A Knowledge-driven Adaptive Collaboration of LLMs for **Enhancing Medical Decision-making**

Anonymous ACL submission

Abstract

Medical decision-making often involves integrating knowledge from multiple clinical specialties, typically achieved through multidisci-005 plinary teams. Inspired by this collaborative process, recent work has leveraged large language models (LLMs) in multi-agent collaboration frameworks to emulate expert teamwork. While these approaches improve reasoning through agent interaction, they are limited by static, pre-assigned roles, which hinder adaptability and dynamic knowledge integration. To address these limitations, we propose KAMAC, a Knowledge-driven Adaptive Multi-Agent Collaboration framework that enables LLM agents to dynamically form and expand expert teams based on the evolving diagnostic context. KAMAC begins with 018 one or more expert agents and then conducts a knowledge-driven discussion to identify and fill knowledge gaps by recruiting additional specialists as needed. This supports flexible, scalable collaboration in complex clinical scenarios, with decisions finalized through reviewing updated agent comments. Experiments on two real-world medical benchmarks demonstrate that KAMAC significantly outperforms both single-agent and advanced multi-agent meth-028 ods, particularly in complex clinical scenarios (i.e., cancer prognosis) requiring dynamic, cross-specialty expertise.

Introduction 1

007

011

017

019

034

042

In healthcare, diagnosis, prognosis, and a variety of clinical treatments are guided by medical decision-making processes that require the application of complex medical knowledge (Sutton et al., 2020). An individual professional medical perspective is not enough to meet the needs of patients. Multidisciplinary teams (MDTs) or integrated care teams may participate in disease treatment in practical clinical processes (Kodner and Spreeuwenberg, 2002).

Recently, large language models (LLMs), owing to their powerful reaLsoning and knowledge synthesis capabilities, have demonstrated promising potential in emulating the roles of clinicians and supporting medical decision-making (Tang et al., 2023; Kim et al., 2024; Chen et al., 2025). Multiagent collaboration (MAC) based on LLMs has emerged as a key paradigm, enhancing the reasoning performance of individual agents through collective deliberation. For instance, (Tang et al., 2023) verified that a training-free collaboration framework in which multiple LLM-based agents simulate a multidisciplinary medical team through role-playing and multi-round discussions, and achieved strong performance across medical question answering (QA) datasets. In addition, (Chen et al., 2025) further leveraged medical multiagents and implemented cumulative consultation strategies using retrieval augmentation generation (RAG), which enhances model outputs by retrieving external medical knowledge to support clinical reasoning and improve diagnostic accuracy. Some multi-LLM debate frameworks are also closely related to collaboration (Kaesberg et al., 2025; Chen et al., 2023b; Abdelnabi et al., 2024; Liang et al., 2023). Among them, a framework for iterative collaboration between agents to make decisions is proposed, which stimulates higher quality answers (compared to a single model) by involving multiple models in the discussion. These works explore the potential application of LLMs and the possibility of their use in medical MDT decisions.

043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

Although these MAC methods enable agents to tackle problems that are difficult or unsolvable by a single agent by learning new contexts and actions through interactions with peers or known information, the challenge remains unresolved. It mainly stems from the use of static, pre-assigned roles based on inherent domain knowledge, which limits the system's adaptability during collaboration. As discussions progress, each agent tends to produce



Figure 1: **Comparison of multi-agent collaboration (MAC) strategies in medical decision-making.** (a) *Problem-driven MAC* (Kim et al., 2024; Yang et al., 2024b) uses predefined question-complexity tiers (easy, medium, hard) to assemble static single- or multi-tier expert teams (b) *Observation-driven MAC* (Chen et al., 2023c,a) dynamically analyzes task and role characteristics from initial observations to optimize expert recruitment for each question. (c) Our proposed *Knowledge-driven MAC* adaptively expands the expert team during discussion by detecting knowledge gaps (KG), enabling scalable and flexible collaboration for complex, cross-domain clinical scenarios.

increasingly fine-grained analyses within its fixed specialty. For example, in the evaluation of a patient presenting with chest pain, a radiology agent may focus solely on imaging findings suggestive of pulmonary embolism, while a cardiology agent may emphasize electrocardiogram (ECG) changes indicative of myocardial infarction. Without a mechanism to reconcile or adapt these perspectives, the collaboration degenerates into a juxtaposition of isolated preferences rather than a convergent diagnostic consensus. This fragmentation undermines the effectiveness of consensus strategies and restricts the system's ability to dynamically incorporate broader context or cross-domain reasoning.

Recent studies have attempted to improve MAC flexibility by incorporating novel expert recruitment strategies. For instance, problem-driven MAC (Kim et al., 2024) (Figure 1a) assigns expert teams based on question complexity, while observation-driven MAC (Chen et al., 2023a,c) (Figure 1b) selects experts according to task and role analysis. However, these methods still rely on static or pre-optimized expert pools and cannot adapt during multi-round interactions. As a result, even when new, fine-grained insights emerge over multiple discussion rounds, no new experts are brought in. The limitations in these works still hinder truly scenario-specific collaboration, especially in dynamic and diverse clinical environments.

