Hierarchical Agent Reflection for Aligning LLM Reasoning with **Clinical Diagnostic Processes**

Anonymous ACL submission

Abstract

Medical diagnosis is a complex, iterative process that relies heavily on clinicians' reasoning and judgment. Traditional models, while able 004 to provide consistent diagnostic results, fail to replicate the reasoning process of clinicians, making their outputs difficult to understand and justify. In this paper, we address this limitation by first generating clinical notes that capture the clinician's diagnostic reasoning. These notes are then used to train a large language model, allowing it to mimic the step-by-step reasoning employed by clinicians during diag-012 nosis. Our method introduces a hierarchical agent reflection mechanism to generate clin-014 ical notes, which deconstructs the diagnostic process into key stages, each handled by specialized agents. This structured approach not 018 only improves the accuracy and reliability of the generated clinical notes but also ensures that the model's reasoning aligns with human clinical practice. Experimental results show that models trained on this data outperform 023 both general-purpose large language models and domain-specific medical models in diagnostic tasks. The proposed method enhances diagnostic transparency and interpretability, offering a valuable tool for AI-assisted clinical decision-making.

1 Introduction

011

034

042

Accurate diagnosis is a critical step in medical practice, as it directly impacts patient outcomes. With the continuous advancement of modern medical technologies, particularly the rise of artificial intelligence and large language models (LLMs), the medical diagnostic process is undergoing transformative change. Research is increasingly focused on leveraging these technologies to assist clinicians in achieving more accurate and efficient diagnoses. Recent studies have demonstrated that LLMs, when operating autonomously, can outperform clinicians in certain diagnostic tasks, underscoring the significant potential of these models (Goh et al., 2024).

However, current LLM-based diagnostic systems primarily offer static responses to clinician inquiries, lacking active engagement in the clinical reasoning process. This limitation restricts their effectiveness as collaborative tools in medical diagnosis, as they do not engage in the dynamic and iterative reasoning processes that clinicians rely on. While recent advancements such as OpenAI's o1 (Jaech et al., 2024) have introduced reasoningaugmented models that excel in domains requiring complex problem-solving, these models still fall short in medical contexts. In particular, diagnostic reasoning in medicine involves nuanced, non-linear decision-making based on a combination of clinical intuition, patient history, and test results. To be truly effective in medical settings, LLMs must not only process vast amounts of data but also replicate the dynamic, step-by-step reasoning that clinicians employ during diagnosis.

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

To bridge this gap and harness the full potential of LLMs, it is essential to align their diagnostic reasoning with clinical reasoning. This alignment can be achieved by fine-tuning the models using clinical notes, which encapsulate the detailed diagnostic processes of clinicians.

A standard clinical note typically includes several key components. First, it contains the patient's medical history, detailing past illnesses, surgeries, allergies, and other relevant health information. Next, it outlines the findings from the physical examination, including the physician's observations and assessments. Following this, the note includes results from auxiliary examinations, such as laboratory tests and imaging studies. After synthesizing and reflecting on this information, the patient's clinical features are documented, culminating in an initial diagnostic impression along with the supporting rationale. Finally, through differential diagnosis, the final diagnosis is determined. However, obtaining high-quality clinical notes is not trivial, mainly due to the significant cost of expert anno-



Figure 1: Previous models cannot reason in a manner akin to clinicians, which makes their decision-making process difficult to interpret. In contrast, our method first generates clinical notes that document the clinician's reasoning process. These notes are then used to train the model, enabling it to reason in a manner similar to clinicians.

tation and the time-intensive nature of the process.This challenge makes it difficult to generate large datasets of clinical notes required for fine-tuning models effectively.

084

091

096

100

101

102

103

104

105

106

108

109

110

111

112

113

To do so, we propose a hierarchical agent reflection mechanism that integrates knowledgeenhancement techniques. We deconstruct the diagnostic process and design agents to simulate the multiple steps a clinician would take when diagnosing with clinical notes. The resulting clinical notes are then used for further training of the model, ensuring that the LLM's diagnostic reasoning resonates with that of clinicians (see Fig. 1). Our framework is designed with a hierarchy of specialized agents, consisting of three foundational agents and one supervisory agent: (1) Information Collection Agent - Extracts and summarizes relevant patient data. (2) Preliminary Diagnosis Agent -Conducts iterative reasoning to generates an preliminary diagnostic hypothesis. (3) Differential **Diagnosis Agent** – Conducts iterative reasoning to refine the differential diagnosis. (4) Coordinator Agent - as the supervisory agent, Oversees and integrates the reasoning outputs of other agents. Our contributions are as follows:

- Simulation of Clinician Reasoning: We introduce a pioneering approach to explicitly simulate clinicians' diagnostic reasoning trajectories using clinical notes, teaching the model the diagnostic thinking of doctors.
- Hierarchical Agent Reflection: We develop an innovative hierarchical agent reflection framework, which enhances clinical note generation through structured iterative refinement. This framework significantly improves the accu-

racy and reliability of the generated data.

119

120

121

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

149

150

151

- Empirical Validation: Our experimental results demonstrate that models trained on datasets generated by our method significantly outperform both general-purpose large language models and domain-specific medical models in diagnostic tasks. Ablation studies further confirm the effectiveness of the hierarchical agent reflection mechanism.
- Enhanced Diagnostic Transparency: The model produces diagnostic pathways that are both interpretable and traceable, effectively aligning with clinicians' reasoning processes. This transparency enhances trust in AI-assisted diagnostics, making it a reliable tool for clinical applications.

2 Related Works

2.1 Medical Large Language Models

In recent years, the application of large language models in the medical field has become a major research focus (Singhal et al., 2023; Thirunavukarasu et al., 2023; Han et al., 2023; Kim et al.; Saab et al., 2024; Truhn et al., 2024; Christophe et al., 2024; Zhou et al., 2023). These models enhance LLM capabilities in medicine through various approaches. For instance, models like BioMedLM (Bolton et al.), OphGLM (Gao et al., 2023), and GatorTronGPT (Peng et al., 2023) absorb extensive medical knowledge during pre-training, enabling strong performance across a range of medical tasks. Given the time and cost associated with developing specialized medical LLMs from scratch, models like Med-

248

249

250

251

252

Gemini (Yang et al., 2024), Med42 (Christophe et al., 2024), MedAlpaca (Han et al., 2023), and MedPaLM-2 (Singhal et al., 2025) opt to build on robust general-purpose base models, fine-tuning them with different strategies to meet the specific needs of the medical domain and ultimately transforming them into specialized medical LLMs.

