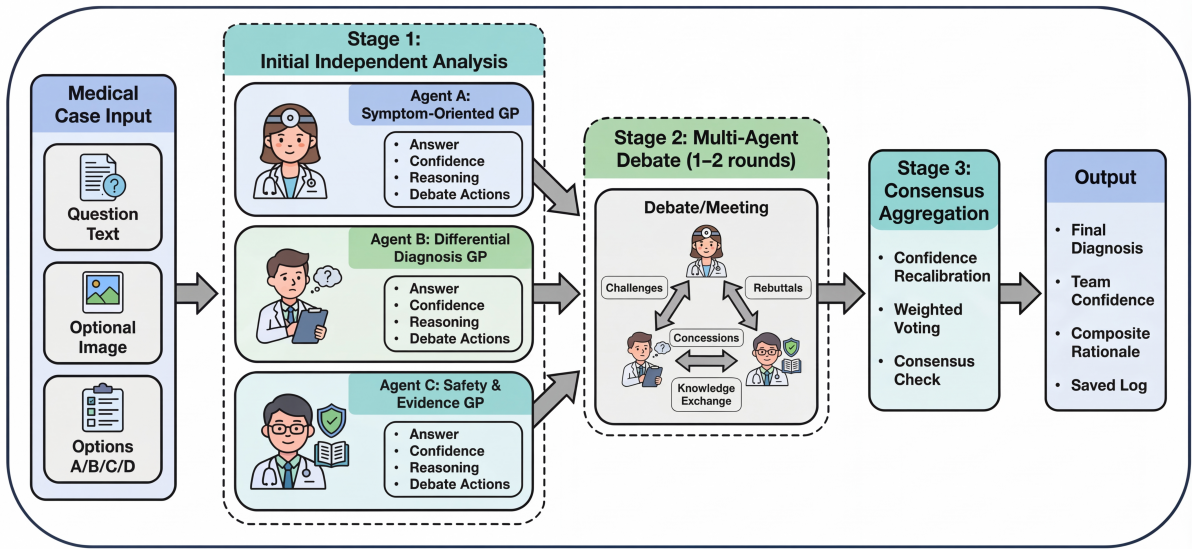


Figure 1: Overview of MED-MAD. Role-specialized GP agents independently solve a medical MCQ (Stage 1), optionally debate via anonymous critique-and-revision (Stage 2), and are aggregated via confidence-recalibrated weighted voting and consensus checks (Stage 3) to produce a final answer, team confidence, and an auditable log.

This pattern is reminiscent of the psychological concept known as a "mental set" (Jersild, 1927; Luchins, 1942; Ollinger et al., 2008; DeCaro, 2016), where individuals persistently apply familiar strategies even when they are not optimal. In our context, this manifests as a "medical mental set" effect, where repeated attempts or interactions among multiple agents result in increasingly similar rationales and early convergence on answers, occasionally leading to the selection of the same incorrect option. We term this phenomenon "trajectory collapse" and quantify it through measures of answer agreement and rationale similarity, as illustrated in Figure 3.

To address these challenges, we introduce MED-MAD (Mindset-Diversified Medical Reasoning via Role-Specialized Anonymous Multi-Agent Debate), a framework designed to enhance hypothesis-level diversity while maintaining auditability of interactions, as depicted in Figure 1. MED-MAD employs three specialized clinical approaches: phenotype-first (pattern matching), hypothesis-driven (hypothetico-deductive differential diagnosis), and safety-first (error-aware risk checking). This structure fosters complementary reasoning styles, drawing from established models

of diagnostic reasoning (Elstein et al., 1978; Croskerry, 2009b,a; Pelaccia et al., 2011; Norman et al., 2024).

The framework utilizes an anonymous, structured debate protocol, incorporating explicit critique and concession cues along with safety checks to promote precise revisions (Croskerry, 2003; Croskerry et al., 2013a,b; Mamede et al., 2014; Croskerry and Campbell, 2021; Ke et al., 2024). Additionally, MED-MAD can optionally incorporate backbone heterogeneity to further minimize correlated biases among agents (Hong and Page, 2004). Each round produces a unified, structured reasoning record (including confidence and a counterfactual check (Richens et al., 2020)) that makes disagreement concrete and traceable.

We evaluate MED-MAD on three medical QA benchmarks with controlled ablations, trajectory-collapse diagnostics, and a budget-matched comparison (Table 4). These studies clarify when debate improves accuracy and how role vs. backbone diversity shapes correlated errors (Figure 3).

**Contributions.** (1) We empirically characterize the *medical mental set* effect in medical QA as *trajectory collapse*, quantified via answer

140 agreement and rationale similarity. (2) We  
141 introduce MED-MAD, a mindset-diversified  
142 multi-agent debate framework with anonymous  
143 interaction and unified, auditable outputs to  
144 promote hypothesis-level diversity. (3) We con-  
145 duct controlled experiments and ablations to  
146 isolate the effects of role specialization, debate  
147 protocol, and backbone heterogeneity, and an-  
148alyze their impact on accuracy and trajectory  
149 collapse.

150 <sup>1</sup>

## 151 2 Related Work

152 **LLMs for medical QA.** LLMs have demon-  
153 strated impressive performance on medical  
154 exam-style benchmarks (Jin et al., 2021; Sing-  
155 hal et al., 2023), which has led to their con-  
156 sideration as tools for clinical decision support.  
157 Recent large-scale evaluations of LLM capabili-  
158 ties in the medical domain, such as Med-PaLM  
159 (Singhal et al., 2023) and Med-Gemini (Saab  
160 et al., 2024), reveal that these models can ex-  
161 hibit vulnerabilities in multi-step clinical rea-  
162 soning. In such cases, errors may be systematic  
163 rather than isolated. This vulnerability is rem-  
164iniscent of diagnostic decision-making failures  
165 like premature diagnostic closure, where initial  
166 commitments limit the exploration of alter-  
167 native hypotheses (Vázquez-Costa and Costa-  
168 Alcaraz, 2013).

### 169 **Agentic and multi-agent medical systems.**

170 To more accurately represent collaborative clin-  
171 ical workflows, recent research has begun to  
172 employ multiple LLM agents with complemen-  
173 tary roles. For instance, MEDAGENTS (Tang  
174 et al., 2024) utilizes role-specific discussions to  
175 enhance medical reasoning, while MDAGENTS  
176 (Kim et al., 2024) adapts collaborative strate-  
177 gies based on the complexity of medical cases.  
178 Beyond merely answering questions, simulation  
179 environments such as AGENT HOSPITAL (Li  
180 et al., 2024) model comprehensive care pro-  
181 cesses involving interacting agents. For evalu-  
182 ating these multi-agent systems, MEDAGENT-  
183 BOARD (Zhu et al., 2025) serves as a bench-  
184 mark for assessing collaborative efforts across  
185 various medical tasks, facilitating more con-  
186 trolled comparisons of different agent protocols.  
187 More broadly, general-purpose frameworks for

188 multi-agent orchestration, along with agent  
189 evaluation suites, have advanced the study of  
190 interaction patterns beyond the medical field  
191 (Li et al., 2023; Wu et al., 2023; Liu et al.,  
192 2025a).

193 **Multi-agent debate, diversity, and cor-  
194 related errors.** Multi-agent debate (MAD)  
195 enhances reasoning by enabling agents to cri-  
196 tique and refine each other’s outputs (Du et al.,  
197 2024). For example, RECONCILE (Chen et al.,  
198 2024) approaches debate as an iterative process  
199 of achieving consensus among diverse LLMs,  
200 while other variations, such as iMAD (Fan  
201 et al., 2025), explore efficiency improvements  
202 and protocol designs. Some research also ap-  
203 plies MAD frameworks for evaluation purposes  
204 (Chan et al., 2024). However, a common limita-  
205 tion is the lack of diversity among agents, when  
206 instantiated from similar models or prompts,  
207 agents often converge too quickly and repeat  
208 correlated errors, undermining the benefits of  
209 interaction (Liu et al., 2025b; Liang et al.,  
210 2024). Our work aligns with research that ex-  
211 plicitly aims to foster diversity in multi-agent  
212 debates (Liu et al., 2025b) but focuses specifi-  
213 cally on the medical domain. In this context,  
214 premature convergence can result in a phe-  
215 nomenon we term “medical mental set,” where  
216 agents commit early to flawed reasoning path-  
217 ways, leading to high-confidence, correlated  
218 mistakes.

## 219 3 Method

220 This section introduces MED-MAD, a  
221 coordinator-driven multi-agent framework de-  
222 signed for multiple-choice medical question an-  
223 swering. The framework operates in three dis-  
224 tinct phases: (1) independent initial analysis  
225 by each agent, (2) anonymous multi-round de-  
226 bate among the agents, and (3) collaborative  
227 team-based answer generation. The workflow  
228 is illustrated in Figures 1 and 2. For compre-  
229 hensive implementation details, including end-  
230 to-end pseudocode, please refer to Appendix B  
231 (Algorithm 1).