113To alleviate this, as illustrated in Figure 1c, we114propose a Knowledge-driven Adaptive Multi-115Agent Collaboration (KAMAC) framework for116enhancing medical decision-making. Specifically,117KAMAC dynamically increases the number of

medical expert team members required for patients by exploring additional expert knowledge during the discussion process. KAMAC begins with an initial consultation involving one or more experts. It then engages in a knowledge-driven collaborative discussion, which assesses whether additional expertise is needed by detecting knowledge gaps (KG) and dynamically recruits appropriate experts to fill the knowledge gaps, enabling scalable and flexible collaboration for complex, cross-domain clinical scenarios. Finally, a moderator is responsible for reviewing updated agent comments to complete the decision-making process. Such progressive collaboration and flexible team expansion allow the model to adaptively allocate resources and produce more accurate, context-aware decisions. In contrast to prior methods, the proposed method enables the system to adapt to the evolving clinical treatment in the real world and provide more nuanced and comprehensive support to patients.

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

Our contributions include three folds:

- 1. We propose the KAMAC framework that dynamically extends a single expert agent into multiple expert agents to form a multi-disciplinary team for medical decisionmaking.
- 2. We design a knowledge-driven collaborative discussion mechanism that enables agents to dynamically expand team to fill knowledge gaps, aiming to improve adaptability and decision accuracy in complex clinical scenarios.
- 3. Extensive experiments on two medical benchmarks, MedQA and Progn-VQA, demonstrat-

ing that our KAMAC improves single-agent
and advanced multi-agent collaboration frameworks.

2 Related Work

154

155

156

158

159

160

161

163

164

Advanced LLMs such as GPT-4 (Achiam et al., 2023), DeepSeek (Liu et al., 2024; Guo et al., 2025), and Gemini (Team et al., 2024) have demonstrated strong reasoning capabilities and have been used as agents with considerable computational investment in various medical tasks such as question answering (Kim et al., 2024; Tang et al., 2023), diagnosis (Zhang et al., 2024), and report generation (Thawakar et al., 2024; Hyland et al., 2023). We list two main related areas of work:

LLM-Based Agentic Medical Decision-Making 165 Medical decision-making systems leverage multiple LLM "experts," each assigned a predefined clinical specialty to mimic real-world multidisci-168 plinary teams. Early work demonstrated that con-169 sensus among expert agents yields higher diagnos-170 tic accuracy than any single model or simple vot-171 ing schemes (Tang et al., 2023; Liang et al., 2023; 172 Chen et al., 2023b). Some recent works mainly 173 focused on diagnostic findings (Kim et al., 2024; 174 Li et al., 2024b) and knowledge integration (Xiong 175 et al., 2024; Nori et al., 2023; Kim et al., 2024). 176 MediQ (Li et al., 2024b) designs a system to seek 177 methods to guide the deepening of interactions between patients and experts. For instance, (Kim 179 180 et al., 2024) verified that expert collaboration has better accuracy for medical decision-making than 181 a single expert, showcasing that consensus is su-182 perior to a voting strategy in various clinical applications. More recently, MDteamGPT (Chen et al., 185 2025) adds a leader agent, historical dialogues, and RAG to integrate information and supplementation 186 strategies to assist in decision-making.

Multi-Agent Collaboration in Medical Decision-188 Making Researchers have demonstrated that 189 multi-agent collaborative research can enhance the reasoning capabilities of these LLMs (Yue et al., 191 2024; Wang et al., 2024; Li et al., 2024a). Well-192 designed strategies can enhance autonomous multi-193 agent systems for task-solving capabilities, such 194 195 as debate (Chan et al., 2023; Abdelnabi et al., 2024), consensus (Kaesberg et al., 2025; Chen 196 et al., 2023b), conflict-solving, generation/evolu-197 tion (Yuan et al., 2024; Chen et al., 2023c,a), and encouragement (Liang et al., 2023; Tran et al., 199

2025). Some multi-LLM debate frameworks are also closely related to collaboration (Kaesberg et al., 2025; Chen et al., 2023b; Abdelnabi et al., 2024; Liang et al., 2023). Among them, a framework for iterative collaboration between agents to make decisions is proposed, which stimulates higher quality answers (compared to a single model) by involving multiple models in the discussion. These works explore the potential application of LLMs and the possibility of their use in medical MDT diagnostics. Although assigning experts can effectively improve the performance of specific tasks, the rationality of expert assignment in multi-agent collaboration is still insufficient. (Chen et al., 2023a,c) introduces optimal expert generation strategies in the initial expert recruitment stage, but it does not consider the relationship between expert knowledge and cooperation between experts. (Yuan et al., 2024) introduces a dynamic evolution strategy for the existing experts but relies on a large initial population and requires additional investment. These limitations make it unsuitable for medical decision-making.

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

3 Method

3.1 Overview

Figure 2 presents the KAMAC framework, which comprises three main stages: (1) Initial Consultation: KAMAC begins with a single/multiple expert agents, which evaluate the case and provide initial feedback for ongoing discussion; (2) Knowledgedriven Collaborative Discussion: Agents engage in a structured, knowledge-guided dialogue to determine whether further expertise is required and then adaptively expands the team and promotes structured discussions among agents, guided by domain knowledge and the evolving diagnostic context, and (3) Decision Making: A designated moderator coordinates the final decision process by initiating a voting mechanism among agents. The pseudocode of KAMAC is shown in Algorithm 1. More details in all prompts refer to Appendix B.