152

153

154

155

156

157

158

159

161

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

181

182

183

186

187

188

190

191

192

193

194

195

196

197

198

199

202

Furthermore, certain models have improved their medical capabilities by implementing preference alignment techniques. For example, HuatuoGPTo1 (Chen et al., 2024b) significantly enhanced its medical reasoning abilities by utilizing verifiable medical reasoning datasets and reinforcement learning. MedFound aligned itself with standard clinical practices by introducing a unified preference alignment framework, while Baichuan-M1 improved its diagnostic capabilities through reinforcement learning and pairwise data optimization.

While fine-tuning reduces computational resources compared to pre-training, it still requires additional model training and high-quality datasets, which can be resource-intensive. In contrast, prompt engineering offers a more efficient method to adapt base models to specific use cases without altering model parameters. Techniques like few-shot learning, in-context learning, chain-ofthought prompting (Wei et al., 2022), and retrievalaugmented generation (RAG) (Lewis et al., 2020) are commonly used. Given the critical importance of accuracy in medical applications, RAG is particularly effective for providing models with reliable information. Models such as Oncology-GPT-4(Ferber et al., 2024), MedRAG (Xiong et al., 2024), and MedGraphRAG (Wu et al., a) enhance overall performance by incorporating external, trustworthy sources of information into the answer generation process.

Our approach combines the strengths of these techniques with the capabilities of general models. We find relying solely on LLMs for direct questionanswering may not sufficiently meet the demands of medical diagnosis. Therefore, we leverage generated patient record data to focus on more complex medical diagnostics, enabling the model to handle intricate medical issues with greater accuracy.

2.2 Multi-Agent Collaboration

A growing body of research demonstrates that collaborative frameworks involving multiple LLM agents can effectively address the limitations of individual models when tackling complex tasks, resulting in more efficient and precise execution across domains such as finance, coding, literature, and mathematics (Li et al., 2023; Wu et al., b; Huot et al., 2024; Hong et al., 2023; Han et al., 2024; Zhang et al., 2023). In the medical field, which is closely tied to everyday life, multi-agent collaboration frameworks are increasingly being recognized for their potential. By leveraging collaboration between different LLM agents, tasks like diagnosis, treatment planning, and patient management can be more effectively handled.

For example, MEDAGENTS (Tang et al., 2023), the first multi-agent framework proposed in the medical domain, has demonstrated exceptional performance in extracting and utilizing medical expertise from LLMs while improving their reasoning capabilities. Agent Hospital (Li et al.), by creating a hospital simulation environment with evolving medical agents, has achieved ongoing improvements in clinician agent performance, both in simulated and real-world settings, thereby laying the groundwork for the use of LLM-driven agent technology in medical applications. Inspired by clinicians' decision-making processes, MDAgents (Kim et al., 2024) has developed an adaptive medical decision-making framework that uses LLMs to simulate hierarchical diagnostic procedures, ranging from individual clinicians to collaborative clinical teams. This has opened new possibilities for enhancing LLM-assisted medical diagnostic systems and advancing automated clinical reasoning.

Building on the success of multi-agent collaboration frameworks, we propose a dual-agent reflection and correction mechanism, augmented by knowledge-enhancement techniques, to further improve the accuracy of generated clinical notes.

3 Method

In this section, we provide a detailed explanation of our method for generating clinical notes using a hierarchical agent reflection mechanism, along with a knowledge enhancement strategy. First, we describe the process of constructing standardized templates for clinical notes. Following this, we outline the specific functions of each agent, as well as the reflection and correction mechanisms and the strategies for knowledge enhancement. Finally, we discuss the generation of high-quality clinical notes, which are used to train models, thereby improving their ability to effectively utilize these notes for medical diagnosis.



Figure 2: The overview of our proposed framework Hierarchical agent reflection. First, the patient inputs a query, the ICA initially collects information. Then, the PDA internally reflects and makes a preliminary diagnosis, followed by the DDA internally reflecting to conduct a differential diagnosis. The resulting raw clinical note is then reviewed by the upper-level CA, which reflects based on the standardized template. If there are any inconsistencies, the CA will notify the three lower-level agents for corrections, ultimately producing a revised, high-quality clinical note. Here, \emptyset refers to the knowledge base we collected as mentioned in Section 3.1, and \emptyset refers to the LLM.

3.1 Template for Clinical Note

254

256

260

261

270

For clinicians, creating high-quality clinical notes is essential for ensuring thorough patient care and accurate diagnoses. These notes provide critical documentation of patient symptoms, preliminary diagnoses, and the rationale behind the differential diagnosis. During the diagnostic process, clinicians typically start with the patient's description of symptoms, conducting further examinations and tests to gather comprehensive information that forms the patient's case profile. This information is meticulously documented in Clinical Notes, which clinicians use for comprehensive assessment, leading to preliminary diagnoses and corresponding rationale. Subsequently, clinicians apply their distinctive differential diagnostic reasoning to verify the preliminary diagnosis and ultimately confirm the disease.

To construct the dataset of clinical notes, we collected 433 diseases and their corresponding key consultation points. We used the state-of-the-art large model, Deepseek-R1, to construct the initial medical course record templates. The prompts for generating templates of clinical notes can be found in Appendix B. After generating the initial templates with R1, we engaged clinical clinicians to annotate and revise the medical course records for each disease, ultimately obtaining a standardized medical course record for each disease. 271

273

274

275

276

277

278

281

282

283

284

3.2 Hierarchical Agent Reflection

In our hierarchical agent reflection process, we configured four agents: the Information Collection Agent (ICA), the Preliminary Diagnosis Agent (PDA), the Differential Diagnosis Agent (DDA) and the Coordinator Agent(CA). We used Knowledge-enhanced methods to assist them in intra-agent and inter-agent reflection and correction. We provided these agents with a knowledge base, covering 433 diseases, collected in 3.1, along with a detailed map of diseases and their differential diagnoses. Specific information about the knowledge graph and the knowledge base can be found in Appendix C.1. The prompts used by each agent can be found in Appendix B. The pseudo code of Generating the Clinical Notes can be found at Appendix A.