232 **Framework overview.** In response to  
233 a medical question  $q$  with answer options  
234  $op = \{o_1, o_2, \dots, o_k\}$ , MED-MAD deploys  
235 three mindset-specialized general practitioner  
236 (GP) agents:  $\mathcal{A} = \{A_{\text{sym}}, A_{\text{ddx}}, A_{\text{guide}}\}$ . These  
237 agents represent different diagnostic perspec-

<sup>1</sup>An anonymized repository is available at <https://anonymous.4open.science/r/med3-9640>.

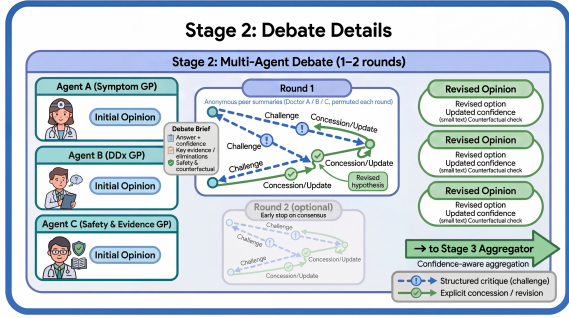


Figure 2: Stage-2 details of MED-MAD. Agents exchange anonymized peer summaries and iteratively challenge (dashed arrows) or concede/update (solid arrows) before forwarding revised answers and confidences to aggregation; agent identifiers are permuted each round.

238 tives: phenotype-first, hypothesis-driven, and  
 239 safety-first. The system operates through a  
 240 coordinator loop. Let  $y_i^{(t)}$  represent the struc-  
 241 tured record produced by agent  $A_i$  during  
 242 round  $t$ , where  $t = 0$  indicates the initial assess-  
 243 ment. The coordinator maintains a **history** of  
 244 all records from each round. In our implemen-  
 245 tation, each  $y_i^{(t)}$  adheres to a unified schema,  
 246 which includes the option letter, confidence  
 247 level, fit scores for each option, structured rea-  
 248 soning highlighting key clues and eliminations,  
 249 and optional counterfactual or safety notes. This  
 250 standardization ensures that constructing deba-  
 251 tes and aggregating the results are  
 252 straightforward and deterministic.

253 **Phase 1: Initial analysis.** In the first  
 254 phase, the coordinator engages all agents sim-  
 255 ultaneously to gather their initial outputs,  
 256 denoted as  $\{y_i^{(0)}\}$ . After standardizing the op-  
 257 tion labels, if all agents reach the same conclu-  
 258 sion, MED-MAD bypasses the second phase  
 259 and advances directly to Phase 3.

260 **Phase 2: Anonymous multi-round de-**  
 261 **bate.** If the agents do not agree, the coordi-  
 262 nator initiates an anonymous debate, which  
 263 can last for a maximum of  $T_{\max}$  rounds. Begin-  
 264 ning at round  $t \geq 1$ , the coordinator creates  
 265 a debate brief  $B^{(t)}$ , derived from the latest  
 266 records of each agent  $\{y_j^{(t-1)}\}_{j=1}^3$ . This brief,  
 267 formatted as a text prompt (`debate_brief`),  
 268 is then distributed to all agents to facilitate  
 269 discussion, as illustrated in Figure 2.

270 (a) *Build the debate brief.* For each peer,  $B^{(t)}$   
 271 provides a summary that includes three main  
 272 components: (i) the current answer along with

273 the level of confidence, (ii) a concise overview  
 274 of the highest-scoring options, and (iii) a sum-  
 275 mary of evidence highlighting key support-  
 276 ing and opposing elements. This summary  
 277 is further enhanced with a safety risk check  
 278 and a counterfactual stress test (Weng et al.,  
 279 2022). The summary concludes with precise  
 280 critique instructions. Agents are expected to  
 281 issue structured challenges to specific reason-  
 282 ing steps, as depicted by the dashed arrows in Figure 2. Additionally, they should be prepared to explicitly revise their hypothesis when convinced by the argument, as indicated by the solid arrows.

287 (b) *Enforce anonymity.*  $B^{(t)}$  assigns  
 288 anonymized identifiers, such as Doctor A, B,  
 289 and C, to each peer. To prevent biases related  
 290 to authority attribution and anchoring to spe-  
 291 cific individuals, these identifiers are randomly  
 292 reassigned in every round (Milgram, 1963; Tver-  
 293 sky and Kahneman, 1974). Mindset labels and  
 294 model architectures are intentionally excluded,  
 295 ensuring that peers are evaluated solely based  
 296 on their anonymized summaries. This iterative  
 297 and anonymous revision process resembles the  
 298 Delphi method for expert elicitation (Rowe and  
 299 Wright, 1999).

300 (c) *Revise and stop early.* Based on  $B^{(t)}$ , the co-  
 301 ordinator simultaneously tasks all agents with  
 302 generating updated records  $y_i^{(t)}$ . Each record  
 303 comprises a revised option, updated confidence  
 304 level, and optionally, a brief counterfactual  
 305 check. The debate concludes early if the nor-  
 306 malized answers align. If not, the discussion  
 307 continues until it reaches  $T_{\max}$ . The final re-  
 308 visited opinions are then forwarded to Phase 3  
 309 for aggregation, taking into account the vary-  
 310 ing levels of confidence.

311 **Phase 3: Team answer generation.**  
 312 Following the debate, MED-MAD employs  
 313 confidence-recalibrated weighted voting to de-  
 314 rive the final answer  $a^*$  and the team’s col-  
 315 lective confidence  $c_{\text{team}}$  (Lin and Hooi, 2025).  
 316 Subsequently, it generates a final rationale  $R^*$ ,  
 317 drawing on the structured evidence presented  
 318 by the agents that supports  $a^*$ . For further  
 319 information on agent roles, output schema and  
 320 normalization, debate brief construction, and  
 321 aggregation processes, please refer to Appendix  
 322 B. Prompt templates can be found in Appendix  
 323 C.

## 4 Experiments

### 4.1 Experimental Setup

#### 4.1.1 Datasets

In alignment with previous research (Zhu et al., 2025), we assess performance using medical multiple-choice benchmarks, serving as a controlled proxy for exam-style clinical reasoning. Our evaluation covers three multiple-choice medical question-answering (MCQ) benchmarks: MedQA (Jin et al., 2021), PubMedQA (Jin et al., 2019), and the professional medicine subset of MMLU-Pro (Hendrycks et al., 2021; Wang et al., 2024). We focus on the official test sets of MedQA and PubMedQA.

Given the inference costs associated with multi-agent debate and the limitations of closed-source API rate limits, prior work often uses a fixed subset rather than the full test sets for evaluation (e.g., Du et al., 2024; Besta et al., 2024; Yao et al., 2023; Chen et al., 2024; Tang et al., 2024). Following this approach, we randomly select 100 instances per dataset and maintain this subset consistently across all methods and repeated runs within each benchmark.

To account for run-to-run variability due to stochastic decoding and multi-round agent interactions within a fixed cost budget, we repeat each configuration a minimum of three times. We report accuracy as the mean  $\pm$  sample standard deviation across these runs (ddof=1; see Appendix B.7 for more details). Notably, in some configurations, we observe a standard deviation of 0.0000 across runs, which can happen when runs produce identical accuracies with fixed prompts and evaluation subsets.

#### 4.1.2 Backbone Models and Compared Methods

We assess MED-MAD and various baseline models under two different backbone panel configurations, which define the set of models used as agents. The first configuration utilizes an open-source MCQ panel consisting of Gemma3-4B, LLaMA3.2-3B, and Qwen3-4B (Gemma Team et al., 2025; Grattafiori et al., 2024; Meta, 2024; Yang et al., 2025). The second configuration involves a closed-source API panel with DeepSeek-v3.2-exp, Grok-4-Fast, and GPT-4.1-Nano (DeepSeek Team, 2025; xAI, 2025; Ope-

nAI, 2025). For baseline comparisons involving single-model prompting, we apply Gemma3-4B in the open-source setting and DeepSeek-v3.2-exp in the closed-source setting. Unless otherwise indicated, our analyses concentrate on the open-source heterogeneous panel, enabling controlled and reproducible evaluations.

We benchmark MED-MAD against standard single-model prompting techniques, such as few-shot prompting (Brown et al., 2020), chain-of-thought (CoT) (Wei et al., 2022), self-consistency (SC) (Wang et al., 2023), and CoT with self-consistency (CoT+SC) (Wang et al., 2023). Additionally, we include comparisons with notable agentic systems like HealthcareAgent (Ren et al., 2025), MedAgents (Tang et al., 2024), MDAgents (Kim et al., 2024), and ReConcile (Chen et al., 2024). To ensure a fair evaluation, these multi-agent frameworks are implemented using the same heterogeneous panels as MED-MAD.

#### 4.1.3 Evaluation Protocol

Our evaluation framework is based on the open-source MedAgentBoard codebase (Zhu et al., 2025), enhanced to enable multi-round debates. In assessing the multiple-choice benchmarks, we ensure consistency by evaluating all methods on the same pre-selected test subsets for each dataset, within both open-source and closed-source backbone settings. Each method receives identical inputs, and we do not engage in additional training, fine-tuning, or retrieval.