3.2 Initial Consultation

Given a clinical problem Q, KAMAC first performs an initial consultation by recruiting one or more expert agents from a predefined expert pool. These agents represent diverse clinical roles (*e.g.*, radiologist, cardiologist) and are selected based on their relevance to the query using an expert recruitment prompt P_1 . Each recruited agent indepen-



Figure 2: Schematic diagram of Knowledge-driven Adaptive Multi-Agent Collaboration (KAMAC) framework for medical decision-making. The KAMAC includes three parts: (a) Initial Consultation: One or more expert agents (*e.g.*, radiologist, pathologist) are selected based on the clinical question to provide initial assessments; (b) Knowledge-driven Collaborative Discussion: Agents iteratively exchange views to refine reasoning. If a knowledge gap is detected, KAMAC dynamically recruits additional specialists, and the expanded team continues the dialogue until consensus is reached or the round limit is met; and (c) Final Decision Making: A moderator reviews all agent responses and produces the final answer. The symbols \checkmark and \checkmark indicate agreement/disagreement with the current expert's comment, respectively. Only when a disagreement occurs, ($\checkmark i \rightarrow j$) or ($\bigstar i <-> j$) is used to denote a one-way or two-way discussion between expert *i* and expert *j*, respectively.

dently analyzes the problem using an initial assessment prompt P_2 , producing diagnostic opinions or treatment suggestions. The individual responses are aggregated into a consolidated feedback signal, which serves as the basis for initiating collaborative discussion in the next stage. This step simulates a typical initial clinical encounter, where specialists offer their perspectives before deliberation begins.

251

262

266

3.3 Knowledge-driven Collaborative Discussion

In this stage, KAMAC facilitates multi-round, knowledge-driven discussions among the recruited expert agents. Each round begins with agents exchanging their views based on the evolving shared context. Using the agent interaction prompt P_3 , they critique each other's responses, resolve inconsistencies, and collaboratively refine their reasoning and comments. At the end of each round, the currently assigned experts are prompted to assess whether a knowledge gap (KG) remains—that is, whether their collective expertise is sufficient to fully address the problem. This self-assessment is facilitated by the KG detection prompt P_4 , which takes as input the current discussion and feedback. If a gap is detected, KAMAC triggers expert recruitment by issuing a targeted recruitment prompt P_5 , allowing the system to enlist additional domain-specific agents to address the identified deficiency.

The newly recruited agents receive contextual examples (*i.e.*, the current discussion history) as few-shot input and respond to the original question using the assessment prompt P_2 , conditioned on the ongoing feedback. Their outputs are appended to the current feedback buffer and integrated into the group discussion in the subsequent round. This

268

269

270

271

Algorithm 1: Knowledge-driven Adaptive Multi-Agent Collaboration (KAMAC) Decision-making **Input:** Problem Q **Result:** Answer ans. 1 Initialize: KAMAC \leftarrow []. ▷ Define prompts. More details in all prompts refer to Appendix B. 2 $r \leftarrow 1$, Consensus \leftarrow False, $KG \leftarrow$ False. **3** P_1 : Expert Recruitment Prompt; P_2 : Initial Assessment Prompt. 4 P_3 : Agent Interaction Prompt; P_4 : KG Prompt for Recruited Experts. 5 P_5 : KG Prompt for Expert Recruitment; P_6 : Agent Update Prompt. 6 P_7 : Final Decision Prompt. > Initial consultation. \triangleright Recruit N clinician agents. 7 (Agent₁, Agent₂, ..., Agent_N) \leftarrow Recruit(Q, KAMAC, P_1) ▷ Clinician agents consist of a multi-disciplinary team. 8 KAMAC \leftarrow (Agent₁, Agent₂, \cdots , Agent_N) ▷ Initial assessment. 9 ($Option_1, Option_2, \dots, Option_N$) $\leftarrow Chat(\mathcal{Q}, KAMAC, P_2)$ ▷ Concat all options as feedback. 10 Feedback \leftarrow Concat(Option₁, Option₂, \cdots , Option_N) > Knowledge-driven collaborative discussion. 11 while $r \leq R$, and not Consensus, and not KG do > Exchange agent's comments and determine the consensus. 12 Consensus, Feedback \leftarrow Chat(Q, KAMAC, Feedback, P_3) ▷ Assess whether any additional specialist is needed to fill a knowledge or diagnostic gap. $KG \leftarrow \text{Chat}(\mathcal{Q}, \text{KAMAC}, Feedback, P_4)$ 13 if KG then 14 > Expert recruitment for recruiting additional experts during the discussion. 15 $(Agent_{N+1}, Agent_{N+2}, \cdots, Agent_M) \leftarrow \text{Recruit}(\mathcal{Q}, \text{KAMAC}, P_5)$ ▷ Review all options and provide comments as feedback. 16 $(Option_{N+1}, Option_{N+2}, \cdots, Option_M) \leftarrow Chat((Agent_{N+1}, Agent_{N+2}, \cdots,$ $Agent_M$, Feedback, P_2) $Feedback \leftarrow Concat(Feedback, Option_{N+1}, Option_{N+2}, \cdots, Option_M)$ 17 ▷ Exchange agent's comments and determine the consensus. Consensus, Feedback \leftarrow Chat(\mathcal{Q} , (Agent_{N+1}, Agent_{N+2}, \cdots , Agent_M), 18 $Feedback, P_3$) 19 ▷ Update KAMAC. $KAMAC \leftarrow (KAMAC, \cdots, Agent_{N+1}, \cdots, Agent_M)$ 20 $KG \leftarrow False$ 21 ▷ Update agent's comments. $Feedback \leftarrow Chat(\mathcal{Q}, KAMAC, Feedback, P_6)$ 22 $r \leftarrow r+1$ 23 ▷ Make the final decision by LLMs. **24** ans \leftarrow Moderator(\mathcal{Q} , Feedback, P_7)