Information Collection Agent. ICA is responsible for recording the patient's basic information, which includes their medical history, physical examination, and auxiliary tests. Subsequently, the ICA conducts a comprehensive analysis, synthesis, and organization of this information to document the characteristics of the case.

299

301

302

303

304

305

306Preliminary Diagnosis Agent. Based on the
case characteristics recorded by the ICA, the Pre-
liminary Diagnosis Agent provides an initial diag-
nosis and its diagnostic basis. Subsequently, the
PDA retrieves diagnostic key points related to the
initial diagnosis from the knowledge base and uses
these to reflect on the initial diagnostic process. It
then evaluates the accuracy of the initial diagno-
sis, and if deemed inaccurate, performs iterative
diagnostic corrections.

Differential Diagnosis Agent. The DDA first 316 retrieves a list of diseases requiring differential di-317 agnosis exclusion from the provided knowledge 318 graph, based on the initial diagnosis provided by 319 the PDA. It then acquires the diagnostic key points 320 for each of these diseases from the knowledge base and performs differential diagnosis for each disease using these points and the patient's case characteristics. Finally, the DDA reflects on the reasonable-324 ness of the entire differential diagnosis process; if found to be unreasonable, it conducts iterative 326 corrections.

Coordinator Agent. In hierarchical agent reflec-328 tion, the Coordinator Agent operates at a higher level than the ICA, PDA, and DDA. Specifically, 330 the CA first receives the raw clinical notes from the CA. It then uses the final diagnosis provided 332 by the DDA to match this note with standardized 334 clinical notes in the knowledge base, obtaining the standardized notes for the corresponding diseases. 335 The CA then reflects on and evaluates whether the outputs of the ICA, PDA, and DDA align with the standards by comparing the raw clinical note with 338

the matched standardized clinical notes. If significant discrepancies are found between an agent's output and the standardized clinical notes, the CA identifies potential errors in that agent's process and notifies it of the reasons for reflection. Conversely, if the raw clinical note is deemed reasonable, the CA integrates and outputs a verified complete clinical note. Throughout this process, the CA leverages knowledge augmentation and the hierarchical agent reflection mechanism to enhance the accuracy of the generated clinical notes.

339

340

341

344

345

346

348

349

350

351

352

353

356

357

358

359

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

385

386

389

The combination of self-reflection in the ICA, PDA, and DDA agents, along with supervisory feedback from the CA agent, enhances accuracy of the generated clinical notes. Self-reflection allows each agent to independently refine its reasoning and detect errors, while the CA agent provides additional oversight to ensure the final output aligns with clinical standards. This dual-layer feedback system improves error detection, enables better generalization across scenarios, and supports continuous adaptation, ultimately leading to more reliable and accurate clinical decision-making.

3.3 Enhance LLM Medical Diagnosis with Clinical Notes

Our Raw dataset is $\mathcal{D}_{Raw} = \{x_i, y_i\}_{i=1}^{|\mathcal{D}_{Raw}|}$, where x_i denotes a patient's question, and y_i denotes the original answer without diagnostic logic. After using our hierarchical agent reflection framework, the data format becomes: $\mathcal{D}_{\text{note}} = \{x_i, (y_{i1}, y_{i2}, y_{i3}) \to y_{i4} \to (y_{i5}, y_{i6}) \to (y_{i7}, y_{i8}, y_{i9})\}_{i=1}^{|\mathcal{D}_{\text{note}}|}$ where x_i denotes a patient's question, while y_{i1} through y_{i9} represent various components of the clinical note: y_{i1} is the medical history, y_{i2} the physical examination, y_{i3} auxiliary examination, y_{i4} clinical features, y_{i5} initial diagnosis, y_{i6} diagnostic basis, y_{i7} disease list, y_{i8} differential diagnosis process, and y_{i9} the final diagnosis. After extracting information from x_i , we obtain (y_{i1}, y_{i2}, y_{i3}) , which are then further organized and summarized to derive y_{i4} . Then generating the initial diagnosis and diagnostic basis (y_{i5}, y_{i6}) . Finally, the process results in a detailed differential diagnosis and the final diagnosis (y_{i7}, y_{i8}, y_{i9}) .

In the standard post-training paradigm, pretrained language models are typically optimized through supervised fine-tuning to better follow user instructions or adhere to specific formats. (Ouyang et al., 2022; Zhou et al., 2024; Fan et al., 2024). We utilize SFT to train the model to gradually generate clinical notes. This approach enables the

484

485

437

438

model to perform thoughtful reasoning, using prior knowledge to generate patient information collection, preliminary diagnosis, and differential diagnosis steps. We randomly sample its prefix (which can be empty), then supervise the model to reason before responding by optimizing the following objective:

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{(x,y)} \bigg[\sum_{t=1}^{T} \log p_{\theta} \big(y_t \mid x \oplus y_{< t} \big) \bigg].$$
(1)

4 Experiments

396

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

4.1 Experimental Setup

Training Data We construct a Chinese clinical note dataset containing 2K notes and an English clinical notes dataset containing 10K notes respectively from the training sets of RJUA-QA (Lyu et al., 2023) and MedQA (Jin et al., 2021) by applying our hierarchical agent reflection framework.

Model Training After obtaining the dataset of clinical notes generated by our framework, we trained the LLM using LLaMA-Factory (Zheng et al., 2024), a widely-used library for LLM training. We conducted all experiments on eight NVIDIA A100 (80G) GPUs. Specifically, we fine-tuned the model using LoRA (Hu et al., 2021) with the DeepSpeed (Rasley et al., 2020) library and Zero Redundancy Optimizer (ZeRO) (Rajbhandari et al., 2020) Stage 2. For SFT, we set the epoch to 3, the learning rate to 5e-5, and the context length to 4096.

Baselines We utilized the generated clinical note 418 data for finetuning of the model and compared the 419 results with two types of LLMs: 1) General LLMs: 420 the GPT series (Achiam et al., 2023), Qwen-2.5 421 (Team, 2024), LLaMA-3.1 (Dubey et al., 2024) 422 and Deepseek (Liu et al., 2024); and 2) Medical-423 Specific LLMs: HuatuoGPT (Chen et al., 2024b), 424 MedFound (Liu et al., 2025), and Baichuan-M1. 425

Benchmarks We evaluate on the standard med-426 ical diagnostic benchmarks: including the RJUA-427 QA(test set) (Lyu et al., 2023) and Urological Dis-428 ease Diagnosis Dataset(UDDD), both of which are 429 Chinese medical diagnosis datasets, using the F1 430 score to assess diagnostic accuracy. Additionally, 431 432 we assessed the Medbullets test set (Chen et al., 2024a), an English single-choice medical diagno-433 sis dataset featuring five-option and four-option 434 questions, using accuracy as the metric for diag-435 nostic performance. To enhance the reliability of 436

the experimental results, we ran every evaluation 3 times and averaged the results and variance.

