

MedReflect: Teaching Medical LLMs to Self-Improve via Reflective Correction

Anonymous ACL submission

Abstract

Medical problem-solving demands expert knowledge and intricate reasoning. Recent studies of large language models (LLMs) attempt to ease this complexity by introducing external knowledge verification through retrieval-augmented generation or by training on reasoning datasets. However, these approaches suffer from drawbacks such as retrieval overhead and high annotation costs, and they heavily rely on substituted external assistants to reach limited performance in medical field. In this paper, we introduce MedReflect, a generalizable framework designed to inspire LLMs with a physician-like reflective thinking mode. MedReflect generates a single-pass reflection chain that includes initial hypothesis generation, self-questioning, self-answering and decision refinement. This self-verified and self-reflective nature releases large language model’s latent capability in medical problem-solving without external retrieval or heavy annotation. We demonstrate that MedReflect enables cost-efficient medical dataset construction. With only a minimal subset of randomly sampled training examples and lightweight fine-tuning, this approach achieves notable absolute accuracy improvements across a series of medical benchmarks while significantly cutting annotation requirements. Our results provide evidence that LLMs can learn to solve specialized medical problems via self-reflection and self-improvement, reducing reliance on external supervision and extensive task-specific fine-tuning data.

1 Introduction

Recent progress of large language models (LLMs) has showcased their immense potential in medical tasks (Lin and Kuo, 2025; Xu et al., 2025). Despite this progress, deploying LLMs in professional medical scenarios offers unique challenges compared to general domains. Medical decision-making is generally narrower in scope and involves complex

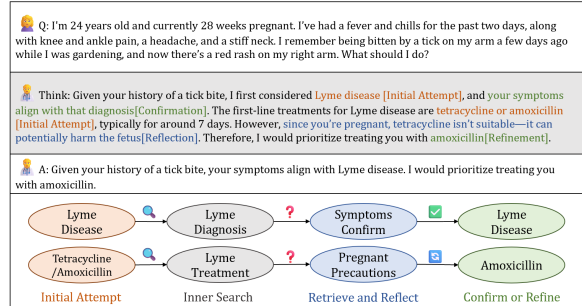


Figure 1: Example of the Physician Clinical Reasoning and Self-Correction Process Simulated by MedReflect. This figure demonstrates, through a specific clinical case, how the MedReflect framework simulates the human physician’s cognitive loop of "Hypothesis-Verification-Reflection-Refinement".

scenarios, which necessitate meticulous thinking to ensure more reliable answers. In addition, specialized medical terminology and intricate clinical narratives may increase hallucinated or unreliable outputs (Dou et al., 2024).

Existing efforts primarily improve medical LLMs by injecting external medical knowledge. Retrieval-Augmented Generation method (RAG), for instance, consults external sources to verify and refine the generated answers (Wu et al., 2024; Lu et al., 2025). However, this method heavily relies on additional external knowledge and incurs extra storage and retrieval overhead.

Recently, reasoning techniques have been explored to guide LLMs through predefined, structured sequence of analytical steps during generation (Kwon et al., 2024). For instance, MedReason (Wu et al., 2025) leverages structured knowledge graphs to synthesize reasoning chains and construct Chain-of-Thought (CoT) data. However, these methods rely on carefully curated datasets to build reasoning paths and require domain experts to predefine comprehensive medical reasoning processes, substantially increasing annotation

068 costs. To alleviate this burden, HuatuoGPT-o1
069 (Chen et al., 2024) use an auxiliary LLM verifier to
070 assess model-generated answers and guide the cor-
071 rection of errors through external feedback signals,
072 highlighting the potential of language models to
073 self-correct under the guidance of external signals.

074 However, these approaches substitute external
075 mechanisms for the model’s internal capabilities
076 of knowledge localization and reasoning planning.
077 This leads to a fundamental question: Can med-
078 ical language models learn to generate hypothe-
079 ses, retrieve relevant medical knowledge, and per-
080 form self-verification and correction—all within
081 a single generation process? The motivation for
082 raising this question is that LLMs already possess
083 extensive medical knowledge during pretraining
084 (Vladika et al., 2024). As discussed above, when
085 faced with complex medical problems, LLMs can
086 be guided by external verifiers to rediscover the
087 correct reasoning path. This resembles clinical
088 reasoning in practice: careful deliberation that in-
089 volves formulating an initial hypothesis, retrieving
090 relevant medical knowledge, and then validating
091 or revising conclusions through iterative reasoning,
092 which reflects the meticulous thinking required in
093 medical decision-making. An illustrative example
094 is shown in Figure 1.

095 In this work, we systematize the key compo-
096 nents of the medical decision-making process and
097 propose MedReflect. Our key idea is to equip the
098 model with a directional reflection-and-correction
099 process. We design a structured reflection mech-
100 anism that guides the model to generate and an-
101 swer its own questions during reasoning, which
102 injects direction into reflection. Through this self-
103 questioning–self-answering procedure, the model
104 explicitly probes uncertain assumptions, retrieves
105 relevant medical knowledge on its own, and re-
106 vises intermediate or final decisions accordingly.
107 As a result, MedReflect does not rely on external
108 knowledge retrieval methods, and instead, it uses
109 self-generated questions and answers to perform
110 internal retrieval and correction within a single gen-
111 eration. Experiments demonstrate that our method
112 yields improvement on medical benchmarks. Our
113 contributions are as follows:

- 114 • We propose **MedReflect**, a reflection-and-
115 correction mechanism to enhance the depth
116 and quality of reasoning in medical LLMs.
- 117 • We develop a practical approach that lever-
118 ages LLMs to construct a low-cost medical

reflection dataset. 119

- Experiments on multiple medical QA bench- 120
marks show that MedReflect consistently im- 121
proves accuracy. The results demonstrate 122
that reflective supervision effectively teaches 123
models to self-reflect and self-correct during 124
generation, outperforming existing chain-of- 125
thought training methods in both performance 126
and training efficiency. 127

2 Related Works 128

Medical LLMs. Prior work has extensively ex- 129
plored developing medical-specific LLMs to excel 130
in the medical domain (Chen et al., 2024, 2023; 131
Team et al., 2025; Labrak et al., 2024; Zhang 132
et al., 2024). While promising, applying LLMs to 133
complex medical cases remains challenging (Ríos- 134
Hoyo et al., 2024), with persistent concerns regard- 135
ing ethics and hallucinations (Soffer et al., 2025). 136
To address diagnostic capabilities, Liu et al. (2025) 137
proposes a generalist medical LLM for diagnostic 138
reasoning across specialties. Several approaches 139
leverage external data or specialized training to 140
bridge these gaps: UltraMedical (Zhang et al., 141
2024) comprises 410K instruction-following ex- 142
amples, while BioMistral (Labrak et al., 2024) uti- 143
lizes PubMed Central for continued pre-training. 144
Jeong et al. (2024) further integrates retrieval with 145
self-reflection to enhance reliability in biomedical 146
tasks. Recent research has increasingly focused on 147
structured and factual reasoning. MedReason (Wu 148
et al., 2025) leverages knowledge graphs to elicit 149
factual reasoning steps, and MedFact-R1 (Li et al., 150
2025a) employs pseudo-label augmentation to bol- 151
ster factual medical reasoning. Similarly, Chen 152
et al. (2025) investigates improving syndrome dif- 153
ferentiation thinking. To enable self-improvement, 154
MuSeR (Zhou et al., 2025) targets medical context- 155
awareness via multifaceted self-refinement, while 156
HuatuoGPT-o1 (Chen et al., 2024) advances medi- 157
cal reasoning through a two-stage approach com- 158
bining SFT, verifiable medical problems, and rein- 159
forcement learning. 160