25 return ans

recursive process allows progressive team expansion, enabling KAMAC to dynamically adapt to the evolving complexity of the diagnostic scenario. ing the initial and newly recruited ones, update their reasoning using the agent update prompt P_6 , which ensures alignment with the current collective context. This process continues until either (1) a 289 290 291

288

Throughout the discussion, all agents, includ-

294

297

301

303

- 305
- 306 307
- ~~~

consensus is reached via iterative agreement checks using P_3 , or (2) a maximum number of discussion rounds R is reached.

3.4 Decision Making

The collaborative discussion continues until either consensus is reached or the maximum number of rounds R is exhausted. In the final stage, KAMAC invokes a moderator agent, typically a generalpurpose LLM, to generate the final decision. The moderator receives the latest set of agent comments and the full discussion history and synthesizes a response via a decision prompt (P_7).

4 Experiments

4.1 Datasets

308To evaluate the proposed KAMAC framework, we309conduct experiments on the testing sets of two pub-310licly available medical question answering (QA)311datasets: MedQA (Jin et al., 2021) and Progn-312VQA (Welch et al., 2023).

313MedQAWe use all 1273 samples in the testing314set. This dataset describes the United States Medi-315cal Licensing Examination and includes questions,316multiple-choice questions, and answers.

Progn-VQA We use all 750 Visual Question An-317 swering (VQA) pairs in the testing set. The dataset 318 319 includes head and neck cancer Computed Tomography (CT) image volumes collected from 2005-320 2017 treated with definitive radiotherapy at the 321 University Health Network in Toronto, Canada. It also contains the corresponding regions of interest 323 (ROIs) and structured patient information in RT-324 STRUCT format with standardized descriptions, 325 including demographic, clinical, and treatment in-326 formation based on the 7th edition TNM staging system and AJCC (American Joint Committee on Cancer). The dataset contains patient information, CT image volumes and ROIs, and the patient's survival status at the last follow-up. Please see Table 4 332 for understanding the clinical and imaging information used in the dataset. For CT input, we selected the axial slice with the largest cross-sectional area of the ROI. More details on the input clinical and imaging variables are provided in Appendix A. 336

4.2 Implementation Details

We use GPT-4.1-mini¹ as the primary model for 338 all experiments, with the temperature set to 0 to 339 ensure deterministic outputs. In addition, we compare our proposed method with DeepSeek-R1 (Guo 341 et al., 2025), as shown in Table 2. For each medical 342 question, we store the corresponding chat history 343 in a local file. When revisiting the same question, 344 the system loads the saved file to regenerate con-345 sistent initial medical comments from each role 346 before resuming the collaborative discussion. The 347 final decision is made solely based on the proposed collaboration method. The maximum number of 349 discussion rounds R is set to 3. The initial number 350 of experts is set to 1. We select GPT-4.1-mini due 351 to its strong medical reasoning capabilities, low 352 latency, predictable computational cost, and fully 353 deterministic behavior. These advantages make it 354 preferable for our controlled evaluation setting, in 355 contrast to larger models such as GPT-4 or retrieval-356 enhanced models like DeepSeek-R1, which often 357 entail higher overhead and less consistent outputs.

337

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

4.2.1 Comparison Methods

The compared methods include: (1) Single-agent, which uses an LLM for decision-making, where the question and the answer template are input to output an answer, (2) Chain of Thought (CoT) (Wei et al., 2022), which combines the single-agent backbone with a step-by-step prompt to conduct analysis and decision-making, (3) Majority Voting, which is used in multi-agent decision-making methods (Chen et al., 2023b; Yang et al., 2024a; Kaesberg et al., 2025) for making final decision with more that 50% votes. (4) Consensus, which is also adopted in (Kaesberg et al., 2025; Kim et al., 2024). (5) MDAgents (Kim et al., 2024) is an advanced multi-agent framework that performs problem-driven expert recruitment, MAC and consensus decision to output the final results.

4.3 Evaluation Metrics

We evaluate the proposed method using four standard metrics: accuracy (Acc), precision (Prec), specificity (Spec), and recall score (Recall).

5 Results and Analysis

5.1 Comparisons with State-of-the-Arts

In Table 1, the proposed method achieves improved results on four metrics compared to multiple meth-

¹https://openai.com/index/gpt-4-1/

Table 1: Main results on four common metrics across MedQA and Progn-VQA datasets, evaluated using GPT-4.1mini. **Bold** values indicate the best performance. Here, 'SA' means the single-agent methods and 'MA' means the multi-agent methods. Gray-highlighted cells indicate the average score.

Methods	Types	MedQA				A	Progn-VQA				A
		Acc	Prec	Spec	Recall	Avg	Acc	Prec	Spec	Recall	Avg
Single-agent	SA	79.50	79.65	94.86	79.36	83.34	86.00	86.28	14.79	97.21	71.07
+CoT	SA	84.21	84.82	96.03	84.02	87.27	84.67	86.29	15.52	97.32	70.95
Majority Voting	MA	86.49	86.93	96.60	86.38	89.10	86.27	86.12	12.17	99.84	71.10
Consensus	MA	80.68	80.70	95.15	80.59	84.28	86.86	86.81	31.85	98.86	76.09
MDAgents	MA	87.74	87.92	96.92	87.55	90.03	87.01	88.83	33.70	96.21	76.44
KAC-MAF	MA	88.14	88.30	97.02	88.11	90.39	87.20	89.79	40.52	95.74	78.31

Table 2: Performance comparison of Baseline and KAC-MAF on MedQA and Progn-VQA using DeepSeek-R1 and GPT-4.1-mini across four metrics and their average. Gray-highlighted cells indicate the average score, with relative improvements shown in small colored text. Where 'Baseline' means single-agent+CoT.