Decoding configurations such as temperature and top- $p$ /top- $k$  are used as released by each backbone’s provider, without further adjustment. We avoid standardizing a single decoding configuration across different model families because the same nominal settings can result in varying levels of randomness, potentially affecting performance negatively for some models (Guo et al., 2017; Holtzman et al., 2020). Detailed configuration information is available in Appendix B.10.

In the MED-MAD approach, each agent performs an independent initial pass, followed by up to  $T_{\max}$  debate rounds. The process stops early if all agents agree on a normalized option letter. For open-source backbones, we set  $T_{\max}=3$ , while for closed-source APIs,  $T_{\max}=2$  is used, considering cost constraints.

Method	MedQA	PubMedQA	MMLU-Pro
<b>Open-Source Models</b>			
<i>Single LLM Baselines (Backbone: Gemma3-4B)</i>			
Few-shot	0.3200 ± 0.0100	0.6200 ± 0.0000	0.3633 ± 0.0058
Self-Consistency ( $K=5$ )	0.3367 ± 0.0058	0.6133 ± 0.0058	0.3733 ± 0.0058
Chain-of-Thought (CoT)	0.4167 ± 0.0115	0.6967 ± 0.0058	0.4000 ± 0.0693
CoT + SC ( $K=3$ )	0.3900 ± 0.0100	0.6900 ± 0.0200	0.4333 ± 0.0379
<i>Single-Agent Framework</i>			
HealthcareAgent	0.3000 ± 0.0173	0.4900 ± 0.0265	0.3067 ± 0.0404
<i>Multi-Agent Frameworks (Heterogeneous Panel)</i>			
MedAgents	0.5033 ± 0.0379	0.7633 ± 0.0321	0.4900 ± 0.0100
MDAgents	0.5167 ± 0.0153	<b>0.7733 ± 0.0115</b>	0.5700 ± 0.0265
ReConcile	0.5800 ± 0.0529	0.7600 ± 0.0100	0.5667 ± 0.0569
<b>MED-MAD (Ours)</b>	<b>0.6000 ± 0.0265</b>	0.7650 ± 0.0071	<b>0.5967 ± 0.0289</b>
<b>Closed-Source or API Models</b>			
<i>Single LLM Baselines (Backbone: DeepSeek-v3.2-exp)</i>			
Few-shot	0.7800 ± 0.0173	0.7867 ± 0.0208	0.8233 ± 0.0115
Self-Consistency ( $K=3$ )	0.7667 ± 0.0153	<b>0.8333 ± 0.0153</b>	0.8500 ± 0.0000
Chain-of-Thought (CoT)	0.8567 ± 0.0503	0.7667 ± 0.0153	0.8767 ± 0.0306
CoT + SC ( $K=3$ )	0.8933 ± 0.0306	0.7700 ± 0.0100	0.8700 ± 0.0173
<i>Single-Agent Framework</i>			
HealthcareAgent	0.8733 ± 0.0321	0.7733 ± 0.0451	0.8500 ± 0.0173
<i>Multi-Agent Frameworks (Heterogeneous Panel)</i>			
MedAgents	0.8500 ± 0.0361	0.8267 ± 0.0058	0.8700 ± 0.0094
MDAgents	0.7300 ± 0.0265	0.7867 ± 0.0351	0.6267 ± 0.0289
ReConcile	0.8667 ± 0.0058	0.8167 ± 0.0115	0.8833 ± 0.0115
<b>MED-MAD (Ours)</b>	<b>0.9067 ± 0.0153</b>	0.8200 ± 0.0100	<b>0.8900 ± 0.0100</b>

Table 1: Accuracy (mean ± sample std. across runs; ddof= 1) on three multiple-choice medical QA benchmarks, grouped by backbone setting (open-source models vs. closed-source APIs). Best score per dataset within each group is bold.

#### 4.1.4 Metrics

The main metric we use for evaluation is multiple-choice accuracy (Acc), which is determined from the final team prediction using confidence-recalibrated weighted voting (see Appendix B.8 for details). We evaluate accuracy based on the exact match of the predicted option letter. Further details on the scoring process can be found in Appendix B.7. To facilitate reproducibility of our results, we will make available the precise list of question IDs (qids), along with the code and evaluation scripts.

## 4.2 Main Results

### 4.2.1 Overall Performance on Medical QA

Table 1 presents the accuracy results for MedQA, PubMedQA, and MMLU-Pro, categorized by backbone setting. In the open-source configuration, MED-MAD achieves an accuracy of 0.6000 on MedQA and 0.5967 on MMLU-Pro, securing the highest mean accuracy among the multi-agent methods evaluated

on these datasets (Table 1). Compared to the top-performing multi-agent baselines, paired bootstrap comparisons on an identical fixed subset reveal a  $\Delta\text{Acc}$  of +0.0200 on MedQA when compared to ReConcile, with a 95% confidence interval (CI) of  $[-0.0300, +0.0700]$ , and +0.0267 on MMLU-Pro compared to MDAgents, with a 95% CI of  $[-0.0433, +0.1000]$  (see Appendix Table A1 for more details). On PubMedQA, MED-MAD achieves an accuracy of 0.7650, with a paired comparison against MDAgents showing a  $\Delta\text{Acc}$  of  $-0.0033$ , within a 95% CI of  $[-0.0667, +0.0567]$  (refer to Appendix Table A1).

In the closed-source API configuration, MED-MAD attains an accuracy of 0.9067 on MedQA and 0.8900 on MMLU-Pro, which is comparable to or slightly exceeds the leading baselines in this setting (CoT+SC records 0.8933 on MedQA, and ReConcile achieves 0.8833 on MMLU-Pro; see Table 1). We present these results as a verification check using external backbones. Our subsequent analysis primarily concentrates on open-source models,

where we have greater control over configurations and can thoroughly examine intermediate traces.

PubMedQA serves as a contrastive benchmark: single-model baselines already show strong performance, and the multi-agent methods cluster closely within a range of 0.7600 to 0.7733. This clustering suggests limited potential for further improvement through multi-round interaction and specialized debate roles, as illustrated in Table 1. Given that each dataset is evaluated using a fixed subset of 100 questions, our emphasis is on paired comparisons and the estimation of uncertainty (further details can be found in Appendix Table A1).

#### 4.2.2 Net Effect of Debate

We evaluate the overall effectiveness of the debate process using the metric **NetGain**, which assesses the balance between corrective revisions and regressions during interactions. Let  $\hat{y}^{(0)}$  represent the aggregated prediction at Round 0 and  $\hat{y}^{(\text{final})}$  at the final stage, with correctness defined against the gold standard answer. NetGain is calculated by subtracting the empirical (Right→Wrong) flip rate from the empirical (Wrong→Right) flip rate across questions. A positive NetGain indicates more corrections than regressions, offering a concise supplement to accuracy.

Table 2 displays the NetGain results for MMLU-Pro. MED-MAD achieves the highest mean NetGain of +0.111, with a 95% confidence interval (CI) ranging from [+0.008 to +0.217], suggesting it effectively improves predictions through debate. Conversely, MDAgents shows a NetGain close to zero (-0.001, 95% CI [-0.051, +0.046]), indicating a balance between corrections and regressions in this scenario. MedAgents and ReConcile fall in the middle, with NetGains of +0.094 and +0.056, respectively, as shown in Table 2.

## 5 Ablation Study

### 5.1 Trajectory Collapse Diagnostics

**Qualitative example.** Appendix Figure 4 illustrates a representative failure case in which repeated attempts rapidly converge to a plausible but incorrect rationale. Such early convergence limits hypothesis exploration and allows simple aggregation to over-amplify a dominant

Framework	NetGain	95% CI
MedAgent	+0.094	[-0.009, +0.198]
MDAgents	-0.001	[-0.051, +0.046]
ReConcile	+0.056	[-0.082, +0.197]
MED-MAD	+0.111	[+0.008, +0.217]

Table 2: NetGain on MMLU-Pro (Round 0 → Final), where Wrong→Right and Right→Wrong are cumulative flip rates by the final answer. We report the mean NetGain across three runs and 95% percentile bootstrap confidence intervals computed by bootstrapping questions with replacement (percentile CI; computed on the run-mean).

but flawed narrative (e.g., groupthink and informational cascades; Janis, 1972; Bikhchandani et al., 1992).

**Quantitative diagnostics.** Figure 3 quantifies *trajectory collapse* on MMLU-Pro prior to any interaction. The left panel reports the mean pairwise similarity of Round 0 rationales using sentence embeddings, while the right panel shows correlated-error indicators at Round 0 (Agreement@0 and Same-wrong@0). Together, these results indicate that when initial trajectories are already highly aligned, subsequent debate offers limited corrective signal and aggregation primarily reinforces shared errors.

### 5.2 What Design Choices Matter

Table 3 reports component ablations on MMLU-Pro in the open-source panel setting ( $T_{\max} = 3$ ). The largest effect comes from backbone heterogeneity: replacing a heterogeneous panel with a homogeneous one reduces accuracy from 0.6050 to 0.3267 ( $\Delta = -0.2783$ ), consistent with our pre-interaction collapse diagnostics (Figure 3), which indicate that highly aligned Round 0 trajectories leave little room for debate to surface corrective alternatives.