**Developing Models for Self-reflection and Self- 161
correction.** Self-reflection has emerged as a crit- 162
ical mechanism for mitigating hallucinations (Ji 163
et al., 2023) and improving reasoning, although its 164
effectiveness is contingent on specific conditions re- 165
garding model capability and task difficulty (Kamoi 166
et al., 2024). Several works have studied backtrack- 167
ing and search as forms of self-correction (Ye et al., 168

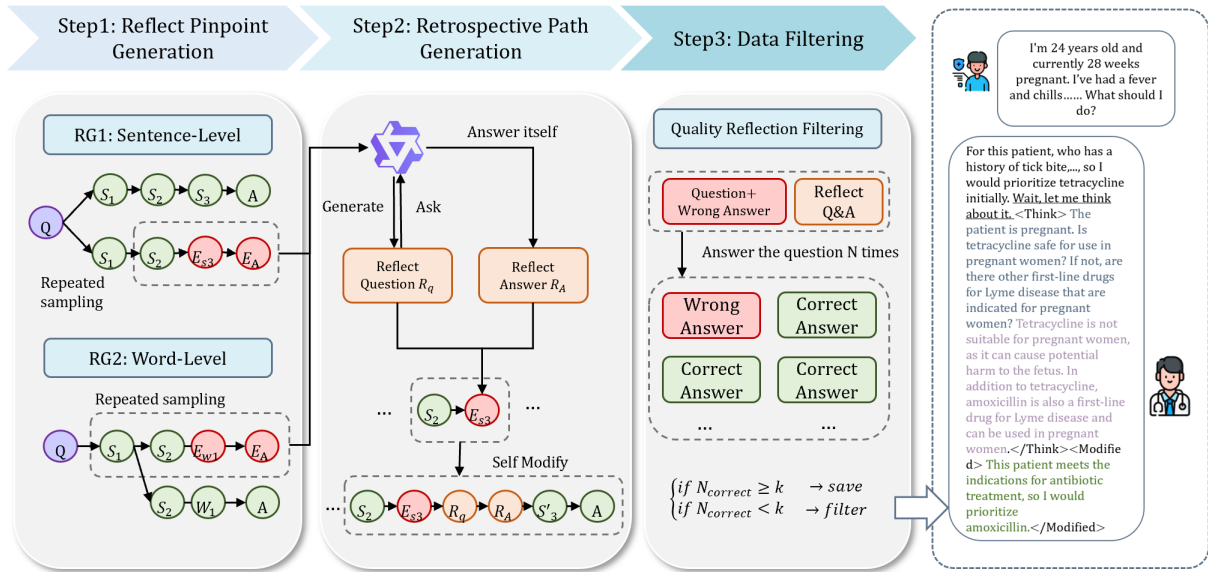


Figure 2: Data Construction Framework

2024; Qin et al., 2025). Notably, the Stream of Search (SoS) framework (Gandhi et al., 2024) enables models to self-correct by searching within language without external components, a phenomenon also explored in recent work on the emergence of thinking processes in LLMs (Ye et al., 2025). Beyond search, critique-based methods have gained traction. Xi et al. (2024) trains an expert critique model using step-level feedback to supervise the reasoning process. Building on this, ReflectEvo (Li et al., 2025b) explores iterative self-reflection to improve meta-introspection, while SaySelf (Xu et al., 2024) teaches LLMs to express confidence calibration within their self-reflective rationales. Additionally, Qu et al. (2024) introduces iterative fine-tuning to alter responses after unsuccessful attempts. More recently, Zhu et al. (2025) has begun to probe the underlying mechanisms for controlling and modulating these self-reflection behaviors.

3 Problem Setup and Preliminaries

We aim to explore whether medical LLMs can generate hypotheses, retrieve, perform self-verify and correction within a single generation process. To achieve this, we introduce MedReflect, focused on teaching LLMs to learn the reflect mode instead of directly injecting knowledge or improving their reasoning skills. The overview of the data generation of MedReflect is shown in Figure 2.

Given a medical question Q and its corresponding response trajectory $T = [S_1, S_2, \dots, S_n, A]$, the LLM is prompted to regenerate either a se-

lected step S_i or the entire answer. This regeneration may introduce errors, resulting in an erroneous trajectory or answer to the whole question $T_{er} = [S_1, S_2, \dots, E_i, \dots, E_A]$. Guided by the error E_i and the erroneous trajectory T_{er} , LLM then generates a targeted reflection question R_{q_i} along with its answer R_{a_i} . Finally, the original question Q , the erroneous trajectory T_{er} , and the reflection pair (R_{q_i}, R_{a_i}) are provided to LLM to facilitate error correction. Upon successful correction, the original correct step S_i and the reflection pair are incorporated into the trajectory, yielding the final reflective trajectory $T_{reflect} = [S_1, S_2, \dots, E_i, R_{q_i}, R_{a_i}, S_i, \dots, A]$.

4 Methodology

4.1 Datasource and Construction Model

We utilize two publicly available medical datasets: ChatDoctor (Li et al., 2023) and training split of MedMCQA (Pal et al., 2022). ChatDoctor contains 100k real-world conversations between patients and doctors sourced from HealthCareMagic.com. MedMCQA is a comprehensive multiple-choice question-answering dataset created specifically from real medical entrance exam questions. We used Qwen2.5-32B-Instruct to complete the entire data construction process.

4.2 Reflect Pinpoint Generation

To better control the diversity of generated pinpoints, we designed multiple reflective pinpoint generation pathways(RG), including reflective pro-

| | Sentence-Level(MedMCQA) | Word-Level(ChatDoctor) |
|--|--|--|
| Question | A 60 yr old chronic smoker presents with painless gross hematuria of 1 day duration. Which is the investigation? A: USG, B: X-ray KUB, C: Urine routine, D: Urine microscopy for malignant cytology cells | I'm experiencing a throbbing ache and severe fatigue behind my right knee, Any ideas what might be causing this? |
| Generate Details (👉 : pinpoint we insert) | LLM Re-Answer: We need to consider the most appropriate and comprehensive diagnostic tool. Ultrasonography (USG) can provide an overview of the urinary system, including the kidneys, bladder, and ureters. My answer is OptionA[truth: Option D] | Mask Medical Entities : Throbbing knee pain can be caused by several issues, such as a ligament tear from a sudden twist [entity1: etiology] To determine the exact cause, you should get MRI [entity2: medical examination and testing] and consult an professional. |
| | Compare and Find Pinpoint: We need to consider the most appropriate and comprehensive diagnostic tool. Ultrasonography (USG) can provide an overview of the urinary system, including the kidneys, bladder, and ureters. 🤔..... | Retry and Find Pinpoint: Throbbing knee pain can be caused by several issues, such as a ligament tear from a [meniscus injury X] 🤔 To determine the exact cause, you should get [MRI scan ✓] and consult an professional. |

Figure 3: Examples of Pinpoint Generation

cesses functioning at two different levels of detail. Figure 3 shows examples of two different levels.