Mothod	MedQA				A-1-0	Progn-VQA				A-1-0	
Methou	Acc	Prec	Spec	Recall	Avg	Acc	Prec	Spec	Recall	Avg	
Baseline (DeepSeek-R1)	88.14	88.12	97.03	88.00	90.32	77.87	88.11	37.07	85.33	72.10	
KAC-MAF (DeepSeek-R1)	89.63	89.53	97.41	89.50	91.52(+1.20)	86.13	88.41	31.03	96.21	75.45 (+3.35)	
Baseline (GPT-4.1-mini)	84.21	84.82	96.03	84.02	87.27	84.67	86.29	15.52	97.32	70.95	
KAC-MAF (GPT-4.1-mini)	88.14	88.30	97.02	88.11	90.39 (+3.12)	87.20	89.79	40.52	95.74	78.31 (+7.36)	

Table 3: Discussion on the number of initial agents on the MedQA and Progn-VQA datasets. Gray-highlighted cells indicate the average score.

Initial Agents	MedQA				A	Progn-VQA				Ana
Number	Acc	Prec	Spec	Recall	Avg	Acc	Prec	Spec	Recall	Avg
1	88.14	88.30	97.02	88.11	90.39	87.20	89.79	40.52	95.74	78.31
5	80.28	80.31	95.06	80.13	83.95	89.10	89.54	35.43	96.69	77.69

ods on the MedQA dataset. For the Progn-VQA dataset, the proposed method achieves better results on the Acc, Prec, and Spec metrics. In addition, KAMAC leverages knowledge-driven prompts to detect KG and expand experts to form multi-agent collaborative discussions. Focusing on multi-agentbased methods, both the majority voting and consensus are set to five experts, while MDAgents adopts a single agent, a multi-disciplinary team with five experts, and an integrated care team with nine experts. Compared with them, the proposed method can achieve better results, which demonstrates that our method overcomes the limitation of knowledge in the single-agent model and has a more suitable multi-disciplinary team to enhance multi-agent reasoning and collaboration.

385

387

389

394

400

401

402

403

404

In Table 2, we further evaluate our method on another model DeepSeek-R1. In our method, the initial number of experts is set to 1, which is consistent with the baseline method (single-agent + CoT), but the experimental results are better than the baseline method. This improvement shows that our method can be generalized to more LLM models. In addition, this improvement aligns with the actual clinical treatment scenario, where clinical treatment allows the dynamic addition of experts according to the patient's clinical treatment situation, thereby carrying out more effective treatment. It contributes to optimizing the best treatment options and limited clinical resources in medical applications.

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

5.2 Discussion

Results of Different Multi-agent Methods In this part, we focus on the main table (Table 1), where our method introduces the knowledge-driven prompt for assessing whether additional expertise is needed by detecting knowledge gaps (KG) and dynamically recruiting appropriate experts to fill the knowledge gaps. The reason why we choose MDAgent as the comparison method is that MDAgent is a problem-driven MAC that adopts an adap-



(a) Number of Initial Experts: 5

Figure 3: Histogram illustrating the impact of initial expert settings on the final top-30 expert distribution in our method on the MedQA dataset. "Count" denotes the total frequency of each expert type. An 80% overlap in expert types is observed between the 1- and 5-expert settings.

tive expert recruitment strategy for multi-agent collaborative discussion. Our method is a knowledge-

425

426

driven adaptive expert recruitment. By recruiting 427 experts during the first consultation and dynam-428 ically detecting KG in the recruited experts, our 429 approach outperforms the MDAgent method on 430 both datasets. On average, 1.28 and 2.41 experts 431 were involved per case, compared to 2.41 and 4.34 432 for MDAgent, resulting in reductions of 67% and 433 56%, respectively. 434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

Distribution of Experts of KAMAC In Figure 3, we show the distribution of the proposed KAMAC under different numbers of initial experts. Among the top 30 expert distributions, there are 24 overlapping experts between the settings of Initial Experts: 1 and Initial Experts: 5, resulting in an 80% overlap. However, the total number of experts involved in the Initial Experts: 1 setting is lower than that in the Initial Experts: 5 setting. Moreover, as shown in Table 3, the performance with Initial Experts: 1 is superior to that with Initial Experts: 5. These two quantitative results validate the effectiveness of the proposed method, and knowledge-driven multiagent framework can effectively identify which specialists are needed to diagnose a patient's condition and promptly recruit the appropriate experts to provide comprehensive treatment.

6 Conclusion

This work presents KAMAC, a knowledge-driven adaptive multi-agent collaboration framework that brings structured, dynamic reasoning into medical decision-making with LLMs. By allowing agents to actively assess their own limitations and request additional expertise when needed, KAMAC overcomes the rigidity of traditional multi-agent setups and more faithfully mirrors real-world clinical workflows. Our experiments on two real-world medical QA benchmarks demonstrate that KA-MAC consistently outperforms both single-agent and existing multi-agent baselines. Beyond accuracy improvements, KAMAC offers deeper insights into AI collaboration: decision quality improves not merely through more parameters or agents, but through adaptive, feedback-driven interaction grounded in knowledge awareness. This framework brings multi-agent LLM systems closer to real-world clinical workflows, where expert composition evolves with case complexity. Future directions include modeling agent uncertainty and integrating clinician-in-the-loop feedback to further support real-time deployment in medical environments.

477 Limitations