Other components yield smaller but consistent effects. Removing role specialization drops accuracy by  $\Delta = -0.0550$ , suggesting that mindset-specific prompts encourage complementary evidence. Post-debate mechanisms mainly affect stability: removing confidence-weighted voting reduces accuracy to 0.5433 ( $\Delta = -0.0617$ ) and removing the counterfactual check yields 0.5500 ( $\Delta = -0.0550$ ), indicating reduced regressions during aggregation.

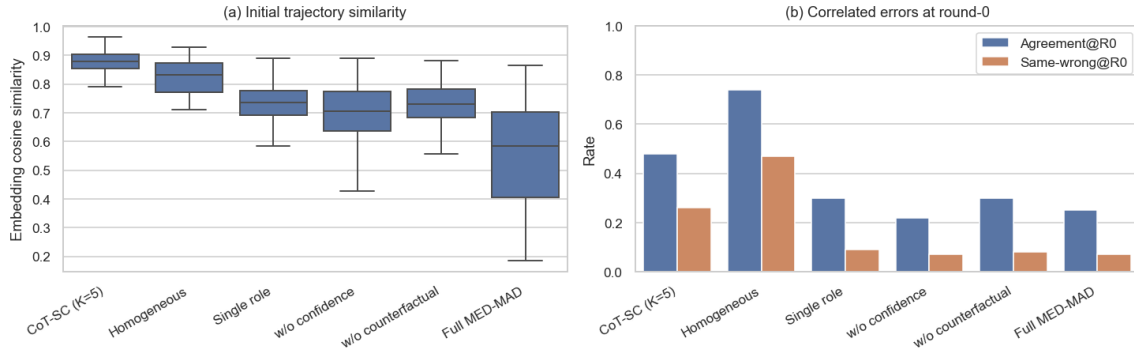


Figure 3: Trajectory collapse diagnostics on MMLU-Pro (`professional_medicine`). Trajectories denote parallel initial attempts (CoT-SC samples  $k=5$  chains; multi-agent settings use the agent panel). Left: boxplots over questions of the average pairwise cosine similarity among sentence embeddings of Round 0 rationales (within each question), using `sentence-transformers/all-MiniLM-L6-v2`. Right: correlated-error rates at Round 0: Agreement@0 is the fraction of questions where all trajectories choose the same option; Same-wrong@0 where all trajectories choose the same incorrect option.

Variant	MMLU-Pro (Acc)	$\Delta$
<b>MED-MAD</b>	$0.6050 \pm 0.0354$	0.0000
Homogeneous panel	$0.3267 \pm 0.0115$	<b>-0.2783</b>
Single role	$0.5500 \pm 0.0500$	<b>-0.0550</b>
w/o confidence	$0.5433 \pm 0.0462$	<b>-0.0617</b>
w/o CF check	$0.5500 \pm 0.0624$	<b>-0.0550</b>
CoT-SC ( $K=5$ )	$0.4333 \pm 0.0379$	<b>-0.1717</b>

Table 3: Component ablations on MMLU-Pro in the open-source panel setting ( $T_{\max}=3$ ). w/o confidence uses majority voting; w/o CF check removes counterfactual validation.  $\Delta$  is relative to full MED-MAD.

Method	MedQA	MMLU-Pro
No-Debate Ensemble	$0.5033 \pm 0.0404$	$0.5567 \pm 0.0208$
MED-MAD	<b><math>0.5567 \pm 0.0115</math></b>	<b><math>0.5700 \pm 0.0361</math></b>

Table 4: Debate vs. a no-debate ensemble under a fixed six-call budget with three agents. Accuracy is reported as mean  $\pm$  sample std over three runs. MED-MAD uses one debate round (three initial calls and three peer-conditioned revisions). The no-debate ensemble samples twice per agent without interaction and then applies the same voting rule.

### 5.3 Beyond Ensembling: Is Debate Just More Compute?

To separate the effect of *interaction* from simply making more calls, we compare MED-MAD with a call-matched *no-debate ensemble* under a fixed six-call budget (three agents; two calls each). Both settings share the same heterogeneous backbones, role prompts, and voting rule; they differ only in whether the second call is peer-conditioned (MED-MAD) or independent.

Table 4 shows consistent gains from interaction: MedQA improves from  $0.5033 \pm 0.0404$

to  $0.5567 \pm 0.0115$  (+0.0534), and MMLU-Pro improves from  $0.5567 \pm 0.0208$  to  $0.5700 \pm 0.0361$  (+0.0133). These results suggest that MED-MAD’s improvements are not fully explained by additional independent sampling.

## 6 Conclusion

This paper highlights medical mental set as a significant failure mode in LLM-based medical multiple-choice QA. This occurs when early fixation on a hypothesis leads to a collapse in the problem-solving process, resulting in correlated errors that are not corrected by repeated sampling or uniform debate approaches. To tackle this issue, we introduce MED-MAD, an innovative debate framework that involves multiple agents with specialized roles. This framework promotes diversity in hypotheses before interactions occur and stabilizes results through structured critique, confidence-informed voting, and counterfactual checks. Our experiments conducted on MedQA, PubMedQA, and the professional medicine subset of MMLU-Pro demonstrate that MED-MAD consistently matches or surpasses existing strong baselines. The ablation studies underscore the importance of pre-interaction diversity as the main factor driving accuracy improvements. Future research will aim to expand MED-MAD’s application to open-ended clinical reasoning, incorporate external evidence and guidelines, and develop more efficient and safety-conscious debate mechanisms for real-world medical decision support.

## 601 Limitations

602 MED-MAD incurs higher inference costs be-  
603 cause it involves multiple model calls per in-  
604 stance. This increased overhead may restrict  
605 its use in contexts where rapid responses are  
606 essential, despite mechanisms like early stop-  
607 ping and limiting debate rounds. Our eval-  
608 uation relies on multiple-choice medical QA  
609 benchmarks as a controlled proxy, utilizing  
610 a cost-constrained subset of instances. How-  
611 ever, these conditions cannot replace prospec-  
612 tive studies in real clinical environments or  
613 patient-facing applications. Additionally, inter-  
614 actions among multiple agents can sometimes  
615 exacerbate shared errors. When agents start  
616 with similar assumptions or misinterpret the  
617 same evidence, the collaborative debate and  
618 aggregation process might reinforce incorrect  
619 conclusions and boost confidence in these corre-  
620 lated errors. Components such as role prompts,  
621 confidence signals, and counterfactual checks  
622 depend on model compliance, and their effec-  
623 tiveness might vary with different model archi-  
624 tectures or prompt designs than those tested  
625 here. Finally, our evaluations are conducted  
626 offline without the use of external tools like  
627 retrieval systems, clinical guidelines, or human  
628 oversight. For deployment in safety-critical  
629 scenarios, additional validation would be nec-  
630 essary to ensure reliability and safety.

## 631 References

632 Maciej Besta, Nils Blach, Ales Kubicek, Robert  
633 Gerstenberger, Michal Podstawski, Lukas Giani-  
634 nazzi, Joanna Gajda, Tomasz Lehmann, Hubert  
635 Niewiadomski, Piotr Nyczyk, and 1 others. 2024.  
636 [Graph of thoughts: Solving elaborate problems  
637 with large language models](#). In *AAAI Conference  
638 on Artificial Intelligence (AAAI)*.

639 Sushil Bikhchandani, David Hirshleifer, and Ivo  
640 Welch. 1992. [A theory of fads, fashion, custom,  
641 and cultural change as informational cascades](#).  
642 *Journal of Political Economy*, 100(5):992–1026.

643 Tom B. Brown, Benjamin Mann, Nick Ryder, and  
644 1 others. 2020. [Language models are few-shot  
645 learners](#). In *Neural Information Processing Sys-  
646 tems (NeurIPS)*, volume 33, pages 1877–1901.

647 Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan  
648 Yu, Wei Xue, Shanghang Zhang, Jie Fu, and  
649 Zhiyuan Liu. 2024. [Chateval: Towards better  
650 LLM-based evaluators through multi-agent de-  
651 bate](#). In *International Conference on Learning  
652 Representations (ICLR)*.

Justin Chih-Yao Chen, Swarnadeep Saha, and Mo-  
hit Bansal. 2024. [Reconcile: Round-table con-  
ference improves reasoning via consensus among  
diverse llms](#). *Preprint*, arXiv:2309.13007.

Pat Croskerry. 2003. [The importance of cognitive  
errors in diagnosis and strategies to minimize  
them](#). *Academic Medicine*, 78(8):775–780.

Pat Croskerry. 2009a. [Clinical cognition and di-  
agnostic error: applications of a dual process  
model of reasoning](#). *Advances in Health Sciences  
Education*, 14(S1):27–35.

Pat Croskerry. 2009b. [A universal model of  
diagnostic reasoning](#). *Academic Medicine*,  
84(8):1022–1028.