RG1: Sentence Level. This approach primarily focuses on the overall rewriting of reasoning sentences generated by the LLM. We utilize this method to construct training data derived from the multi-choice MedMCQA dataset. Specifically, we perform repeated sampling by prompting the LLM to respond to the original question Q multiple times, generating new answers that include their reasoning processes. For each sampled response, we employ heuristic pattern matching to extract the multi-choice decision statement and compare it with the ground truth. If the extracted decision answer A' is incorrect, we locate the specific sentence within the response that led to the wrong option, designating it as a reflection pinpoint E_{si} .

RG2: Word Level This approach focuses on medical texts rich in entities. We used this method to construct the data of the consultation dataset ChatDoctor. We extract entity W_i from medical consultation reasoning sentence S_i and mask it with its corresponding entity type, such as diseases, etiology, and treatment. The LLM is prompted to predict and fill in the masked entity multiple times. If an attempt by the LLM is assessed as incorrect and such errors occur frequently, the sentence containing this entity serves as our reflection pinpoint E_{wi} . In contrast, responses that are assessed as correct will not be treated as such pinpoints. At this level of detail, a single physician’s response can produce one to three pinpoints for reflection.

4.2.1 Retrospective Path Generation

Reflect QA Generation To create reflection questions R_q and their matching answers R_a , we provide LLM with its original answer and emphasize that its answer was wrong. Through targeted prompting, we guide LLM to formulate a focused reflective question R_q . The goal of this question is to help LLM identify the key knowledge components needed to develop a better solution.

Afterward, the LLM is then instructed to answer the reflection question using only its own knowledge, without referring to the original question or adding extra details. The purpose of this method is to engage and utilize LLM’s internal knowledge, promoting thorough self-awareness and reflection.

Modification Based on Reflection Using Q , R_q and R_a as input, we guide the LLM to generate the modified statement M .

(1) For RG1, LLM regenerates sentences S'_i to replace original sentences with expression flaws, performing adaptive optimization to resolve any contextual coherence issues caused by the new sentences.

(2) For RG2, LLM tries to use the new response word W'_i to substitute inaccurate vocabulary.

These modified parts will be inspected in the next step to ensure they are qualified.

4.2.2 Data Filtering

To guarantee data quality, we first filtered out samples with insufficient reasoning or irrelevant content. We then implemented a secondary validation loop to mitigate confirmation bias by feeding the generated reflections (R_q , R_a) and the initial erroneous trajectory back into the model. A reflection

was deemed valid only if it guided the model to the objective ground truth—specifically, selecting the correct option for RG1 or restoring the masked entity for RG2. Finally, to exclude stochastic reasoning, we enforced a robustness constraint based on repeated trials. We conducted $[k]$ independent inference trials for each validated sample and calculated the success rate. Only instances exceeding a consistency threshold of $\tau = 0.8$ were retained. This yielded a final dataset of 36,413 medical consultation records and 21,107 multiple-choice questions.

4.3 Training Strategy

Following the construction of the generated reflection dataset, we investigate strategies to leverage these samples to enhance the model’s reasoning capabilities. We formalize the medical reasoning process as a sequential decision-making problem within a discrete semantic space. In conventional foundation models, the policy is often dominated by next-token prediction objectives derived from pre-training corpus statistics. This approach frequently results in logical disconnects or hallucinations when addressing complex medical cases. As illustrated in previous works, merely extending the Chain of Thought (CoT) is analogous to performing a directionless random walk on a semantic manifold; it induces divergent thinking but relies on stochastic retries to identify a superior solution.

To address this limitation, the **MedReflect** framework introduces a reflection operator designed to provide directed self-correction to the otherwise undirected reasoning process. This is formally defined as:

$$\tau_{reflect} = \mathcal{R}(\tau_{err}, \mathcal{K}_{int}) \Rightarrow (R_q, R_a) \quad (1)$$

Here, τ_{err} represents the initial trajectory containing error breakpoints, and \mathcal{K}_{int} denotes the model’s implicit medical knowledge. The reflection pair (R_q, R_a) serves as a critical control signal to rectify τ_{err} . Specifically, R_q facilitates *error attribution* by compelling the model to pinpoint knowledge deficits, while R_a leverages \mathcal{K}_{int} to generate a rectified factual basis.

To instantiate this theoretical operator within the language model, we structure the reasoning process as a token-augmented trajectory. We introduce four special tokens—`<Think>`, `</Think>`, `<Modified>`, and `</Modified>`—to explicitly delineate the semantic boundaries of the reflection process. Consequently, the training instance is transformed from a

standard input-output pair into a *reflective sequence* y_{seq} , constructed as:

$$y_{seq} = [\tau_{err}, \langle \text{Think} \rangle, R_q, R_a, \langle / \text{Think} \rangle, \langle \text{Modified} \rangle, \tau_{corrected}, \langle / \text{Modified} \rangle] \quad (2)$$

This structure forces the model to treat reflection not as an external module, but as an intrinsic part of its generative policy.

Finally, we align the reflective medical LLM with this structured reasoning pattern through supervised fine-tuning (SFT). We interpret SFT here as behavioral cloning of the reflection operator. Given the dataset $\mathcal{D}_{med} = \{(\mathbf{x}, y_{seq})\}$, the optimization objective is to maximize the likelihood of the augmented reflective sequence. The loss function for the model π_θ is defined as:

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{(\mathbf{x}, y_{seq}) \sim \mathcal{D}_{med}} \left[\sum_{t=1}^T \log \pi_\theta(y_t | \mathbf{x}, y_{<t}) \right] \quad (3)$$

By minimizing this loss, the model learns to sequentially generate the error analysis and correction steps before outputting the final answer, effectively internalizing the $\tau_{reflect}$ operator into its parameters.

5 Experiments

5.1 Benchmark

We evaluate on standard medical benchmarks: MedQA(USMLE test set) (Jin et al., 2021), MedMCQA(validation set), and PubMedQA(test set) (Jin et al., 2019). To thoroughly assess the model’s abilities, we also included medical sections from multi-domain evaluation frameworks, specifically the Health and Biology topics in MMLU-Pro (Wang et al., 2024), as well as the Genetics and Molecular Biology areas in GPQA (Rein et al., 2024). These benchmarks cover tasks in various aspects such as diagnosis, treatment, knowledge, and medical reasoning.

5.2 Baselines and Models

We fine-tuned Qwen2.5-7B-Instruct and Llama-3.1-8B-Instruct using 2k training examples each, and Qwen2.5-32B-Instruct using 30k examples; all training sets were randomly sampled from the MedReflect dataset. To strictly evaluate the efficacy of our approach, we benchmarked MedReflect against a comprehensive suite of baselines, ranging from standard open-source medical and general

389 LLMs to closed-source models. We also explicitly
390 compared our model against reasoning-oriented
391 baselines (e.g., HuatuoGPT-o1) under comparable
392 settings.

393 5.3 Implementation Details

394 We fine-tuned for 3 epochs with $lr = 1e - 4$. Addi-
395 tionally, we applied Low-Rank Adaptation (LoRA)
396 (Shen et al., 2022) with $\alpha = 8$ to fine-tune the LLM.
397 The experiments were carried out on 4 NVIDIA
398 A800 GPUs. More implementation details of train-
399 ing can be found in appendix C.

400 5.4 Main Results

401 We conducted a comprehensive evaluation of open-
402 source LLMs across diverse medical benchmarks,
403 with detailed results presented in Table 1. The
404 experimental findings indicate that MedReflect out-
405 performs comparable open-source models on the
406 majority of evaluated benchmarks, establishing
407 new state-of-the-art results for models of equivalent
408 scale, particularly in reasoning-intensive tasks.