While KAMAC demonstrates promising results, it 478 has several limitations. The current framework fo-479 cuses on textual and imaging inputs; future work 480 could incorporate additional modalities such as ge-481 nomic or longitudinal clinical data to support a 482 483 wider range of medical tasks. Although KAMAC achieves strong performance without fine-tuning 484 the underlying LLMs, domain-specific fine-tuning 485 may further improve accuracy and agent-role fi-486 delity. However, this would introduce significant 487 computational overhead and is challenged by the 488 scarcity of high-quality, labeled medical data. Bal-489 ancing accuracy gains with efficiency and data 490 availability remains an important direction for fu-491 ture fine-tuning efforts. 492

References

493

494

495

496

497

498

499

501

502

503

504

508

509

510

511

512

513

514

515

516

517

518

519

520

521

523

524

525

526

527

- Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. 2024. Cooperation, competition, and maliciousness: Llm-stakeholders interactive negotiation. *Advances in Neural Information Processing Systems*, 37:83548–83599.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.
- Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F Karlsson, Jie Fu, and Yemin Shi. 2023a. Autoagents: A framework for automatic agent generation. *arXiv preprint arXiv:2309.17288*.
- Justin Chih-Yao Chen, Swarnadeep Saha, and Mohit Bansal. 2023b. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. *arXiv preprint arXiv:2309.13007*.
- Kai Chen, Xinfeng Li, Tianpei Yang, Hewei Wang, Wei Dong, and Yang Gao. 2025. Mdteamgpt: A self-evolving llm-based multi-agent framework for multi-disciplinary team medical consultation. *arXiv* preprint arXiv:2503.13856.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, and 1 others. 2023c. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. arXiv preprint arXiv:2308.10848, 2(4):6.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *arXiv preprint arXiv:2501.12948*. 528

529

530

531

532

534

535

536

537

538

539

540

541

542

543

544

545

546

547

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

- Stephanie L Hyland, Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Mercy Ranjit, Anton Schwaighofer, Fernando Pérez-García, Valentina Salvatelli, Shaury Srivastav, Anja Thieme, and 1 others. 2023. Maira-1: A specialised large multimodal model for radiology report generation. arXiv preprint arXiv:2311.13668.
- 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.
- Lars Benedikt Kaesberg, Jonas Becker, Jan Philip Wahle, Terry Ruas, and Bela Gipp. 2025. Voting or consensus? decision-making in multi-agent debate. *arXiv preprint arXiv:2502.19130*.
- Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik S Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, and Hae W Park. 2024. Mdagents: An adaptive collaboration of llms for medical decision-making. *Advances in Neural Information Processing Systems*, 37:79410–79452.
- Dennis L Kodner and Cor Spreeuwenberg. 2002. Integrated care: meaning, logic, applications, and implications–a discussion paper. *International journal of integrated care*, 2:e12.
- Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. 2024a. More agents is all you need. *Transactions on Machine Learning Research*.
- Stella Li, Vidhisha Balachandran, Shangbin Feng, Jonathan Ilgen, Emma Pierson, Pang Wei W Koh, and Yulia Tsvetkov. 2024b. Mediq: Question-asking Ilms and a benchmark for reliable interactive clinical reasoning. *Advances in Neural Information Processing Systems*, 37:28858–28888.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2023. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*.

582 583 Reed T Sutton, David Pincock, Daniel C Baumgart,

Daniel C Sadowski, Richard N Fedorak, and Karen I

Kroeker. 2020. An overview of clinical decision

support systems: benefits, risks, and strategies for

Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming

Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and

Mark Gerstein. 2023. Medagents: Large language

models as collaborators for zero-shot medical reason-

Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan

Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,

Damien Vincent, Zhufeng Pan, Shibo Wang, and 1

others. 2024. Gemini 1.5: Unlocking multimodal

understanding across millions of tokens of context.

Omkar Chakradhar Thawakar, Abdelrahman M

Shaker, Sahal Shaji Mullappilly, Hisham Cholakkal,

Rao Muhammad Anwer, Salman Khan, Jorma Laak-

sonen, and Fahad Khan. 2024. Xraygpt: Chest radiographs summarization using large medical visionlanguage models. In Proceedings of the 23rd work-

shop on biomedical natural language processing,

Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O'Sullivan, and Hoang D Nguyen. 2025. Multi-agent collaboration mech-

Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang,

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten

Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,

and 1 others. 2022. Chain-of-thought prompting elic-

its reasoning in large language models. Advances

in neural information processing systems, 35:24824-

ML Welch, S Kim, A Hope, SH Huang, Z Lu, J Marsilla, M Kazmierski, K Rey-McIntyre, T Patel,

Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong

Joshua C Yang, Damian Dalisan, Marcin Korecki, Ca-

rina I Hausladen, and Dirk Helbing. 2024a. Llm voting: Human choices and ai collective decisionmaking. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, volume 7, pages

Zhang. 2024. Benchmarking retrieval-augmented

generation for medicine. In Findings of the Association for Computational Linguistics ACL 2024, pages

(radcure). The Cancer Imaging Archive.

B O'Sullivan, and 1 others. 2023. Computed tomography images from large head and neck cohort

and James Zou. 2024. Mixture-of-agents enhances

large language model capabilities. arXiv preprint

arXiv preprint

success. NPJ digital medicine, 3(1):17.