4.2 Experimental Results

Main Results We evaluated various LLMs on medical benchmarks, as shown in Table 1. The results indicate that foundational models, which have not undergone enhanced training with specialized medical knowledge, perform rather poorly in the medical diagnosis domain. This is evident in models such as gwen2.5-7B and llama3.1-8B. Even when the parameter scale of these models is increased, the improvements in performance remain guite limited. For the currently popular reasoning models, such as gpt-40 and DeepSeek-R1, their performance on medical diagnosis datasets is still insufficient, which further indicates that having only reasoning ability cannot achieve ideal results in the medical field. In contrast, the Huatuo series models and the Baichuan-M1 model demonstrate more significant diagnostic capabilities in the field of medical diagnosis.

After being trained on our clinical note datasets, our models have demonstrated outstanding performance across all test datasets. In the comprehensive evaluation, all versions of our models exhibited improvements exceeding 5 points compared to their respective base models. Additionally, our Qwen series models have even surpassed other similar open-source large language models, including those of comparable scale to the Huatuo models. These results clearly validate the effectiveness of our approach. It is noteworthy that the Baichuan-M1 model, specifically developed for medical diagnostic capabilities, also performed exceptionally well, approaching the performance of our trained Qwen-72B model, despite our training data volume being significantly smaller than that used for Baichuan-M1.

5 Ablations on Training Data

Table 2 presents the results of ablation experiments on the qwen2.5-7B-Instruction model for diagnostic tasks after fine-tuning with different training corpora. Among them, (1) **w/o Data** means no training data is used, (2) **Raw Data** refers to the fine-tuning data consisting of the original RJAU-QA training set. (3) **w/o HAR** denotes the dataset generated without hierarchical agent reflection, involving only the ICA, PDA, and DDA without reflection. (4) **w/o CA** indicates the dataset gener-

	RJUA-QA	UDDD	Medbullets-5op	Medbullets-4op				
\sim 7-8B Large Language Models								
LLaMA3.1-8B-Instruct	$\textbf{28.33} \pm 0.16$	58.17 ± 0.54	$\textbf{30.63} \pm 0.50$	45.56 ± 3.56				
Qwen2.5-7B-Instruct	$\textbf{36.39} \pm 0.25$	65.84 ± 0.34	37.23 ± 0.19	42.97 ± 0.50				
DeepSeek-R1-Distill-Qwen-7B	34.67 ± 0.90	60.19 ± 0.10	26.73 ± 1.32	28.79 ± 2.39				
🕹 Huatuo-o1-Qwen-7B	43.71 ± 1.72	$\textbf{72.51} \pm 0.25$	50.33 ± 2.82	52.69 ± 2.53				
🕹 Huatuo-o1-LLaMA-8B	34.17 ± 1.13	$\textbf{60.19} \pm 0.10$	$\textbf{45.07} \pm 1.39$	$\textbf{52.49} \pm 1.60$				
🕹 MedFound-7B	16.48 ± 5.33	$\textbf{33.08} \pm 6.44$	35.19 ± 2.37	29.67 ± 3.21				
local Science &	29.64 ± 1.35	53.51 ± 1.87	18.07 ± 2.31	25.00 ± 5.77				
🕹 🗟 LLaMA3.1-8B-Instruct	$\underline{44.71} \pm 0.70$	74.16 ± 0.51	50.76 ± 1.14	$\underline{57.47} \pm 1.75$				
🕹 🗟 Qwen2.5-7B-Instruct	$\textbf{46.50} \pm \textbf{0.06}$	$\textbf{74.75} \pm \textbf{0.02}$	$\textbf{53.35} \pm \textbf{1.14}$	$\textbf{60.17} \pm \textbf{1.42}$				
> 10B Large Language Models								
GPT-3.5-Turbo	28.92 ± 0.56	55.21 ± 0.66	35.71 ± 0.57	42.75 ± 0.68				
GPT-4-turbo	31.14 ± 0.02	$\textbf{59.95} \pm 0.31$	58.23 ± 0.18	65.26 ± 0.97				
GPT-40	$\textbf{33.98} \pm 0.05$	65.32 ± 0.64	69.48 ± 0.56	$\underline{75.00 \pm 0.65}$				
Deepseek-V3	$\textbf{37.34} \pm 0.01$	66.58 ± 0.09	56.71 ± 0.49	61.69 ± 0.65				
Deepseek-R1	44.31 ± 0.65	64.56 ± 0.26	68.51 ± 7.21	63.85 ± 1.96				
🕹 HuatuoGPT2-13B	33.13 ± 0.79	59.70 ± 0.99	$\textbf{37.77} \pm 1.35$	37.23 ± 3.47				
🕹 Baichuan-M1-14B	50.01 ± 1.05	$\textbf{75.60} \pm 1.62$	55.52 ± 1.17	61.58 ± 0.49				
LLaMA3.1-70B-Instruct	$\overline{\textbf{35.54} \pm 0.67}$	$\overline{\textbf{66.65} \pm 1.05}$	57.58 ± 1.05	64.29 ± 0.33				
Qwen2.5-72B-Instruct	38.54 ± 0.65	66.08 ± 1.07	54.76 ± 0.49	62.88 ± 0.99				
🕹 Huatuo-o1-LLaMA-70B	$\textbf{38.11} \pm 0.96$	68.07 ± 0.17	68.83 ± 0.65	$\textbf{73.38} \pm 1.72$				
🕹 🗟 LLaMA3.1-70B-Instruct	$\textbf{44.93} \pm 0.42$	73.25 ± 0.39	$\underline{70.24 \pm 0.65}$	$\textbf{74.78} \pm 1.21$				
🕹 🗟 Qwen2.5-72B-Instruct	$\textbf{50.61} \pm \textbf{2.10}$	$\textbf{77.29} \pm \textbf{0.91}$	$\textbf{71.21} \pm 1.14$	$\textbf{76.52} \pm \textbf{0.65}$				

Table 1: Main Results on Medical Benchmarks. LLMs with B are specifically trained for the medical domain, and B indicates LLMs training for clinical note dataset. **bold** highlights the best scores, and <u>underlines</u> indicate the second-best.