409 **MedReflect Demonstrates Robust Reasoning Ca-**
410 **pabilities** MedReflect-7B (Qwen) achieves sig-
411 nificant improvements over its base model and out-
412 performs other specialized models on key bench-
413 marks such as MedQA and MedMCQA. Notably, it
414 exhibits superior performance on complex reason-
415 ing benchmarks. On MMLU-Pro (Health/Biology)
416 and GPQA, MedReflect-7B surpasses all other
417 models in its size category. MedReflect demon-
418 strates a more balanced and robust performance
419 across both knowledge retrieval and complex clinical
420 reasoning tasks.

421 **Scalability and Competitiveness with Propri-**
422 **etary Models** Scaling up to 32B parameters,
423 MedReflect shows outstanding performance, sur-
424 passing significantly larger open-source models.
425 Crucially, MedReflect-32B remains highly com-
426 petitive with leading proprietary models on spe-
427 cific tasks. While a performance gap remains on
428 MedQA compared to the strongest proprietary sys-
429 tems, MedReflect-32B significantly bridges the gap
430 between open-source and commercial SOTA mod-
431 els (Achiam et al., 2023; Team et al., 2024; Guo
432 et al., 2025).

433 Generalizability Across Model Architectures

434 To verify the universality of our approach, we ap-
435 plied the reflection mechanism to the Llama-3.1-
436 8B-Instruct architecture. The results show that

MedReflect-8B (Llama) achieves a substantial gain
over its base model, confirming that the reflection
mechanism is model-agnostic. However, consistent
with the base models’ capabilities, the Qwen-based
implementation (MedReflect-7B) retains an overall
performance edge over the Llama-based variant.

The Value of Reflection These advancements
underscore the efficacy of the reflection chain.
MedReflect achieves these results using efficient su-
pervised fine-tuning on a curated dataset,
without relying on the massive computational re-
sources typically required for pre-training or rein-
forcement learning at the scale of 70B+ models.
This suggests that the self-reflection mechanism
efficiently unlocks the latent reasoning capabilities
of LLMs in the medical domain.

453 5.5 Ablation Study

454 We conducted further ablation studies on
455 MedReflect-7B across the MedQA, PubMedQA,
456 MMLU, and GPQA datasets to investigate how the
457 reflection mechanism and training methodologies
458 specifically affect performance. All variants are
459 built upon the Qwen2.5-7B-Instruct backbone and
460 were fine-tuned on the same fixed set of 2,000
461 samples. These samples were randomly selected
462 from the MedReflect dataset and subsequently
463 reformatted for training. This controlled setup
464 ensures that any observed performance differences
465 stem solely from the architectural modifications
466 to the reasoning chain. As shown in Table 2, the
467 results support the following analyses:

Directional reflection significantly outperforms
blind correction We compared MedReflect with
a fine-tuning strategy that excludes the reflection
step (SFT w/o Reflect). It removes the intermediate
reflection process, essentially forcing the model to
learn to jump directly from an error to the correct
answer.

475 Experimental results demonstrate that
476 MedReflect exhibits significant advantages
477 across all benchmarks. The **SFT w/o Reflect**
478 baseline resembles a form of random retry within
479 the semantic space; while the model observes the
480 correct outcome, it lacks the logical pathway to
481 deduce the correct state from the erroneous one.
482 While **SFT w/o Reflect** enhances performance
483 through domain knowledge injection and direct
484 answer supervision, MedReflect achieves superior
485 results by equipping the model with an autonomous
486 error-correction mechanism during the reasoning

| Model | MedQA | MedMCQA | PubMedQA | MMLU-Pro | | GPQA | |
|--------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------------|
| | | | | Health | Biology | Genetics | Molecular Biology |
| <i><=8B Language Models</i> | | | | | | | |
| BioMistral-7B | 45.0 | 40.2 | 66.9 | 27.4 | 49.2 | 28.6 | 38.5 |
| UltraMedical-8B | 71.1 | 58.3 | 77.4 | 55.1 | 66.7 | 41.2 | 48.4 |
| Qwen2.5-7B-Instruct | 57.0 | 55.6 | 55.6 | 50.6 | 70.2 | 36.2 | 49.7 |
| LLaMA-3.1-8B-Instruct | 58.7 | 56.0 | 75.2 | 52.7 | 64.6 | 33.8 | 46.8 |
| HuatuoGPT-o1-8B | 72.6 | 60.4 | 79.2 | 58.7 | 68.2 | 48.8 | 59.7 |
| MedReflect-8B(LLaMa) | 64.9 | 72.4 | 75.3 | 53.1 | 64.2 | 56.7 | 60.1 |
| MedReflect-7B(Qwen) | 75.5 | 77.1 | 75.3 | 62.8 | 75.8 | 65.0 | 60.3 |
| <i>> 8B Language Models</i> | | | | | | | |
| Deepseek-R1 | 90.1 | 78.8 | 77.2 | 79.2 | 90.8 | 65.0 | 75.3 |
| Gemini2.5-Flash | 92.0 | 79.7 | 76.2 | – | 98.6 | – | – |
| Gemini2.5-Pro | 92.6 | 71.1 | 75.8 | – | 98.6 | – | – |
| GPT-4o | 86.5 | 76.1 | 78.4 | – | 98.4 | – | – |
| GPT-o3 | <u>93.3</u> | <u>83.3</u> | 80.0 | – | <u>98.6</u> | – | – |
| UltraMedical-70B | 82.2 | 71.8 | 78.4 | 64.8 | 71.1 | 33.8 | 62.9 |
| OpenBioLLM-70B | 78.3 | 74.0 | 79.0 | – | 93.8 | – | – |
| Qwen2.5-72B-Instruct | 72.7 | 66.2 | 71.7 | 65.3 | 78.8 | 41.2 | 56.8 |
| QwQ-32B-Preview | 72.3 | 65.6 | 73.7 | 62.0 | 78.1 | 37.5 | 64.5 |
| HuatuoGPT-o1-70B | 83.3 | 73.6 | 80.6 | 71.0 | 82.8 | 56.2 | 66.5 |
| MedReflect-32B(Qwen) | 83.5 | 80.1 | 82.3 | 78 | 90.8 | 68.3 | 70.6 |

Table 1: Model performance on biomedical QA benchmarks.

| Model | MedQA | PubMedQA | MMLU-Pro | | GPQA | |
|---------------------------------|-------------|-------------|-------------|-------------|-------------|-------------------|
| | | | Health | Biology | Genetics | Molecular Biology |
| Qwen2.5-7B-Instruct | 57.0 | 55.6 | 50.6 | 70.2 | 36.2 | 49.7 |
| SFT w/o Reflect | 57.3 | 63.3 | 54.3 | 72.8 | 48.3 | 55.5 |
| SFT w/o Reflect Question | 61.8 | 67.3 | 63.2 | 73.1 | 55.0 | 58.4 |
| SFT w/o Reflect Answer | 63.7 | 68.3 | 57.9 | 73.3 | 55.0 | 57.3 |
| MedReflect-7B | 75.5 | 75.3 | 62.8 | 75.8 | 65.0 | 60.3 |

Table 2: Performance comparison across biomedical QA datasets using various training strategies.

process. Furthermore, MedReflect does not simply memorize correction data but rather employs a diagnostic thinking pattern analogous to that of human physicians.

The integrity of the reflection chain is crucial

To deconstruct the functional necessity of different components within the reflection chain, we compared three training configurations, all deriving from the same 2,000-sample dataset:

(1) **SFT w/o Reflect Answer:** Retains only the reflection questions (R_q), removing the detailed reflective answers.