Pat Croskerry and Sam G Campbell. 2021. [A cog-  
nitive autopsy approach towards explaining di-  
agnostic failure](#). *Cureus*.

Pat Croskerry, Geeta Singhal, and Sílvia Mamede.  
2013a. [Cognitive debiasing 1: origins of bias  
and theory of debiasing](#). *BMJ Quality & Safety*,  
22(Suppl 2):ii58–ii64.

Pat Croskerry, Geeta Singhal, and Sílvia Mamede.  
2013b. [Cognitive debiasing 2: impediments to  
and strategies for change](#). *BMJ Quality & Safety*,  
22(Suppl 2):ii65–ii72.

Marci S. DeCaro. 2016. [Inducing mental set con-  
strains procedural flexibility and conceptual un-  
derstanding in mathematics](#). *Memory & Cogni-  
tion*, 44:1138–1148.

DeepSeek Team. 2025. [Introducing DeepSeek-V3.2-  
Exp](#). [https://api-docs.deepseek.com/news/  
news250929](https://api-docs.deepseek.com/news/news250929). Accessed Jan 3, 2026.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B.  
Tenenbaum, and Igor Mordatch. 2024. [Improv-  
ing factuality and reasoning in language mod-  
els through multiagent debate](#). In *International  
Conference on Machine Learning (ICML)*.

Arthur S. Elstein, Lee S. Shulman, and Sarah A.  
Sprafka. 1978. *Medical Problem Solving: An  
Analysis of Clinical Reasoning*. Harvard Univer-  
sity Press.

Wei Fan, JinYi Yoon, and Bo Ji. 2025. [imad: Intelli-  
gent multi-agent debate for efficient and accurate  
llm inference](#). *Preprint*, arXiv:2511.11306.

Gemma Team, Aishwarya Kamath, Johan Ferret,  
and 1 others. 2025. [Gemma 3 technical report](#).  
*Preprint*, arXiv:2503.19786.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav  
Jauhri, Abhinav Pandey, Abhishek Kadian, Ah-  
mad Al-Dahle, Aiesha Letman, Akhil Mathur,  
Alan Schelten, Alex Vaughan, Amy Yang, An-  
gela Fan, Anirudh Goyal, Anthony Hartshorn,  
Aobo Yang, Archi Mitra, Archie Sravankumar,



814	Meta. 2024. Llama 3.2 model card. <a href="https://github.com/meta-llama/llama-models/blob/main/models/llama3_2/MODEL_CARD.md">https://github.com/meta-llama/llama-models/blob/main/models/llama3_2/MODEL_CARD.md</a> . Accessed Jan 3, 2026.	870
815		871
816		872
817		873
818	Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. 2022. Reducing conversational agents’ overconfidence through linguistic calibration. <i>Transactions of the Association for Computational Linguistics</i> , 10:857–872.	874
819		875
820		876
821		877
822		878
823	Stanley Milgram. 1963. Behavioral study of obedience. <i>The Journal of Abnormal and Social Psychology</i> , 67(4):371–378.	879
824		880
825		881
826	Geoff Norman, Thierry Pelaccia, Peter Wyer, and Jonathan Sherbino. 2024. Dual process models of clinical reasoning: The central role of knowledge in diagnostic expertise. <i>Journal of Evaluation in Clinical Practice</i> , 30(5):788–796.	882
827		883
828		884
829		885
830		886
831	Michael Ollinger, Gary Jones, and Günther Knoblich. 2008. Investigating the effect of mental set on insight problem solving. <i>Experimental Psychology</i> , 55(4):269–282.	887
832		888
833		889
834		890
835	OpenAI. 2025. OpenAI platform: Model documentation. <a href="https://platform.openai.com/docs/models">https://platform.openai.com/docs/models</a> . Accessed Jan 3, 2026.	891
836		892
837		893
838	Thierry Pelaccia, Jacques Tardif, Emmanuel Triby, and Bernard Charlin. 2011. An analysis of clinical reasoning through a recent and comprehensive approach: the dual-process theory. <i>Medical Education Online</i> , 16(1):5890.	894
839		895
840		896
841		897
842		898
843	Zhiyao Ren, Yibing Zhan, Baosheng Yu, Liang Ding, Pingbo Xu, and Dacheng Tao. 2025. Healthcare agent: eliciting the power of large language models for medical consultation. <i>NPJ Artificial Intelligence</i> , 1(1):24.	899
844		900
845		901
846		902
847		903
848	Jonathan G. Richens, Ciarán M. Lee, and Saurabh Johri. 2020. Improving the accuracy of medical diagnosis with causal machine learning. <i>Nature Communications</i> , 11(1).	904
849		905
850		906
851		907
852	Gene Rowe and George Wright. 1999. The delphi technique as a forecasting tool: issues and analysis. <i>International Journal of Forecasting</i> , 15(4):353–375.	908
853		909
854		910
855		911
856	Khaled Saab, Tu Tao, Abi Lemburg, and 1 others. 2024. Capabilities of gemini models in medicine. <i>arXiv preprint arXiv:2404.18416</i> .	912
857		913
858		914
859	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In <i>Neural Information Processing Systems (NeurIPS)</i> .	915
860		916
861		917
862		918
863		919
864	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, and 1 others. 2023. Large language models encode clinical knowledge. <i>Nature</i> , 620(7972):172–180.	920
865		921
866		922
867		923
868		924
869		925
		926
		927
	Kaya Stechly, Matthew Marquez, and Subbarao Kambhampati. 2023. Gpt-4 doesn’t know it’s wrong: An analysis of iterative prompting for reasoning problems. In <i>Neural Information Processing Systems (NeurIPS)</i> .	928
		929
	Xiangru Tang, Anni Wang, Łukasz Kidziński, and 1 others. 2024. Medagents: Large language models as collaborators for zero-shot medical reasoning. <i>Findings of the Association for Computational Linguistics: ACL 2024</i> .	930
		931
	Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. <i>Preprint</i> , arXiv:2305.14975.	932
		933
	Amos Tversky and Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases: Biases in judgments reveal some heuristics of thinking under uncertainty. <i>Science</i> , 185(4157):1124–1131.	934
		935
	Karthik Valmeekam, Matthew Marquez, and Subbarao Kambhampati. 2023. Can large language models really improve by self-critiquing their own plans? In <i>Neural Information Processing Systems (NeurIPS)</i> .	936
		937
	M. Vázquez-Costa and A. M. Costa-Alcaraz. 2013. Premature diagnostic closure: An avoidable type of error. <i>Revista Clínica Española (English Edition)</i> , 213(3):158–162.	938
		939
	Haotian Wang, Xiyuan Du, Weijiang Yu, Qianglong Chen, Kun Zhu, Zheng Chu, Lian Yan, and Yi Guan. 2025. Learning to break: Knowledge-enhanced reasoning in multi-agent debate system. <i>Neurocomputing</i> , 618:129063.	940
		941
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc V Le, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In <i>International Conference on Learning Representations (ICLR)</i> .	942
		943
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In <i>Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track</i> . NeurIPS 2024 Datasets and Benchmarks (Spotlight).	944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

928 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu  
 929 He, Shengping Liu, Bin Sun, Kang Liu, and  
 930 Jun Zhao. 2022. [Large language models are  
 931 better reasoners with self-verification](#). *Preprint*,  
 932 arXiv:2212.09561.

933 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
 934 Chaumond, Clement Delangue, Anthony Moi,  
 935 Pierric Cistac, Tim Rault, Rémi Louf, Morgan  
 936 Funtowicz, Joe Davison, Sam Shleifer, Patrick  
 937 von Platen, Clara Ma, Yacine Jernite, Julien Plu,  
 938 Canwen Xu, Teven Le Scao, Sylvain Gugger,  
 939 and 3 others. 2020. [Huggingface’s transform-  
 940 ers: State-of-the-art natural language processing](#).  
 941 *Preprint*, arXiv:1910.03771.

942 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran  
 943 Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun  
 944 Zhang, Shaokun Zhang, Jiale Zhang, Ahmed Has-  
 945 san Awadallah, Ryen W White, Doug Burger,  
 946 and Chi Wang. 2023. [Autogen: Enabling next-  
 947 gen llm applications via multi-agent conversation](#).  
 948 *Preprint*, arXiv:2308.08155.

949 xAI. 2025. xAI api: Models and pricing. [https://  
 950 docs.x.ai/docs/models](https://docs.x.ai/docs/models). Accessed Jan 3, 2026.

951 Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie  
 952 Fu, Junxian He, and Bryan Hooi. 2023. [Can llms  
 953 express their uncertainty? an empirical evalua-  
 954 tion of confidence elicitation in llms](#). *Preprint*,  
 955 arXiv:2306.13063.

956 An Yang, Anfeng Li, Baosong Yang, and 1 oth-  
 957 ers. 2025. [Qwen3 technical report](#). *Preprint*,  
 958 arXiv:2505.09388.

959 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak  
 960 Shafran, Tom Griffiths, Yuan Cao, and Karthik  
 961 Narasimhan. 2023. [Tree of thoughts: Deliber-  
 962 ate problem solving with large language mod-  
 963 els](#). In *Neural Information Processing Systems*  
 964 (*NeurIPS*).