ing. arXiv preprint arXiv:2311.10537.

arXiv preprint arXiv:2403.05530.

anisms: A survey of llms.

pages 440-448.

arXiv:2501.06322.

arXiv:2406.04692.

24837.

6233-6251.

1696-1708.

- 585

- 599
- 601
- 604

610

611

612

613 614

615 616

618

619 620

- 623

- 627

628

630 631

632

636

Zhe Yang, Yichang Zhang, Yudong Wang, Ziyao Xu, Junyang Lin, and Zhifang Sui. 2024b. Confidence vs critique: A decomposition of self-correction capability for llms. arXiv preprint arXiv:2412.19513.

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

- Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Dongsheng Li, and Deqing Yang. 2024. Evoagent: Towards automatic multi-agent generation via evolutionary algorithms. In NeurIPS Workshop on Open-World Agents.
- Ling Yue, Sixue Xing, Jintai Chen, and Tianfan Fu. 2024. Clinicalagent: Clinical trial multi-agent system with large language model-based reasoning. In Proceedings of the 15th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics, pages 1–10.
- Kai Zhang, Rong Zhou, Eashan Adhikarla, Zhiling Yan, Yixin Liu, Jun Yu, Zhengliang Liu, Xun Chen, Brian D Davison, Hui Ren, and 1 others. 2024. A generalist vision-language foundation model for diverse biomedical tasks. Nature Medicine, pages 1-13.

10

A More Details for Progn-VQA Dataset

According to the settings of (Welch et al., 2023),
we provide the clinical and imaging information
required for prognosis in Table 4. This information
can fully describe the situation of patients with
head and neck cancer.

B Prompt Template

657

663

664

665

666

667

668

We provide all prompts in our multi-agent medical decision-making framework, including expert recruitment, initial comments, collaborative discussion, and knowledge-driven prompts. For a singleagent setting, you can refer to (Kim et al., 2024).

Variable	Description			
Age	Patient age			
Sex	Patient sex			
ECOG PS	ECOG Performance Status			
Smoking PY	Cumulative smoking exposure (pack-years)			
Smoking Status	Smoking status at initial consultation			
Ds Site	Primary disease (cancer) site			
Subsite	Subsite of the primary tumor			
Т	Tumor size and extent (AJCC 7th edition T category)			
Ν	Regional lymph node involvement (AJCC 7th edition N category)			
М	Distant metastasis (AJCC 7th edition M category)			
Stage	Overall stage group (AJCC 7th edition)			
Path	Pathological diagnosis or histological subtype			
HPV	HPV status of the tumor, determined by p16 IHC with or without			
	confirmation by HPV DNA PCR (blank if unavailable)			
Tx Modality	Treatment modality			
Chemo?	Whether concurrent chemoradiotherapy was administered			
Dose	Total radiotherapy dose delivered (in Gy)			
Fx	Number of radiotherapy fractions			
Local	Indicator of local recurrence			
Regional	Indicator of regional recurrence			
Distant	Indicator of distant metastasis			
2nd Ca	Indicator of second primary cancer			
ContrastEnhanced	Indicator of whether contrast-enhanced imaging was used			

Table 4: Descriptions of clinical and imaging variables included in the Progn-VQA dataset (Welch et al., 2023).

Expert Recruitment Prompt for MedQA (P_1)

System: You are an experienced medical expert who recruits a group of experts with diverse identities and asks them to discuss and solve the given medical query.

User:

Question: {{QUESTION}}

You can recruit {{NUM_AGENTS}} experts in different medical expertise.

Considering the medical question and the options for the answer, what kind of experts will you recruit to better make an accurate answer?

Also, you need to specify the communication structure between experts (e.g., Pulmonologist == Neonatologist == Medical Geneticist == Pediatrician > Cardiologist), or indicate if they are independent.

For example, if you want to recruit five experts, your answer can be like:

1. Pediatrician - Specializes in the medical care of infants, children, and adolescents. - Hierarchy: Independent

2. Cardiologist - Focuses on the diagnosis and treatment of heart and blood vessel-related conditions. - Hierarchy: Pediatrician > Cardiologist

3. Pulmonologist - Specializes in the diagnosis and treatment of respiratory system disorders. - Hierarchy: Independent

4. Neonatologist - Focuses on the care of newborn infants, especially those who are born prematurely or have medical issues at birth. - Hierarchy: Independent

5. Medical Geneticist - Specializes in the study of genes and heredity. - Hierarchy: Independent Please answer in the above format, and do not include your reason.

Expert Recruitment Prompt for Progn-VQA (P_1)

System: You are an experienced medical expert who recruits a group of experts with diverse identity and ask them to discuss and solve the given medical query.

User:

Question: {{QUESTION}}

Considering the medical question and the options for the answer, what kinds of experts will you recruit to better make an accurate decision? You also need to clearly specify the communication structure between experts or indicate if they are independent.

You must recruit exactly the following {{NUM_AGENTS}} experts, with no substitutions, no additional experts, and no omissions:

(e.g., Radiation Oncologist == Medical Oncologist == Pathologist == Surgical Oncologist (Recurrence/Secondary Cancers) == Targeted Therapy Expert),

Please strictly follow the format shown below, without adding any extra explanation or reasoning. Format example if recruiting {{NUM_AGENTS}} experts:

1. Radiation Oncologist - Your expertise is strictly limited to radiation therapy planning and dosing for head and neck squamous cell carcinoma, especially HPV-positive cases.

- Hierarchy: Radiation Oncologist == Medical Oncologist

2. Medical Oncologist - Your expertise is strictly limited to systemic therapy decisions, including chemotherapy and immunotherapy in head and neck cancers.

- Hierarchy: Medical Oncologist == Radiation Oncologist