Data	w/o Data	Raw Data	w/o HAR	w/o CA	Com. HAR
RJUA-QA	$\textbf{36.39} \pm 0.25$	41.24 ± 1.16	42.95 ± 0.57	44.41 ± 0.26	$\textbf{46.50} \pm \textbf{0.06}$
UDDD	$\textbf{65.84} \pm 0.34$	68.73 ± 2.38	69.33 ± 0.21	$\textbf{72.76} \pm 1.01$	$\textbf{74.75} \pm \textbf{0.02}$

Table 2: **Ablation results on training data.** w/o Data indicates not using any data, Raw Data refers to raw training data, w/o HAR refers to without hierarchical agent reflection, w/o CA refers to without CA, and Com. HAR indicates the use of our complete framework. We report the average diagnostic accuracy and variance over three runs of the model fine-tuned with different training data.

ated by removing the CA for upper-level reflection.
(5) Com. HAR represents the dataset generated through the complete hierarchical agent reflection framework.

486

487

488

489

490

491

492

493

494

495

496 497

498

499

500

501

The model exhibits poor diagnostic performance without fine-tuning on any medical data. However, after fine-tuning with the original format of the RJUA-QA training set, the model's diagnostic accuracy improves. Further enhancement is achieved by converting the training data into a clinical note format, enabling the model to mirror a doctor's diagnostic logic, which results in increased accuracy. Incorporating knowledge for self-reflection by the Agent improves the quality of the generated clinical note data, further boosting diagnostic accuracy. Ultimately, employing the complete hierarchical agent reflection framework, combined with standardized clinical note templates annotated by clinicians, enhances the quality of the clinical note data again, elevating the model's diagnostic capability to new heights. 502

503

504

506

507

508

509

510

511

512

513

514

515

516

6 Human Evaluation

To get a deeper understanding of the differences in diagnostic accuracy and transparency between the model fine-tuned using the hierarchical agent reflection framework-generated clinical note dataset, and other medical models and base models, we manually compared the diagnostic results of different models on the same medical issue. The diagnostic results are depicted in Fig. 3 We can observe that the base model is quite disorganized



Figure 3: Case study on RJUA-QA. We examined a patient case requiring the diagnosis of two diseases, with key symptoms highlighted for emphasis. Panel (a) displays the zero-shot diagnostic result from the base model, Qwen2.5-7B-Instruction. Panel (b) shows the output from Huatuo-o1-7B, with its reasoning process omitted for brevity. Panel (c) presents the diagnostic result from Baichuan-M1-14B. Panel (d) illustrates the diagnostic outcome from the Qwen2.5-7B-Instruction model after fine-tuning with our high-quality clinical note data. Sections marked in red indicate errors in the model's responses, while those in green highlight areas where the model accurately used key symptoms to diagnose diseases.

and fails to adequately utilize patient information. The HuatuoGPT-o1 model demonstrates some medical knowledge errors, such as not recognizing that Mirabegron is a β 3-adrenergic receptor agonist used for overactive bladder symptoms. The Baichuan-M1 model has difficulty in distinguishing similar urinary incontinence diseases. In contrast, our fine-tuned model delivers a clear diagnostic process that better aligns with the clinical reasoning of clinicians. This is achieved by using our hierarchical agent reflection framework during the generation of the clinical note training dataset, which injects the model with the correct inquiry points related to diseases, enabling it to recognize specialized medical knowledge. Furthermore, our model effectively employs differential diagnosis techniques to exclude similar diseases.

7 Conclusion and Future Work

In this paper, we propose a hierarchical agent reflection framework to generate high-quality clinical notes. By training LLMs with clinical notes that reflect the reasoning processes of clinicians, we aim to enhance the model's ability to engage in medical reasoning and improve diagnostic accuracy. The model's output not only mirrors the reasoning clinicians use in diagnosis but also assists them by offering a similar thought process during clinical decision-making. Experimental results demonstrate that simulating clinicians' use of clinical notes for diagnosis significantly boosts the model's diagnostic performance. Moving forward, we plan to extend the framework to cover rare diseases and refine the model's reasoning capabilities for even greater diagnostic accuracy.

534

535

536

537

538

539

540

542

543

544

545

546

547

549

645

646

647

648

649

650

651

652

653

654

655

656

602

603

604

Limitations

551

566

567

569

573

574

575

576

577 578

579

581

582

583

584

585

586

588

590 591

592

593

595

596

597

598

In this paper, we aim to develop a hierarchical agent reflection framework that narrows the gap 553 between model-based diagnostic processes and the 554 diagnostic logic used by clinicians by generating 555 high-quality clinical note data. Despite our best ef-556 557 forts, certain limitations remain. First, our current work is limited to text-based medical diagnoses, while the medical field often involves a wealth of multimodal information that aids in diagnosis. Second, when it comes to rare and complex diseases, 561 562 our framework lacks the capability to compose the discussion section of challenging cases. We plan to address these limitations in future work.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- E Bolton, A Venigalla, M Yasunaga, D Hall, B Xiong, T Lee, R Daneshjou, J Frankle, P Liang, M Carbin, et al. Biomedlm: A 2.7 b parameter language model trained on biomedical text, arxiv, 2024. *arXiv preprint arXiv:2403.18421*.
- Hanjie Chen, Zhouxiang Fang, Yash Singla, and Mark Dredze. 2024a. Benchmarking large language models on answering and explaining challenging medical questions. *arXiv preprint arXiv:2402.18060*.
- Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. 2024b. Huatuogpt-o1, towards medical complex reasoning with llms. *arXiv preprint arXiv:2412.18925*.
- Clément Christophe, Praveen K Kanithi, Tathagata Raha, Shadab Khan, and Marco AF Pimentel. 2024. Med42-v2: A suite of clinical llms. *arXiv preprint arXiv:2408.06142*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Run-Ze Fan, Xuefeng Li, Haoyang Zou, Junlong Li, Shwai He, Ethan Chern, Jiewen Hu, and Pengfei Liu. 2024. Reformatted alignment. *arXiv preprint arXiv:2402.12219*.
- Dyke Ferber, Isabella C Wiest, Georg Wölflein, Matthias P Ebert, Gernot Beutel, Jan-Niklas Eckardt, Daniel Truhn, Christoph Springfeld, Dirk Jäger, and Jakob Nikolas Kather. 2024. Gpt-4 for information

retrieval and comparison of medical oncology guidelines. *NEJM AI*, 1(6):AIcs2300235.