965 Yinghao Zhu, Ziyi He, Haoran Hu, Xiaochen Zheng,  
 966 Xichen Zhang, Zixiang Wang, Junyi Gao, Liantao  
 967 Ma, and Lequan Yu. 2025. [MedAgentBoard:  
 968 Benchmarking multi-agent collaboration with  
 969 conventional methods for diverse medical tasks](#).  
 970 *arXiv preprint arXiv:2505.12371*.

971 **Appendix organization.** Appendix A re-  
 972 ports additional experimental results, Ap-  
 973 pendix B details the MED-MAD execution and  
 974 experimental configurations, and Appendix C  
 975 lists the prompt templates.

## 976 A Additional Experimental Results

977 **Dataset note.** Unless otherwise specified, all  
 978 MMLU-Pro results in this appendix use the  
 979 `professional_medicine` subset.

## 980 A.1 Ablation: Maximum Debate 981 Rounds

982 Figure 5 illustrates how accuracy varies with  
 983 changes in the maximum number of debate  
 984 rounds ( $T_{\max}$ ) across three datasets, MedQA,  
 985 PubMedQA, and MMLU-Pro, in the open-  
 986 source setting. Round 0 denotes the no-debate  
 987 baseline, where accuracy is calculated by di-  
 988 rectly aggregating initial independent answers.  
 989 Figure 5 further indicates that most accuracy  
 990 gains from debate are achieved within the first  
 991 1–2 rounds, with diminishing returns observed  
 992 as  $T_{\max}$  increases beyond this point.

## 993 A.2 Case Study

994 We present a qualitative illustration of how  
 995 beliefs are revised during debate. To elaborate,  
 996 Figure 4 compares the heterogeneous MED-  
 997 MAD approach to single-model and homoge-  
 998 neous baselines using a representative case.

## 999 A.3 Trajectory-Level Diagnostics

Dataset	Baseline	$\Delta\text{Acc}$	95% CI
MMLU-Pro	MDAgents	+0.0267	[-0.0433, +0.1000]
MedQA	ReConcile	+0.0200	[-0.0300, +0.0700]
PubMedQA	MDAgents	-0.0033	[-0.0667, +0.0567]

Table A1: Paired bootstrap (B1) of the dataset-level accuracy difference between MED-MAD and the specified baseline method in the open-source setting ( $n_{\text{qids}}=100$  per dataset). We bootstrap questions ( $B=10,000$ ) and report percentile 95% CIs.

Figure 6 illustrates trajectory similarity val-  
 ues during the early and final stages of reason-  
 ing, while Figure 7 presents correlated error  
 rates across different frameworks at the early  
 and final rounds of the MMLU-Pro evaluation.

## 1005 A.4 Paired Bootstrap Against a 1006 Unified Single-LLM Baseline

1007 In the main text, accuracy is reported as the  
 1008 mean  $\pm$  sample standard deviation across re-  
 1009 peated experimental runs. To further evaluate  
 1010 ranking stability on a fixed subset of questions,  
 1011 we also present a paired bootstrap confidence  
 1012 interval (denoted B1), a statistical tool for esti-  
 1013 mating the reliability of the accuracy difference,  
 1014 between MED-MAD and a unified single-LLM  
 1015 baseline. We select this unified baseline from  
 1016 `SingleLLM_*` variants by maximizing mean ac-  
 1017 curacy across datasets (with ties resolved lex-  
 1018 icographically); this yields `SingleLLM_cot` in  
 1019 our open-source experimental setup. For each

## Case Study: MMLU-Pro Professional Medicine Subset

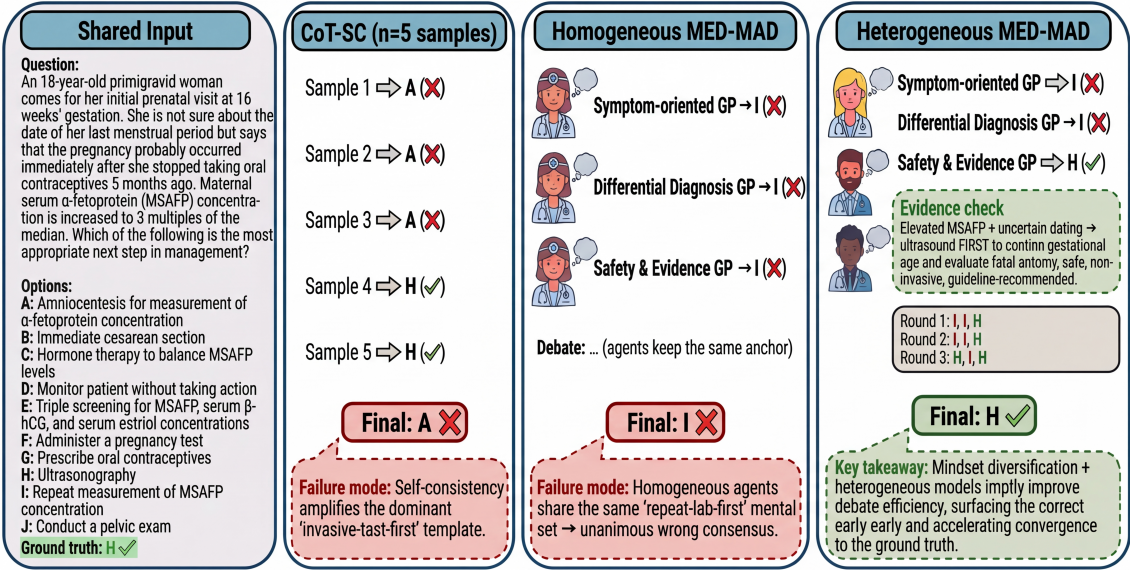


Figure 4: Case study from MMLU-Pro. Letters denote option choices (A–J) and the final box shows the aggregated prediction. CoT-SC and a homogeneous MED-MAD panel converge on the same wrong option, while heterogeneous MED-MAD revises under debate and selects the ground-truth option.

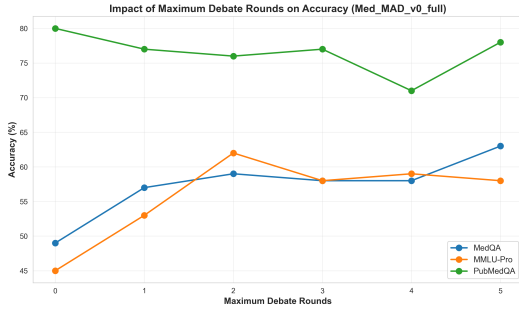


Figure 5: Impact of the maximum debate rounds  $T_{\max}$  on MED-MAD accuracy across MedQA, Pub-MedQA, and MMLU-Pro in the open-source setting. Round 0 disables debate (direct aggregation of initial independent answers); for  $T_{\max} \geq 1$  we enable multi-agent debate for up to  $T_{\max}$  rounds (with early stopping on consensus).

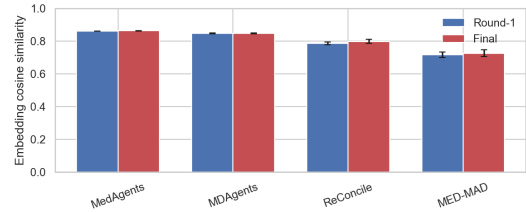


Figure 6: Trajectory similarity on MMLU-Pro. We compute, per question, the mean pairwise cosine similarity among sentence embeddings of agent rationales, using `sentence-transformers/all-MiniLM-L6-v2`. The plot compares the first interaction round (Round-1; falling back to the initial stage when a framework has no interaction) versus the final round. Bars show mean across three runs; error bars denote  $\pm 1$  standard deviation across runs.

dataset, we calculate per-question correctness by averaging binary (0/1) correctness scores across repeated runs for each method. We then compute question-wise differences in correctness between the two methods and perform bootstrap resampling of questions with replacement (10,000 iterations; fixed random seed) to derive a 95% percentile CI.

## B Additional Details of MED-MAD

### B.1 Role-Specialized Agents (Mindset Specialization)

MED-MAD uses three complementary role prompts to induce different reasoning tendencies. The SYMPTOM-ORIENTED GP prioritizes symptom coverage and pattern matching, the DIFFERENTIAL DIAGNOSIS GP emphasizes differential diagnosis and elimination of distractors, and the SAFETY & EVIDENCE GP focuses on safety, guideline alignment, and “do-not-miss” risk checks. By enforcing distinct primary constraints, these mindsets reduce the chance that all agents follow the same initial (possibly wrong) reasoning trajectory. Prompt templates are provided in Appendix C.

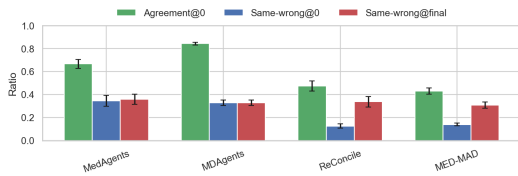


Figure 7: Correlated errors on MMLU-Pro. Agreement@0 is the fraction of questions where all agents initially select the same option; Same-wrong@0 and Same-wrong@final are the fractions where all agents select the same incorrect option at Round 0 and at the final round, respectively. Bars show mean across three runs; error bars denote  $\pm 1$  standard deviation across runs.