(2) **SFT w/o Reflect Question:** Retains only the reflection answers (R_a), removing the guiding self-inquiry.

(3) **MedReflect:** Includes the full pair of reflection questions and answers (R_q, R_a).

The experimental data indicates that the absence

of any single component leads to performance degradation, although configurations retaining partial reflection chains still outperform the baseline without reflection (**SFT w/o Reflect**). The reflection question (R_q) serves the role of error attribution, explicitly guiding attention toward potential defects in the reasoning chain and setting a clear direction for correction. The reflection answer (R_a) serves the role of knowledge retrieval and consolidation, utilizing a declarative tone to invoke the internal knowledge base and provide a factual basis for the correction. Therefore, the complete MedReflect mechanism constructs a closed cognitive loop, where the combination of R_q and R_a plays a critical role in effectively enhancing the model’s reasoning capabilities.

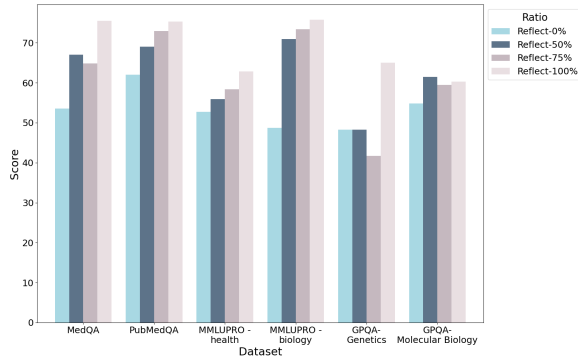


Figure 4: Results of the Reflection Data Proportion Experiment

5.6 Analysis

5.6.1 Data Efficiency and Marginal Utility

To investigate the marginal utility of reflection data during training, we conducted experiments with four different proportions of reflection data (0%, 50%, 75%, and 100%) while maintaining a constant total volume of 2,000 training samples. As shown in Figure 4, the results reveal a significant positive correlation between the proportion of reflection data and model performance.

Significant Marginal Returns of Reflection Data

As the proportion of reflection data in the training set increases, the model’s accuracy across the vast majority of medical and biological reasoning tasks demonstrates a steady upward trend. Although minor fluctuations are observed in isolated datasets, the overall trajectory indicates that a high proportion of reflection data effectively activates the model’s deep reasoning capabilities. Notably, this benefit is maximized when the reflection proportion reaches 100%.

Efficiency Analysis Under the constraint of a constant total sample size, we progressively replaced standard Correct Data with Reflection Data containing detailed reasoning processes. The results show that this substitution yields significant performance gains. This implies that, per unit data sample, data containing reflection processes possesses higher information density and training efficiency. By explicitly demonstrating error correction and logical deduction, reflection data compels the model to learn reflective thinking rather than mere rote memorization.

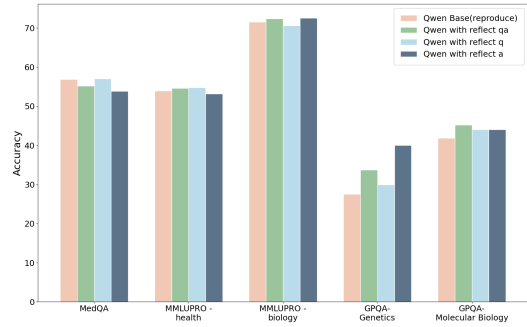


Figure 5: Results of the Attribution Analysis Experiment

5.6.2 Effectiveness of Reflection Content

To study the effectiveness of the reflective content produced during MedReflect generation, we collected the reflection questions R_Q , reflection answers R_A , and reflection pairs R_{QA} generated by the MedReflect-7B model, and provided them as additional context to the base Qwen2.5-7B-Instruct model. The results are shown in Figure 5, and additional experimental details are provided in Appendix C.5. As illustrated in Figure 5, overall, augmenting the original Qwen2.5-7B-Instruct model with reflection information consistently obtains better performance than the base model alone, indicating that the reflections produced by MedReflect are beneficial. On highly specialized datasets such as GPQA, both reflection questions and reflection answers lead to substantial improvements. Moreover, using only reflection questions (without answers) still brings consistent gains, suggesting that MedReflect’s reflection questions help steer the direction of reasoning.

6 Conclusion

In this study, we propose MedReflect, a framework designed to enable LLMs to autonomously perform reflection and revision during medical tasks. By leveraging a lightweight LLM to construct a low-cost, diverse reflection training dataset, we trained models to acquire physician-like reflective thinking characterized. This training paradigm itself exhibits conspicuous cost-effectiveness. Notably, fine-grained reflection design proves crucial for fully activating the model’s reflective capabilities. Finally, we demonstrate that MedReflect significantly enhances model performance across multiple medical QA benchmarks, robustly validating the efficacy of the reflection mechanism in improving LLM accuracy in medical tasks.

591 Limitation

592 While MedReflect significantly enhances diagnostic
593 accuracy by mimicking the physician’s reflective
594 thinking process, this performance gain comes
595 with an inherent trade-off in inference efficiency.
596 As described in our method, the framework explicitly
597 inserts a reflection chain before generating the
598 final modified decision. Consequently, this mechanism
599 inevitably increases the total number of generated
600 tokens per query compared to direct-answer models.
601 The increased sequence length results in higher
602 computational costs and longer inference latency,
603 which may pose challenges for deployment in real-time
604 or resource-constrained medical consultation scenarios.
605 Future work could explore methods to trigger reflection
606 adaptively only when the model detects high uncertainty,
607 thereby balancing accuracy and efficiency.
608

609 Ethics Statement

610 MedReflect is built upon the Qwen2.5-Instruct
611 architecture and, like all Large Language Models
612 (LLMs), inherits intrinsic limitations such as the
613 susceptibility to hallucinations and the potential
614 generation of counterfactual medical advice.
615 In high-stakes medical decision-making scenarios,
616 any inaccuracies or misinterpretations of clinical
617 narratives could lead to adverse outcomes.

618 Therefore, MedReflect is intended solely for
619 research purposes to explore the potential of autonomous
620 reflective thinking in medical AI. It is not designed
621 for direct clinical deployment without expert human
622 oversight. Researchers and developers utilizing this
623 framework must be cognizant of these risks and
624 implement robust safeguards, such as rigorous
625 secondary validation loops. We further declare
626 that the training data derived from public datasets
627 (ChatDoctor and MedMCQA) was processed strictly
628 for academic research, adhering to the data
629 filtering protocols described in this work to ensure
630 quality and privacy compliance.

631 References

632 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
633 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
634 Diogo Almeida, Janko Altenschmidt, Sam Altman,
635 Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical
636 report. *arXiv preprint arXiv:2303.08774*.

637 C. Chen, X. Wang, and M. Guan. 2025. Evaluating and
638 improving syndrome differentiation thinking ability
639 in llms. *arXiv preprint*.

Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang,
Wanlong Liu, Rongsheng Wang, Jianye Hou, and
Benyou Wang. 2024. Huatuogpt-o1, towards medical
complex reasoning with llms. *arXiv preprint
arXiv:2412.18925*. 640
641
642
643
644

Junying Chen, Xidong Wang, Ke Ji, Anningzhe Gao,
Feng Jiang, Shunian Chen, Hongbo Zhang, Dingjie
Song, Wenya Xie, Chuyi Kong, and 1 others. 2023.
Huatuogpt-ii, one-stage training for medical adaption
of llms. *arXiv preprint arXiv:2311.09774*. 645
646
647
648
649

Chengfeng Dou, Ying Zhang, Yanyuan Chen, Zhi Jin,
Wenpin Jiao, Haiyan Zhao, and Yu Huang. 2024. *Detection, diagnosis, and explanation: A benchmark for Chinese medial hallucination evaluation*. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4784–4794, Torino, Italia. ELRA and ICCL. 650
651
652
653
654
655
656
657

Kanishk Gandhi, Denise Lee, Gabriel Grand, Muxin
Liu, Winson Cheng, Archit Sharma, and Noah D
Goodman. 2024. Stream of search (sos): Learning to
search in language. *CoRR*. 658
659
660
661

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
llms via reinforcement learning. *arXiv preprint
arXiv:2501.12948*. 662
663
664
665
666
667

Minbyul Jeong, Jiwoong Sohn, Mujeen Sung, and Jae-
woo Kang. 2024. Improving medical reasoning
through retrieval and self-reflection with retrieval-
augmented large language models. *Bioinformatics*,
40(Supplement_1):i119–i129. 668
669
670
671
672

Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko
Ishii, and Pascale Fung. 2023. Towards mitigating
llm hallucination via self reflection. In *Findings
of the Association for Computational Linguistics:
EMNLP 2023*, pages 1827–1843. 673
674
675
676
677

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. 678
679
680
681
682

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W
Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset
for biomedical research question answering. *arXiv
preprint arXiv:1909.06146*. 683
684
685
686

Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han,
and Rui Zhang. 2024. When can llms actually correct
their own mistakes? a critical survey of self-
correction of llms. *Transactions of the Association
for Computational Linguistics*, 12:1417–1440. 687
688
689
690
691

Taeyoon Kwon, Kai Tzu-iunn Ong, Dongjin Kang,
Seungjun Moon, Jeong Ryong Lee, Dosik Hwang,
Beomseok Sohn, Yongsik Sim, Dongha Lee, and Jinyoung
Yeo. 2024. Large language models are clinical 692
693
694
695

| | | | |
|-----|--|--|---|
| 696 | reasoners: Reasoning-aware diagnosis framework with prompt-generated rationales. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 38, pages 18417–18425. | A. Ríos-Hoyo, N.L. Shan, and A. Li. 2024. Evaluation of llms as a diagnostic aid for complex medical cases. <i>arXiv preprint arXiv:2410.xxxxx</i> . | 750 751 752 |
| 700 | Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. <i>arXiv preprint arXiv:2402.10373</i> . | Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and 1 others. 2022. Lora: Low-rank adaptation of large language models. | 753 754 755 |
| 705 | Gengliang Li, Rongyu Chen, Bin Li, Linlin Yang, and Guodong Ding. 2025a. Medfact-r1: Towards factual medical reasoning via pseudo-label augmentation. <i>arXiv preprint arXiv:2509.15154</i> . | S. Soffer, V. Sorin, and G.N. Nadkarni. 2025. Pitfalls of large language models in medical ethics reasoning. <i>arXiv preprint</i> . | 756 757 758 |
| 709 | Jiaqi Li, Xinyi Dong, Yang Liu, Zhizhuo Yang, Quansen Wang, Xiaobo Wang, Song-Chun Zhu, Zixia Jia, and Zilong Zheng. 2025b. Reflectevo: Improving meta introspection of small llms by learning self-reflection. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 16948–16966. | 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. <i>arXiv preprint arXiv:2403.05530</i> . | 759 760 761 762 763 764 |
| 715 | Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. <i>Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge</i> . Preprint, arXiv:2303.14070. | Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. <i>arXiv preprint arXiv:2503.19786</i> . | 765 766 767 768 769 |
| 720 | Chihung Lin and Chang-Fu Kuo. 2025. <i>Roles and potential of large language models in healthcare: A comprehensive review</i> . <i>Biomedical Journal</i> , page 100868. | Juraj Vladika, Phillip Schneider, and Florian Matthes. 2024. Medrequal: Examining medical knowledge recall of large language models via question answering. <i>arXiv preprint arXiv:2406.05845</i> . | 770 771 772 773 |
| 724 | X. Liu and 1 others. 2025. Medical large language model for diagnostic reasoning across specialties. <i>arXiv preprint</i> . | Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. <i>Advances in Neural Information Processing Systems</i> , 37:95266–95290. | 774 775 776 777 778 779 780 |
| 727 | Yuxing Lu, Gecheng Fu, Wei Wu, Xukai Zhao, Sin Yee Goi, and Jinzhuo Wang. 2025. <i>Doctorrage: Medical rag fusing knowledge with patient analogy through textual gradients</i> . Preprint, arXiv:2505.19538. | Juncheng Wu, Wenlong Deng, Xingxuan Li, and 1 others. 2025. Medreason: Eliciting factual medical reasoning steps in llms via knowledge graphs. <i>arXiv preprint arXiv:2504.00993</i> . | 781 782 783 784 |
| 731 | Ankit Pal, Logesh Kumar Umapathi, and Malaikandan Sankarasubbu. 2022. Medmqca: A large-scale multi-subject multi-choice dataset for medical domain question answering. In <i>Conference on health, inference, and learning</i> , pages 248–260. PMLR. | Junde Wu, Jiayuan Zhu, Yunli Qi, Jingkun Chen, Min Xu, Filippo Menolascina, and Vicente Grau. 2024. <i>Medical graph rag: Towards safe medical large language model via graph retrieval-augmented generation</i> . Preprint, arXiv:2408.04187. | 785 786 787 788 789 |
| 736 | Tian Qin, David Alvarez-Melis, Samy Jelassi, and Eran Malach. 2025. To backtrack or not to backtrack: When sequential search limits model reasoning. <i>arXiv preprint arXiv:2504.07052</i> . | Zhiheng Xi, Dingwen Yang, Jixuan Huang, Jiafu Tang, Guanyu Li, Yiwen Ding, Wei He, Boyang Hong, Shihan Dou, Wenyu Zhan, and 1 others. 2024. Enhancing llm reasoning via critique models with test-time and training-time supervision. <i>CoRR</i> . | 790 791 792 793 794 |
| 740 | Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. 2024. Recursive introspection: Teaching language model agents how to self-improve. <i>Advances in Neural Information Processing Systems</i> , 37:55249–55285. | Dexuan Xu, Yanyuan Chen, Zhongyan Chai, Yifan Xiao, Yandong Yan, Weiping Ding, Hanpin Wang, Zhi Jin, Wenpin Jiao, Weihua Yue, and 1 others. 2025. Knowledge fusion in deep learning-based medical vision-language models: A review. <i>Information Fusion</i> , page 103455. | 795 796 797 798 799 800 |
| 745 | David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In <i>First Conference on Language Modeling</i> . | Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. 2024. Saysself: Teaching llms to express confidence with self-reflective rationales. <i>arXiv preprint arXiv:2405.20974</i> . | 801 802 803 804 805 |

806 Guanhao Ye, Khiem Duc Pham, Xinzhi Zhang,
807 Sivakanth Gopi, Baolin Peng, Beibin Li, Janardhan
808 Kulkarni, and Huseyin A Inan. 2025. On the emer-
809 gence of thinking in llms i: Searching for the right
810 intuition. *arXiv preprint arXiv:2502.06773*.