- Weihao Gao, Zhuo Deng, Zhiyuan Niu, Fuju Rong, Chucheng Chen, Zheng Gong, Wenze Zhang, Daimin Xiao, Fang Li, Zhenjie Cao, et al. 2023. Ophglm: Training an ophthalmology large language-andvision assistant based on instructions and dialogue. arXiv preprint arXiv:2306.12174.
- Ethan Goh, Robert Gallo, Jason Hom, Eric Strong, Yingjie Weng, Hannah Kerman, Joséphine A Cool, Zahir Kanjee, Andrew S Parsons, Neera Ahuja, et al. 2024. Large language model influence on diagnostic reasoning: a randomized clinical trial. *JAMA Network Open*, 7(10):e2440969–e2440969.
- Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bressem. 2023. Medalpaca–an open-source collection of medical conversational ai models and training data. arXiv preprint arXiv:2304.08247.
- Xuewen Han, Neng Wang, Shangkun Che, Hongyang Yang, Kunpeng Zhang, and Sean Xin Xu. 2024. Enhancing investment analysis: Optimizing ai-agent collaboration in financial research. In *Proceedings of the 5th ACM International Conference on AI in Finance*, pages 538–546.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Fantine Huot, Reinald Kim Amplayo, Jennimaria Palomaki, Alice Shoshana Jakobovits, Elizabeth Clark, and Mirella Lapata. 2024. Agents' room: Narrative generation through multi-step collaboration. *arXiv preprint arXiv:2410.02603*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. arXiv preprint arXiv:2412.16720.
- 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.
- Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won Park. 2024. Mdagents: An adaptive collaboration

713

of llms for medical decision-making. In *The Thirty*eighth Annual Conference on Neural Information Processing Systems.

657

658

670

671

672

673

674

675

676

677

678

679

686

693

701

702

703

704

706 707

711

- Yubin Kim, Xuhai Xu, Daniel McDuff, Cynthia Breazeal, and Hae Won Park. Health-llm: Large language models for health prediction via wearable sensor data. arxiv 2024. arXiv preprint arXiv:2401.06866.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for" mind" exploration of large language model society. Advances in Neural Information Processing Systems, 36:51991–52008.
- J Li, S Wang, M Zhang, W Li, Y Lai, X Kang, et al. Agent hospital: A simulacrum of hospital with evolvable medical agents [internet]. arxiv; 2024 [cited 2024 may 10].
- Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, et al. 2024.
 Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*.
- Xiaohong Liu, Hao Liu, Guoxing Yang, Zeyu Jiang, Shuguang Cui, Zhaoze Zhang, Huan Wang, Liyuan Tao, Yongchang Sun, Zhu Song, et al. 2025. A generalist medical language model for disease diagnosis assistance. *Nature Medicine*, pages 1–11.
- Shiwei Lyu, Chenfei Chi, Hongbo Cai, Lei Shi, Xiaoyan Yang, Lei Liu, Xiang Chen, Deng Zhao, Zhiqiang Zhang, Xianguo Lyu, et al. 2023. Rjua-qa: A comprehensive qa dataset for urology. *arXiv preprint arXiv:2312.09785*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Cheng Peng, Xi Yang, Aokun Chen, Kaleb E. Smith, Nima PourNejatian, Anthony B. Costa, Cheryl Martin, Mona G. Flores, Ying Zhang, Tanja Magoc, Gloria Lipori, Duane A. Mitchell, Naykky S. Ospina, Mustafa M. Ahmed, William R. Hogan, Elizabeth A. Shenkman, Yi Guo, Jiang Bian, and Yonghui Wu. 2023. A study of generative large language model for medical research and healthcare. *npj Digital Medicine*, 6(1).
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. Zero: Memory optimizations

toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1– 16. IEEE.

- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th* ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 3505–3506.
- Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, et al. 2024. Capabilities of gemini models in medicine. *arXiv preprint arXiv:2404.18416*.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. 2025. Toward expert-level medical question answering with large language models. *Nature Medicine*, pages 1–8.
- 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 reasoning. *arXiv preprint arXiv:2311.10537*.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930– 1940.
- Daniel Truhn, Jan-Niklas Eckardt, Dyke Ferber, and Jakob Nikolas Kather. 2024. Large language models and multimodal foundation models for precision oncology. *NPJ Precision Oncology*, 8(1):72.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- J Wu, J Zhu, and Y Qi. a. Medical graph rag: Towards safe medical large language model via graph retrievalaugmented generation. arxiv 2024. *arXiv preprint arXiv:2408.04187*.
- Q Wu, G Bansal, J Zhang, Y Wu, B Li, E Zhu, L Jiang, X Zhang, S Zhang, J Liu, et al. b. Autogen: Enabling next-gen llm applications via multi-agent conversation,(2023). *arXiv preprint arXiv:2308.08155*.

Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrievalaugmented generation for medicine. *arXiv preprint arXiv:2402.13178*.

768

769

770

772 773

774

775

777

781

782

783 784

785

786

787

788

789

790

791

794

797

- Lin Yang, Shawn Xu, Andrew Sellergren, Timo Kohlberger, Yuchen Zhou, Ira Ktena, Atilla Kiraly, Faruk Ahmed, Farhad Hormozdiari, Tiam Jaroensri, et al. 2024. Advancing multimodal medical capabilities of gemini. *arXiv preprint arXiv:2405.03162*.
- Jintian Zhang, Xin Xu, Ningyu Zhang, Ruibo Liu, Bryan Hooi, and Shumin Deng. 2023. Exploring collaboration mechanisms for llm agents: A social psychology view. *arXiv preprint arXiv:2310.02124*.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations), Bangkok, Thailand. Association for Computational Linguistics.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.
- Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu Zou, Jinfa Huang, Jinge Wu, Yiru Li, Sam S Chen, Peilin Zhou, Junling Liu, et al. 2023. A survey of large language models in medicine: Progress, application, and challenge. arXiv preprint arXiv:2311.05112.