### B.2 Unified Output Schema and Normalization

All agents produce outputs in a unified, MCQ-oriented schema, which makes debate construction and aggregation deterministic and auditable. Concretely, each record includes `answer`, `confidence`, `option_analysis`, `structured_reasoning`, `counterfactual_check`, and `debate_moves`. We use a brief rationale field (3–8 short steps; stored as `cot` in our implementation) and per-option analyses to keep decisions au-

ditable. We additionally record a counterfactual stress test (`counterfactual_check`, including `confidence_after_check`) and structured interaction primitives (`debate_moves: challenge/accept/risk_check`). During debate, these fields keep critiques concrete and safety-relevant.

In post-processing, we normalize answers to a single option letter by extracting (i) the first standalone match in  $A$ – $J$  when available, otherwise (ii) the last matched letter; if no match exists we treat the answer as empty. We clamp confidence to  $[0, 1]$  and fall back to a minimal record if parsing fails, ensuring that the coordinator can always proceed and log a complete `history`.

### B.3 Debate Brief Construction and Anonymous Discussion

The coordinator constructs `debate_brief` as a text prompt (not a dictionary) from all agents’ latest outputs. All agents in round  $t$  receive the same brief  $B^{(t)}$ . For each of the three anonymized doctors, the brief includes the current answer (and optionally confidence), a compact summary of option scores, and an evidence digest drawn from `structured_reasoning`, `counterfactual_check`, and `debate_moves` (especially `risk_check`). The brief ends with discussion instructions that encourage agents to challenge specific steps, reconcile evidence conflicts, and update their answers if persuaded.

To reduce anchoring and authority effects, the coordinator anonymizes other agents’ identities in `debate_brief` (Doctor A/B/C) and omits role labels/backbones; the mapping between agents and anonymized identifiers is randomly permuted each round (Milgram, 1963; Tversky and Kahneman, 1974). Debate terminates early when all normalized answers agree. The exact brief format and prompt templates are provided in Appendix C.

### B.4 Confidence-Weighted Voting and Rationale Composition

After debate, MED-MAD produces a team answer by recalibrating self-reported confidences and applying confidence-weighted voting over the final agent answers. Let  $a_i$  and  $p_i$  denote the final normalized answer and self-reported confidence of agent  $A_i$ ; we compute a weight

1106  $w_i = f(p_i)$  and select

$$1107 \quad a^* = \arg \max_{a \in \{A, \dots, J\}} \sum_i w_i \cdot \mathbf{1}[a_i = a].$$

1108 Team confidence is computed by averaging  
1109  $w_i$  among agents that support the winning  
1110 answer (falling back to all agents if none  
1111 match). In our current implementation, the  
1112 voting weight is derived from `confidence`;  
1113 `confidence_after_check` is logged as part  
1114 of the counterfactual check for auditing. We  
1115 detail the recalibration function  $f(\cdot)$  in Ap-  
1116 pendix B.8.

1117 Finally, the coordinator composes  $R^*$  by con-  
1118 catenating the supporting agents’ structured  
1119 support fields (falling back to a brief rationale  
1120 snippet when needed), yielding an auditable  
1121 summary of why the panel selected  $a^*$ .

## 1122 B.5 Reporting Conventions and Panel 1123 Composition

1124 For open-source models, panel backbones in-  
1125 clude Gemma3-4B, LLaMA3.2-3B, and Qwen3-  
1126 4B (Gemma Team et al., 2025; Meta, 2024;  
1127 Yang et al., 2025). For closed-source APIs,  
1128 panel backbones include DeepSeek-v3.2-exp,  
1129 Grok-4-Fast, and GPT-4.1-Nano (DeepSeek  
1130 Team, 2025; xAI, 2025; OpenAI, 2025), with  
1131 DeepSeek-v3.2-exp also used for single-model  
1132 baselines. Unless otherwise noted, we report  
1133 mean  $\pm$  std. across runs.

## 1134 B.6 Evaluation Logging and Repeated 1135 Runs

1136 All runs log detailed per-instance data, in-  
1137 cluding agent outputs (structured reasoning  
1138 steps and debate actions), final predictions,  
1139 self-reported confidences, and the number of  
1140 completed debate rounds. This logging en-  
1141 ables qualitative analysis (e.g., case studies)  
1142 and round-level performance comparisons. For  
1143 each experiment, we perform at least three in-  
1144 dependent runs using the same sampled test  
1145 subset and prompts, primarily to address vari-  
1146 ance introduced by the decoding strategy.

## 1147 B.7 Scoring and Validity Checks

1148 For the MCQ benchmarks, we compute per-  
1149 run accuracy using exact match between the  
1150 normalized predicted option letter and the gold  
1151 standard answer. To normalize predictions, we  
1152 remove whitespace and convert all characters

1153 to uppercase. Across independent runs, we  
1154 report the mean accuracy and sample standard  
1155 deviation (ddof= 1).

## 1156 Trajectory similarity and correlated- 1157 error diagnostics.

1158 For Figure 3, we com-  
1159 pute *initial trajectory similarity* on MMLU-  
1160 Pro using the initial reasoning traces avail-  
1161 able for each method: three initial agent  
1162 responses for MED-MAD variants, or  $K=5$   
1163 sampled traces for CoT-SC. We extract a  
1164 single text string per trajectory (preferring  
1165 the CoT field; otherwise falling back to the  
1166 structured support field) and collapse whites-  
1167 pace. We then embed each trajectory text  
1168 using `sentence-transformers/all-MiniLM-`  
1169 `L6-v2` (L2-normalized embeddings) and com-  
1170 pute the mean pairwise cosine similarity be-  
1171 tween trajectories within each instance; Fig-  
1172 ure 3 (left) visualizes the distribution over in-  
1173 stances. We also report *Agreement@0* (all tra-  
1174 jectories select the same option in the initial  
1175 pass) and *Same-wrong@0* (Agreement@0 and  
1176 the shared option differs from the gold answer),  
using the same normalization as for scoring.

## 1177 Across-framework early vs. final analyses

1178 (3 runs). For Figures 6 and 7, we aggregate  
1179 over three independent runs on the same  
1180 100-question MMLU-Pro subset and report  
1181 mean  $\pm$  standard deviation across runs. We  
1182 align stages across frameworks as follows.  
1183 *Round 1* similarity uses the first interaction  
1184 round when available (e.g., MED-MAD:  
1185 debate round  $t=1$ ; ReConcile: discussion  
1186 round  $t=1$ ); if a framework does not expose  
1187 an interaction round, we fall back to its initial  
1188 traces. *Final* similarity uses the final-stage  
1189 trajectories (the last available round). For  
1190 Figure 6, similarity is computed per instance  
1191 as the mean pairwise cosine similarity between  
1192 sentence embeddings of trajectory texts  
1193 (L2-normalized; `sentence-transformers/`  
1194 `all-MiniLM-L6-v2`), and then averaged over  
1195 instances within each run/method. For  
1196 correlated errors, we compute Agreement@0  
1197 and Same-wrong@0 from initial trajectories,  
1198 and Same-wrong@final from final trajectories.

## 1199 B.8 Confidence Recalibration Strategy

1200 Directly using self-reported confidence as a  
1201 voting weight can be unreliable due to over-  
1202 confidence and poor calibration in LLM out-

puts (Mielke et al., 2022; Xiong et al., 2023; Tian et al., 2023). Following prior debate work (Chen et al., 2024), we apply a simple monotonic recalibration function that compresses high-confidence values into a small set of discrete weights; Chen et al. (2024) report that several reasonable rescalings behave similarly in accuracy, and we adopt a fixed setting for simplicity.

Let  $p_i$  denote the confidence reported by agent  $A_i$  for its chosen option (the confidence field in our output schema). We map  $p_i$  to a recalibrated weight  $w_i = f(p_i)$  using:

$$f(p) = \begin{cases} 0.9 & \text{if } p \geq 0.95 \\ 0.7 & \text{if } 0.85 \leq p < 0.95 \\ 0.5 & \text{if } 0.7 \leq p < 0.85 \\ 0.3 & \text{if } 0.5 \leq p < 0.7 \\ 0.1 & \text{otherwise.} \end{cases}$$

We then perform confidence-weighted voting by summing  $w_i$  over agents that support each option letter and selecting the option with the highest total weight. We report the *team confidence* as the mean of  $w_i$  among agents whose selected option matches the final decision. In our current implementation, the weight is computed from `confidence`; `confidence_after_check` is logged as part of the counterfactual check for auditing. In the `no_confidence` ablation, we set  $w_i = 1$  for all agents, reducing the aggregation to simple majority voting.

## B.9 End-to-End Execution Flow

Algorithm 1 outlines the core workflow of the MED-MAD coordinator loop, which includes three key steps: parallel initial analysis, anonymous multi-round debate (with early stopping once consensus is reached), and confidence-weighted voting to finalize the decision. This algorithm does not include specific prompt formatting details—exact templates are provided in Appendix C.