811 Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-
812 Zhu. 2024. Physics of language models: Part 2.2,
813 how to learn from mistakes on grade-school math
814 problems. *arXiv preprint arXiv:2408.16293*.

815 Kaiyan Zhang, Sihang Zeng, Ermo Hua, Ning Ding,
816 Zhang-Ren Chen, Zhiyuan Ma, Haoxin Li, Ganqu
817 Cui, Biqing Qi, Xuekai Zhu, and 1 others. 2024.
818 Ultramedical: Building specialized generalists in
819 biomedicine. *Advances in Neural Information Pro-
820 cessing Systems*, 37:26045–26081.

821 Yuxuan Zhou, Yubin Wang, Bin Wang, Chen Ning, Xien
822 Liu, Ji Wu, and Jianye Hao. 2025. Enhancing the
823 medical context-awareness ability of llms via mul-
824 tifaceted self-refinement learning. *arXiv preprint
825 arXiv:2511.10067*.

826 Xudong Zhu, Jiachen Jiang, Mohammad Mahdi Khalili,
827 and Zhihui Zhu. 2025. From emergence to control:
828 Probing and modulating self-reflection in language
829 models. *arXiv preprint arXiv:2506.12217*.

A Implementation Details for Data Generation

As outlined in the Methods section, we utilized the MedMCQA dataset to construct reflections at the sentence level and the ChatDoctor dataset to build reflections at the word level. The implementation details of this process are described below, with all corresponding prompts provided in Table 5 and Table 6.

Implementation Detail for reflect pinpoint generation

RG1: MedMCQA Given a multiple-choice question, we employ Prompt I to a LLM, eliciting both an initial response and its reasoning process. The model is required to strictly adhere to a predefined format specification for its final answer. Subsequently, we compare the model’s generated final answer option against the standard answer: samples where the options differ are retained, while those that match are discarded. Then, using Prompt II, we identify where in the explanation the erroneous option manifests, isolating the corresponding segment of text containing the critical reasoning flaw. This segment serves as a pinpointed reflection target.

RG2: ChatDoctor Following text preprocessing, we employ Prompt III to guide the model in performing the entity extraction task. Subsequently, we randomly select 1 to 3 recognized entities within a sentence to mask, replacing each with entity type. We then utilize Prompt IV to instruct the model to fill these masked entities. Through multiple retries, we calculate a post-retry similarity score. Entities yielding a score below 0.8 are deemed incorrectly after retries and are retained in their masked state entity type. We repeat this step five times to obtain diverse data.

Implementation Detail for Reflection Generation

The generation of reflection questions and answers then proceeds as follows: Generating Reflection Questions: The erroneous masked entity along with the original question are input to the model using Prompt V to generate guiding reflection questions. These questions specifically focus on the relationship between the erroneous entity and the correct answer.

Generating Reflective Answers: Prompt VI is used to direct the model to generate reflective answers based on its own knowledge and capabilities,

| | | | |
|-----|---|--|-----|
| 879 | addressing the questions generated in the previous | table 4, each sample in the MedMCQA dataset | 929 |
| 880 | step. | (21,107 samples) contains only a single reflection | 930 |
| 881 | Generating Corrected Candidate Answers: | instance (Reflect avg = 1.0), with reflections primar- | 931 |
| 882 | Building on the reflection content from steps 1 and | ily occurring at the sentence level. In contrast, the | 932 |
| 883 | 2, Prompt VII guides the model to produce multiple | ChatDoctor dataset is significantly larger (36,413 | 933 |
| 884 | candidate answers incorporating corrected entities. | samples) and features richer, more granular reflec- | 934 |
| 885 | Each generation step incorporates a retry mecha- | tion content. It encompasses a total of 55,460 re- | 935 |
| 886 | nism to ensure valid output. Finally, the corrected | flexion instances, averaging 1.52 reflections per | 936 |
| 887 | candidate answers undergo strict validation requir- | sample. This indicates a substantial number of | 937 |
| 888 | ing inclusion of all ground-truth correct entities to | dialogues within ChatDoctor involve diverse reflec- | 938 |
| 889 | be accepted. | tive processes. Furthermore, ChatDoctor focuses | 939 |
| 890 | Data Filtering To guarantee the quality of the re- | specifically on reflecting medical entities, enabling | 940 |
| 891 | flexion data, we first preprocessed the raw datasets | precise identification of cognitive biases at the en- | 941 |
| 892 | to remove samples with short reasoning paths or | tity level and capable of generating 1 to 3 distinct | 942 |
| 893 | content irrelevant to medical question-answering. | reflection points per response. This structural dis- | 943 |
| 894 | Following this, we conducted a secondary quality | parity between the datasets provides the essential | 944 |
| 895 | assessment to mitigate potential confirmation bias | data foundation for training models to understand | 945 |
| 896 | in the model-generated reflections. | and generate reflections at varying levels of granu- | 946 |
| 897 | This process entailed feeding the generated re- | larity. | 947 |
| 898 | flexion questions (R_q) and answers (R_a), together | | |
| 899 | with the model’s initial erroneous trajectory, back | | |
| 900 | into the model to generate a revised answer. Cru- | C More Implementation Details | 948 |
| 901 | cially, the validity of a reflection was determined | | |
| 902 | by verifying the revised answer against the objec- | To ensure the reproducibility of MedReflect and | 949 |
| 903 | tive ground truth provided by the datasets, ensuring | address the specifics of our training and inference | 950 |
| 904 | that the correction was not merely endorsed by the | mechanisms, we provide detailed configurations of | 951 |
| 905 | model itself but aligned with the gold standard: | training. | 952 |
| 906 | (1) For RG1, we utilized the ground-truth option | | |
| 907 | of the multiple-choice question as the standard. A | C.1 Training Details. | 953 |
| 908 | reflection instance was considered valid only if | | |
| 909 | it guided the model to select the correct option | Training Data Statistics We tailored the scale of | 954 |
| 910 | matching the dataset label. | the fine-tuning dataset to the capacity of the target | 955 |
| 911 | (2) For RG2, the objective was to restore masked | models to ensure optimal adaptation. Specifically, | 956 |
| 912 | medical entities. A reflection was deemed valid if | we utilized 2k training examples for the Qwen2.5- | 957 |
| 913 | the model, guided by the reflection, correctly pre- | 7B-Instruct model. For the larger Qwen2.5-32B- | 958 |
| 914 | dicted the original entity present in the raw medical | Instruct model, we scaled up the training set to 30k | 959 |
| 915 | text. | examples to prevent under-fitting and ensure robust | 960 |
| 916 | To further ensure the robustness of the reflec- | performance. | 961 |
| 917 | tion data, we repeated this evaluation process 10 | | |
| 918 | times for each instance. We retained a reflection | Configuration of Supervised Fine-Tuning (SFT). | 962 |
| 919 | instance only if the model successfully reached the | We fine-tuned the base models (Qwen2.5-7B/32B- | 963 |
| 920 | correct ground truth in at least 8 out of these 10 | Instruct) using the generated MedReflect dataset. | 964 |
| 921 | trials. Instances failing to meet this threshold were | Table 7 details the specific Configurations used in | 965 |
| 922 | filtered out to exclude unstable or stochastic reason- | our experiments. We utilized the AdamW opti- | 966 |
| 923 | ing paths. Ultimately, this filtering process yielded | mizer with a cosine learning rate scheduler. The | 967 |
| 924 | a dataset comprising 36,413 medical consultation | training was performed on $4 \times$ NVIDIA A800 | 968 |
| 925 | records and 21,107 multiple-choice questions. | (80GB) GPUs. | 969 |
| 926 | B Dataset Analysis | | |
| 927 | The key datasets underpinning this study are | | |
| 928 | MedMCQA and ChatDoctor. As shown in the | | |

| Metric | MedQA | PubMedQA | MMLU-Pro | | GPQA | |
|--|--------|----------|----------|---------|----------|-------------------|
| | | | Health | Biology | Genetics | Molecular Biology |
| Avg. Generated Tokens (All) | 711.18 | 321.31 | 641.96 | 721.19 | 840.35 | 730.50 |
| Avg. Generated Tokens (Correct) | 711.39 | 320.14 | 638.30 | 674.78 | 904.73 | 697.90 |
| Avg. Generated Tokens (Incorrect) | 710.90 | 322.20 | 645.94 | 810.58 | 761.67 | 761.58 |