3. Surgical Oncologist (Recurrence/Secondary Cancers)—Your expertise is strictly limited to evaluating surgical options for recurrent or secondary malignancies in head and neck cancers.

- Hierarchy: Surgical Oncologist == Pathologist"

4. Pathologist - Your expertise is strictly limited to pathological diagnosis of head and neck squamous cell carcinoma, HPV status evaluation, and margin assessment post-surgery.

- Hierarchy: Pathologist == Surgical Oncologist

5. Targeted Therapy Expert - Your expertise is strictly limited to clinical application of EGFR inhibitors and novel agents targeting HPV-positive tumors.

- Hierarchy: Targeted Therapy Expert -> Medical Oncologist

Your answer must conform exactly to the format above.

Chain-of-thought Prompt for Initial Assessment (P_2)

System: You are a {{ROLE}} who {{DESCRIPTION}}. Your job is to collaborate with other medical experts in a team.

User: {{VISUAL COT INSTRUCTION}} (Optional)

Given the examplers, as a {{ROLE}}, please return your answer to the medical query among the options provided. You are not allowed to switch to any other medical specialty.

{{FEWSHOT_EXAMPLERS}}

Question: {{QUESTION}}

Your answer should be in the format below.

{{answer_template}}

Visual Chain-of-thought Prompt for Head and Neck CT Scan (Optional, only be used when input data include images.)

User: You will be provided with a head and neck CT scan that includes one or more masked regions of interest (ROIs). Alongside the scan, one or more 3D bounding box coordinates will be supplied, each defining specific volumetric regions within the scan. These coordinates identify either organs, disease regions, or cellular structures. Each bounding box is defined by its minimum and maximum values along the z, y, and x axes and is normalized relative to the original image size.

The given bounding box coordinates are: {{BBOX_COORDS}}.

Task Instructions:

1. **Initial Assessment**: Carefully analyze the CT scan image (without using the bounding box data). Describe any visible anatomical structures, patterns, abnormalities, and note the characteristics of the masked regions of interest (ROIs).

Do not use the bounding box data at this stage.

2. **Mapping Bounding Boxes**: Consider the bounding box coordinates and map them to the corresponding areas within the scan.

3. **Clinical Reasoning**: Summarize the patient's clinical context and findings in a clear, structured bullet-point format and reason through the patient's condition step by step.

4. **Integrated Conclusion**: Combine your findings from the image analysis, bounding box mapping, and masked ROI to concisely synthesize your final clinical impression.

Be thorough and precise in both your image-based observations and your clinical reasoning.

Agent Interaction Prompt (P_3)

User: Earlier in this conversation, a set of discussion opinions from other medical experts on your team was provided. Please do not forget those earlier opinions.

Now, additional new opinions have been provided. Considering both the earlier and the latest opinions together, please indicate whether you want to talk to any additional expert (yes/no). Opinions: {{ASSESSMENT}}

Knowledge-driven Prompt for Recruited Experts (P_4)

User: You are part of the team: {{AGENTS}}. Earlier in this conversation, a set of discussion opinions from one or more medical experts on your team was provided. Please carefully review that information now. Based on your professional boundaries, determine whether there is a knowledge limitation or missing perspective that requires support from another specialist.

Please answer yes or no.

If yes, specify the type of expert needed and provide a short reason. Be specific and consider the multidisciplinary needs involved in managing complex patient information (e.g., diagnostic imaging, supportive care, pathology review, and other medical expertise).

It is acceptable to recognize areas of expertise already covered by current team members ({{AGENTS}}).

Do not recommend a specialist if their expertise is already represented in the team.

Knowledge-driven Prompt for Expert Recruitment (P_5)

User: Considering the medical question, discussion options, and the current expert team {{AGENTS}}, identify any that require recruiting new types of experts to ensure an accurate decision (exclude {{AGENTS}}).

You also need to clearly specify the communication structure between experts (e.g. Targeted Therapy Expert -> Medical Oncologist, Medical Oncologist == Radiation Oncologist)" or indicate if the new expert(s) will work independently.

Do not suggest removing, substituting, or duplicating existing experts. Only add new experts if necessary.

Format example if recruiting experts:

1. Medical Oncologist - Your expertise is strictly limited to systemic therapy decisions, including chemotherapy and immunotherapy in head and neck cancers. - Hierarchy: Independent 2. Other Medical Experts.

Your answer must conform exactly to the format above. If the existing expert team comprehensively have covered the necessary expertise for accurate decision, answer: <skip recruitment>

Agent Update Comments after Discussion Prompt (P_6)

User: Now that you've interacted with other medical experts, remind your expertise and the comments from other experts and make your final answer to the given question:{{QUESTION}} Answer: {{ANSWER_TEMPLATE}}

Only output your final answer in the format below:

{{FINAL_ANSWER_TEMPLATE}}

Question: {{QUESTION}}

Final Decision Prompt (P_7)

System: You are a final medical decision maker who reviews all opinions from different medical experts and make final decision.

User: Given each agent's final answer, please review each agent's opinion and make the final answer to the question by taking a majority vote.

Only output your final answer in the format below:

{{FINAL_ANSWER_TEMPLATE}}

Question: {{QUESTION}}