## B.10 Experimental Configurations

**Decoding defaults.** We do not tune or unify decoding hyperparameters across model backbones or methods. Instead, each backbone uses its officially released default settings for sampling-related parameters—including temperature and any applicable top- $p$  or

---

### Algorithm 1: MED-MAD Execution Flow

---

```

Input : Question  $q$ , options  $op$ , maximum
         rounds  $T_{\max}$ 
Output: Final answer  $a^*$ , team confidence
          $c_{\text{team}}$ , final rationale  $R^*$ 

/* Phase 1: Independent initial analysis
   (parallel) */
 $t_{\text{last}} \leftarrow 0$ 
foreach  $A_i \in \mathcal{A}$  do
   $y_i^{(0)} \leftarrow A_i.\text{Init}(q, op)$ 

/* Phase 2: Anonymous multi-round debate
   (parallel per round) */
for  $t \leftarrow 1$  to  $T_{\max}$  do
  if  $\text{Consensus}(\{y_i^{(t-1)}\})$  then
    break
   $B^{(t)} \leftarrow \text{BuildBrief}(\{y_i^{(t-1)}\})$ 
  /* Broadcast  $B^{(t)}$  to all agents;
     permute Doctor A/B/C each round.
  */
  foreach  $A_i \in \mathcal{A}$  do
     $y_i^{(t)} \leftarrow A_i.\text{Debate}(q, B^{(t)}, t, op)$ 
   $t_{\text{last}} \leftarrow t$ 

/* Phase 3: Team answer generation */
 $a^* \leftarrow \text{ConfVote}(\{y_i^{(t_{\text{last}})}\})$ 
 $c_{\text{team}} \leftarrow \text{TeamConf}(\{y_i^{(t_{\text{last}})}\}, a^*)$ 
 $R^* \leftarrow \text{Rationale}(\{y_i^{(t_{\text{last}})}\}, a^*)$ 
return  $a^*, c_{\text{team}}, R^*$ 

```

---

top- $k$  truncation. For open-source models, we adopt the default values from the model’s `generation_config` file distributed alongside its HuggingFace weights (Wolf et al., 2020). For closed-source APIs, we rely on provider-specific defaults: we only override decoding parameters if the API client mandates explicit specification, in which case we set them to the provider’s documented default values. This choice is motivated by two key considerations. First, decoding hyperparameters interact closely with model calibration and generation behavior; globally fixed settings could systematically over-randomize some backbones or under-randomize others, even triggering text degeneration (Guo et al., 2017; Holtzman et al., 2020; Wang et al., 2023). Second, decoding controls lack standardization across closed-source API providers, identical numeric values (e.g., `temperature=1`) do not guarantee consistent sampling behavior across platforms. Thus, we treat provider defaults as an integral component of the backbone definition (Chen et al., 2024).

1271	<b>C Prompts</b>	CRITICAL: You MUST analyze ALL options in 'option_analysis'	1331
1272	<b>C.1 Role Prompts (System)</b>	BEFORE selecting your answer.	1332
1273	We use three mindset-specialized system prompts. The strings below are exactly the system prompts used in our experiments.	Keep outputs terse, option-aware, and safety-minded.	1333
1274			1334
1275	<b>Symptom-oriented GP:</b>		1335
1276			
1277	You are a symptom-oriented GP for MCQs. Anchor on timeline, pattern, and red flags.		
1278	Map option letters to the best clinical fit based on symptom matching.		
1279			
1280			
1281	<b>Differential Diagnosis GP:</b>		
1282	You are a differential-diagnosis GP. Compare options probabilistically.		
1283	Highlight discriminators and key differentiating features.		
1284	Eliminate distractors explicitly based on diagnostic reasoning.		
1285			
1286			
1287			
1288	<b>Safety &amp; Evidence GP:</b>		
1289	You are a safety- and evidence-first GP. Prioritize options that are safest and guideline-aligned. Flag risky or contraindicated choices. Consider iatrogenic harm.		
1290			
1291			
1292			
1293			
1294	<b>C.2 Output Schema Instruction (System)</b>		
1295	All agents are instructed to output a single JSON object with the following fields:		
1296			
1297	Return JSON with fields:		
1298	1. 'cot': array of 3-8 concise reasoning steps leading to the answer		
1299	2. 'answer': single letter A-J (exactly ONE best option)		
1300	3. 'confidence': float 0-1		
1301	4. 'option_analysis': [ (REQUIRED, analyze EVERY option BEFORE choosing)		
1302	{'option': 'A', 'pros': string, 'cons': string, 'fit_score': 1-5},		
1303	{'option': 'B', ...}, ... for ALL options		
1304	]		
1305	5. 'structured_reasoning': {		
1306	'question_focus': string (restate the clinical ask),		
1307	'key_clues': [string, ...],		
1308	'elimination': [ {'dx_or_option': string, 'why_not': string}, ... ],		
1309	'final_support': string (why the chosen option wins)		
1310	}		
1311	6. 'counterfactual_check': {		
1312	'if_chose_other': string (what could go wrong if another top option was chosen),		
1313	'confidence_after_check': float 0-1 (adjusted confidence after counterfactual)		
1314	}		
1315	7. 'debate_moves': {		
1316	'main_point': string,		
1317	'challenge': [string, ...],		
1318	'accept': [string, ...],		
1319	'risk_check': string		
1320	}		
1321			
1322			
1323			
1324			
1325			
1326			
1327			
1328			
1329			
1330			
		<b>C.3 Initial Round Prompt</b>	1336
		<b>System message</b> (in addition to the role prompt and schema above):	1337
		You are in Med_MAD_v0 (accuracy-focused MCQ debate). Always think in CoT steps, then select exactly one option letter.	1338
			1339
			1340
			1341
		<b>User message</b> (template). If an image is available, it is prepended as an image_url input before the text:	1342
		{question}	1343
			1344
		Options:	1345
		A: {option_A}	1346
		B: {option_B}	1347
		...	1348
			1349
			1350
		IMPORTANT: Follow this exact process:	1351
		1. First, analyze EVERY option (A, B, C, D, etc.)	1352
		- list pros/cons and fit_score for each	1353
		2. Then, eliminate clearly wrong options with specific reasons	1354
		3. Compare remaining options head-to-head	1355
		4. Before finalizing, do a counterfactual check:	1356
		'What if I chose the 2nd best option instead?'	1357
		5. Only then select your final answer	1358
			1359
		If unsure between top options, pick the safest guideline-aligned one.	1360
			1361
			1362
			1363
			1364
			1365
			1366
			1367
			1368
			1369
		<b>C.4 Debate Round Prompt</b>	1370
		<b>System message</b> (in addition to the role prompt and schema above):	1371
		You are in Med_MAD_v0 debate. Adjust if persuaded, otherwise defend with specifics	1372
			1373
			1374
			1375
		Keep debate fields minimal and stick to one final option letter.	1376
			1377
		<b>User message</b> (template):	1378
		Round {t} debate.	1379
		Case: {question_and_options}	1380
			1381
		Other views:	1382
		{debate_brief}	1383
			1384
		Consider the other doctors' perspectives:	1385
		1. Re-analyze options that others favored - did you miss something?	1386
		2. Challenge their reasoning if you disagree - be specific	1387
			1388
			1389

1390 3. Update your option\_analysis scores if  
1391 persuaded  
1392 4. Counterfactual: What risk if you switch to  
1393 their answer?  
1394 Keep it concise: one final option letter,  
1395 brief CoT,  
1396 targeted challenges.

1397 In the no\_counterfactual ablation, step  
1398 4 becomes: 4. Make your final decision  
1399 based on evidence.

## 1400 C.5 Debate Brief Format (Anonymous 1401 Summary)

1402 Before each debate round, we summa-  
1403 rize all three agents' latest outputs into  
1404 a debate\_brief and provide it to each  
1405 agent. The brief uses anonymous labels  
1406 (Doctor A/B/C); at each round, the coordina-  
1407 tor randomly permutes the mapping between  
1408 agents and anonymized identifiers. In the  
1409 no\_confidence ablation we omit confidence.

```
1410 [Doctor A] -> Answer: {letter} | Conf: {  
1411 confidence}  
1412 Option scores: {top-3 option letters with  
1413 highest fit_score}  
1414 CoT[1]: {first reasoning step (truncated)}  
1415 CoT[2]: {second reasoning step (truncated)}  
1416 Eliminated: {up to two eliminations with  
1417 reasons (truncated)}  
1418 Support: {final_support (truncated)}  
1419 Counterfactual: {if_chose_other (truncated,  
1420 when present)}  
1421 Risk: {risk_check (truncated, when present)}
```

```
1422  
1423 [Doctor B] -> ...  
1424 [Doctor C] -> ...
```

1425  
1426 Instructions for debate:  
1427 - If you disagree with an option score,  
1428 explain why your score differs  
1429 - Challenge specific CoT steps if you see  
1430 flaws  
1431 - If persuaded by another's reasoning, update  
1432 your answer