A Algorithm

799

In this algorithm, the MH_i , PE_i and AE_i respectively denote Medical History, Physical Examination, and Auxiliary Examination, while CF stands for Clinical Features. The D and B represent Preliminary Diagnosis and Diagnostic Basis, respectively, and K refers to key inquiry points associated with a disease, sourced from the knowledge base. The Dia_{flag} and Dia_{error} are used to indicate the correctness of the preliminary diagnosis reflection and areas for improvement in subsequent iterations should errors occur. The DisList represents the list of diseases that require differential diagnosis to be excluded based on the preliminary diagnosis results. The process of differential diagnosis is denoted by Diff, while Diff_{flag} and Diff_{error} indicate the correctness of the differential diagnosis reflection and the necessary improvements for future iterations in case of errors. Finally, T_{flag} and T_{error} signify the correctness evaluation after CA reflection and highlight the aspects that need to be communicated to other agents for improvement if errors are detected.

Algorithm 1 Generate the Clinical Notes

- 1: Input: Question Q, Knowledge Graph KG, Clinical Note Template T, Disease Knowledge DK {Initialization}
- 2: Initialize Information Collection Agent ICA, Preliminary Diagnosis Agent PDA, Differential Diagnosis Agent DDA and Coordinator Agent CA, Maximum Attempts N {Generate the Progress Notes}
- 3: for try count i = 0 to N 1 and T_{flag} is False do
- 4: $MH_i, PE_i, AE_i \leftarrow \text{Extraction of ICA}(Q)$
- 5: $CF_i \leftarrow \text{Summarization of ICA}(MH_i, PE_i, AE_i)$
- 6: $\text{Dia}_{\text{error}} \leftarrow \text{None}, \text{Diff}_{\text{error}} \leftarrow \text{None}$
- 7: for Dia_{flag} is *False* and j = 0 to N 1 do
- 8: $\mathcal{D}_j, \mathcal{B}_j \leftarrow \text{Diagnosis of PDA}(CC_i, \text{Dia}_{error})$
- 9: $\mathcal{K}_j \leftarrow \text{Retrieve of PDA}(\mathcal{D}_j, \text{DK})$
- 10: $\text{Dia}_{\text{flag}}, \text{Dia}_{\text{error}} \leftarrow \text{Reflection of } \text{PDA}(\mathcal{D}_j, \mathcal{B}_j, \mathcal{K}_j)$
- 11: **end for**

```
12: for Diff<sub>flag</sub> is False and j = 0 to N - 1 do
```

- 13: DisList_j \leftarrow Differential Diagnosis List of DDA(\mathcal{D}_j , KG)
- 14: $\mathcal{K}_j \leftarrow \text{Retrieve of DDA}(\text{DisList}_j, \text{DK})$
- 15: $\text{Diff}_j \leftarrow \text{Differential Process of PDA}(\text{DisList}_j, \mathcal{K}_j, \text{Diff}_{\text{error}})$
- 16: $\operatorname{Diff}_{\operatorname{flag}}, \operatorname{Diff}_{\operatorname{error}} \leftarrow \operatorname{Reflection of DDA}(\operatorname{DisList}_j, \operatorname{Diff}_j, \mathcal{K}_j)$
- 17: **end for**
- 18: RawNote_i \leftarrow Output Raw Clinical Note(ICA, PDA, DDA)
- 19: $\mathcal{T}_i \leftarrow \text{Retrieval Standardized Clinical Note Template of CA}(\mathcal{D}_i, T)$
- 20: $T_{\text{flag}}, T_{\text{error}} \leftarrow \text{Reflection of CA}(\text{ICA}, \mathcal{T}_i)$
- 21: **if** T_{flag} is *False* then
- 22: ICA \leftarrow Corrective of ICA(T_{error})
- 23: **end if**
- 24: $T_{\text{flag}}, T_{\text{error}} \leftarrow \text{Reflection of CA}(\text{PDA}, \mathcal{T}_i)$
- 25: **if** T_{flag} is *False* then
- 26: $PDA \leftarrow Corrective of PDA(T_{error})$
- 27: **end if**
- 28: $T_{\text{flag}}, T_{\text{error}} \leftarrow \text{Reflection of CA}(\text{DDA}, \mathcal{T}_i)$
- 29: **if** T_{flag} is *False* then
- 30: DDA \leftarrow Corrective of DDA(T_{error})
- 31: **end if**
- 32: **end for**
- 33: **Return** Revised Clinical Note (ICA, PDA, DDA)

B Prompt Templates

B.1 Generate Raw Clinical Note

Generate Raw Clinical Note Prompt

You are an experienced medical expert skilled in drafting standardized medical course records based on diseases and key consultation points. Please use the provided disease information and corresponding consultation points, along with the given template and supplied knowledge, to compose a standardized medical course record for this disease.

Below is the knowledge to this disease:

{{disease}}

{{Diagnostic key points}}

Below is the template for the clinical note:

Medical history:\n\n Physical examination:\n\n Auxiliary examination:\n\n Case characteristics:\n\n Initial diagnosis:\n\n Diagnostic basis:\n\n Diseases List:\n\n Differential diagnosis process:\n\n Final diagnosis:

B.2 Information Collection Agent Setting

Patient Information Extrction Prompt

You are an experienced clinical note specialist, adept at extracting the medical history, physical examination, and auxiliary examination information from data provided by patient. Please use the information provided by the patient to systematically consider and itemize the medical history, physical examination, and auxiliary examinations. If certain data are not provided, mark the corresponding section as 'None' without making additional assumptions.

Below is the patient's question:

{{question}}

Analysis and Summarize Prompt

You are an experienced medical analysis expert, skilled in comprehensively analyzing, summarizing, and organizing a patient's medical history, physical examination, and auxiliary examination to document the patient's clinical features. Please carefully review the patient's issues and itemize the clinical features, including positive findings and negative symptoms and signs relevant for differential diagnosis. Be sure to use only the provided information, without referencing external data.

Below is the medical history, physical examination, and auxiliary examination to this patient:

{{Medical history}}

{{Physical examination}}

{{Auxiliary examination}}

Below is the patient's question: {{question}}

B.3 Preliminary Diagnosis Agent Setting

Make Preliminary Diagnosis Prompt

You are an experienced clinical diagnosis expert, skilled in making preliminary diagnoses and analyses based on provided patient clinical features. Please provide a preliminary diagnosis based on the patient's case features and detail the diagnostic basis point by point.

Below is the clinical features to this patient:

{{Clinical features}}

Below is the patient's question:

{{question}}

812

813

814

Reflect Preliminary Diagnosis Prompt

You are an experienced clinical review expert, skilled in evaluating the diagnostic validity of clinical notes based on key inquiry points for diseases. Please thoroughly review the key inquiry points of the preliminary diagnosis provided and assess whether the preliminary diagnosis and diagnostic basis in the clinical note align with these points. If deemed unreasonable, output the result as a JSON-formatted Dict{"flag": false, "diagnosis_error": Str(Reasons for diagnostic errors)}.

Below is the preliminary diagnosis and diagnostic basis: {{Preliminary Diagnosis}} {{Diagnostic Basis}} Below is the key inquiry points:

{{key inquiry points}}

B.4 Differential Diagnosis Agent Setting

Differential Diagnosis Prompt

You are an experienced differential diagnosis expert, skilled in systematically analyzing key inquiry points to rule out diseases. Please carefully review the inquiry points of the diseases requiring differentiation and conduct a step-by-step differential diagnosis based on the patient's clinical note.

Document the differential diagnosis process point by point and output it in JSON format as Dict{"diff_process": Str(differential diagnosis process)}.

Below is the list of diseases to be ruled out through differential diagnosis:

{{Diseases List}}

Below is the key inquiry points to these diseases:

{{key inquiry points}}

Reflect Differential Diagnosis Process Prompt

You are an experienced clinical differential diagnosis expert, skilled in reflecting on and evaluating the rationality of differential diagnosis processes. Please reflect on the differential diagnosis process and assess whether the differentiation for each disease is reasonable.

If it is reasonable, output in JSON format as Dict{"flag":true, "Final_Diagnosis": Str(final diagnosis)}. Otherwise, output in JSON format as Dict{"flag":false, "diff_error": Str(Diseases requiring rediagnosis)}.

Below is the list of diseases to be ruled out through differential diagnosis, along with the corresponding diagnostic process.

{{Diseases List}}
{{Differential Diagnosis Process}}

822

B.5 Coordinator Agent Setting

Reflect and Correct ICA Prompt

You are an experienced expert in reviewing clinical notes, skilled in comparing raw clinical note with a given standardized template. Now, please compare the obtained raw clinical note with the given standardized clinical note template. The part that needs to be analyzed is the medical history, physical examination, auxiliary examination, and clinical features. If you find any part to be unreasonable, provide suggestions for improvement, and output in JSON format as Dict{"flag":false, "ICA_error": Str(suggestions for improvement)}.

Below is the raw clinical note.

{{Raw Clinical Note}}

Below is a standardized template for a standardized clinical note of the final diagnosis. {{Standardized Clinical Note}}

821

Reflect and Correct PDA Prompt

You are an experienced expert in reviewing clinical notes, skilled in comparing raw clinical note with a given standardized template. Now, please compare the obtained raw clinical note with the given standardized clinical note template. The part that needs to be analyzed is the preliminary diagnosis and diagnostic basis. If you think this part is unreasonable, please give suggestions for improvement. , and output in JSON format as Dict{"flag":false, "PDA_error": Str(suggestions for improvement)}.

Below is the raw clinical note.

{{Raw Clinical Note}}

Below is a standardized template for a standardized clinical note of the final diagnosis. {{Standardized Clinical Note}}

Reflect and Correct DDA Prompt

You are an experienced expert in reviewing clinical notes, skilled in comparing raw clinical note with a given standardized template. Now, please compare the obtained raw clinical note with the given standardized clinical note template. The part that needs to be analyzed is the diseases list and differential diagnosis process. If you think this part is unreasonable, please give suggestions for improvement. , and output in JSON format as Dict{"flag":false, "DDA_error": Str(suggestions for improvement)}.

Below is the raw clinical note.

{{Raw Clinical Note}}

Below is a standardized template for a standardized clinical note of the final diagnosis.

{{Standardized Clinical Note}}

C Knowledge Base and Knowledge Graph

C.1 Knowledge Graph



Figure 4: A knowledge graph of diseases and those requiring differential diagnosis, with \bigcirc refers to the diseases used in this paper.

825

826

C.2 Knowledge Base

Standardized Clinical Note

Medical history : 1. The 68-year-old patient experiences urinary leakage when coughing, sneezing, or during urgency. 2. She used Mirabegron for one month, with symptom improvement during treatment, but symptoms recurred within two days after she stopped the medication. 3. She underwent coronary intervention two months ago. **Physical examination :** None

Auxiliary examination : 1. In urinalysis, the microscopic white blood cell count was 27.7/HPF two months ago and 2.1/HPF in the most recent analysis.

Clinical features : 1. The patient is a 68-year-old female who has recently experienced frequent nighttime urination and incontinence, with normal urination frequency during the day but requiring three trips at night. 2. She experiences urinary leakage when coughing, sneezing, and during urgency. 3. She used Mirabegron for one month, which improved symptoms, but they recurred after discontinuation. 4. She underwent coronary intervention two months ago. 5. Urinalysis showed a high white blood cell count of 27.7/HPF two months ago, which has since decreased to normal levels at 2.1/HPF in the most recent analysis.

Initial diagnosis : stress incontinence, overactive bladder

Diagnostic basis : 1. The patient experiences urinary leakage during coughing and sneezing, which is indicative of typical stress urinary incontinence. 2. The patient exhibits urgency and increased nighttime urination, consistent with overactive bladder, but lacks other symptoms such as frequency and dysuria. The effectiveness of Mirabegron, a medication primarily used for overactive bladder, further supports this diagnosis.

Diseases List : urge incontinence, overflow incontinence, lower urinary tract syndrome

Differential diagnosis process : 1. The patient experiences urinary leakage during urgency without symptoms like frequency or dysuria, and shows improvement with Mirabegron, allowing us to preliminarily rule out urge incontinence. 2. Overflow incontinence is often caused by lower urinary tract obstruction, such as prostatic hyperplasia. This patient has no relevant history, and the white blood cell count in the urinalysis has returned to normal, largely excluding this possibility. 3. Lower urinary tract syndrome encompasses various symptoms like frequency, urgency, and dysuria. The patient only exhibits urgency and leakage, and responds well to Mirabegron, which does not strongly align with the characteristics of lower urinary tract syndrome.

Final diagnosis : stress incontinence, overactive bladder