Table 3: Statistics of average response length (number of tokens) across different medical benchmarks.

| dataset | Num | Total_reflections | Length |
|------------|-------|-------------------|--------|
| MedMCQA | 21107 | 21107 | 394.56 |
| ChatDoctor | 36413 | 55460 | 305.05 |

Table 4: Data Analysis

| Configuration | Value |
|---------------------------|------------------------------------|
| Global Batch Size | 4 |
| Learning Rate | $1e - 4$ |
| Epochs | 3 |
| Max Sequence Length | 4096 |
| Optimizer | AdamW |
| Weight Decay | 0.01 |
| Warmup Ratio | 0.03 |
| LR Scheduler | Cosine |
| Gradient Clipping | 1.0 |
| LoRA Configuration | |
| LoRA Rank (r) | 16 |
| LoRA Alpha (α) | 8 |
| LoRA Dropout | 0.05 |
| Target Modules | q, k, v, o, gate, up, down_proj |

Table 7: Detailed configuration for SFT

```

<|im_start|>system
You are a medical expert skilled in answer
questions while reflecting. Your reflecting
goal is to generate a good reflection to
assist you to improve your previous answer
and continue to answer the question.
<|im_end|>
<|im_start|>user
{Input Question}
<|im_end|>
<|im_start|>assistant
{Target Response}
<|im_end|>

```

Figure 6: The training prompt template used for MedReflect.

C.2 Evaluation Details

Inference Configuration During the inference stage, we utilized the vllm library. We employed **Nucleus Sampling** to ensure the diversity and robustness of the generated responses. The specific generation parameters are listed in Table 8.

| Parameter | Value |
|--------------------|-------|
| Temperature | 0.7 |
| Top-p | 0.9 |
| Repetition Penalty | 1.05 |
| Max New Tokens | 4096 |

Table 8: Inference Configurations

C.3 Evaluation Protocol

To mitigate the variance arising from generation randomness and ensure the reliability of our results, we adopted a multi-run evaluation strategy. For standard benchmarks including MedQA, MMLU, PubMedQA, and MedMCQA, we conducted 3 independent runs and reported the average accuracy. For the GPQA dataset, given its relatively small test set size which is susceptible to higher statistical fluctuation, we increased the number of independent runs to 5 and reported the average performance to ensure a more stable assessment.

Prompt Template To ensure the model adheres to the desired instruction-following format and medical reasoning style, we employed a specific template during the supervised fine-tuning stage. The input sequence is formatted as follows:

993 C.4 Answer Extraction and Verification 994 Strategy

995 To standardize the evaluation process and mitigate
996 discrepancies arising from formatting inconsisten-
997 cies in model outputs, we implemented a robust
998 two-stage protocol for answer extraction across all
999 datasets.

1000 **Stage 1: Deterministic Pattern Matching.** Ini-
1001 tially, we employed a strict rule-based extraction
1002 mechanism using regular expressions. We defined
1003 a set of heuristic patterns to capture conclusive
1004 answers, such as matches for "The answer is
1005 option [X]" or "My answer is [X]". If a dis-
1006 tinct option key (A, B, C, or D) was successfully
1007 extracted, it was directly compared with the ground
1008 truth label.

1009 **Stage 2: LLM-Assisted Semantic Parsing.** To
1010 address instances where the model’s response struc-
1011 ture was irregular or failed the regex matching (i.e.,
1012 the answer was implicitly embedded in the reason-
1013 ing chain without a standard prefix), we utilized
1014 a superior-capability LLM (GPT-4) as a semantic
1015 parser. This acts as a fallback mechanism to inter-
1016 pret the semantic intent of the model’s output.

1017 Specifically, we utilized the gpt-4-0613 version
1018 via API with a temperature setting of 0 to ensure
1019 deterministic and reproducible outputs. The model
1020 was tasked to act as an objective evaluator. We con-
1021 structed a specific prompt template (detailed in Ta-
1022 ble 10) that includes the original question, the can-
1023 didate options, and the model’s raw response. GPT-
1024 4 was instructed to identify the final chosen option
1025 and output a single character. This hybrid approach
1026 ensures that valid responses are not penalized due
1027 to minor formatting deviations, thereby providing a
1028 more accurate assessment of the model’s reasoning
1029 capabilities.

1030 C.5 Evaluation setup of Experiment in 1031 Section 5.6.2

1032 In the experiments presented in Section 5.6.2, the
1033 reflections (comprising reflection questions and an-
1034 swers) within the MedReflect-7B responses were
1035 extracted. This reflection information was subse-
1036 quently utilized as context and concatenated with
1037 the original query to directly input into the original
1038 Qwen2.5-7B-Instruct. The prompt template used
1039 for this inference is shown in Table 9.

| Component | Content Details |
|---------------------------|---|
| System Role | You are an excellent medical exam assis- tant, skilled at extracting correct informa- tion from complex responses. |
| Instruction | I will provide you with a question and the response from a large model. Please help determine which answer the large model believes is correct. If you think it does not answer the multiple-choice question, return None. No explanation is needed, just return the option itself.question: questionanswer: ans |
| Input For- mat | Question: [Insert Question Text] Options: (A) [Option A] (B) [Option B] (C) [Option C] (D) [Option D] |

Table 10: The prompt template employed for the GPT-4 based answer extraction stage. The model acts as a semantic parser to extract the final answer key from unstructured reasoning chains.

1040 D More Qualitative Analysis

1041 D.1 Analysis of the Efficiency Boundary for 1042 Reflection Data

| Dataset | Reflect_2k | Reflect_30k |
|------------------------|------------|-------------|
| MedQA | 75.5 | 74.3 |
| PubMedQA | 75.3 | 74.0 |
| MMLU-health | 62.8 | 61.0 |
| MMLU-biology | 75.8 | 72.9 |
| GPQA-Genetics | 65.0 | 63.0 |
| GPQA-Molecular Biology | 60.3 | 59.5 |

Table 11: Results of the Data Efficiency Experiment of MedReflect-7B

1043 We investigated the impact of training data magni-
1044 tude on model performance by conducting a com-
1045 parative analysis between a compact dataset of
1046 2,000 samples and a significantly larger dataset
1047 of 30,000 samples. As indicated by the results
1048 shown in Table 11, the model trained on only 2,000
1049 samples achieves a performance level compar-
1050 able to that of the model trained on 30,000 sam-
1051 ples. Across various metrics, the smaller dataset
1052 yields results commensurate with those of the larger
1053 dataset. This observation underscores the excep-
1054 tional data efficiency of the proposed reflection
1055 mechanism, suggesting that the model can effec-
1056 tively internalize the core patterns of reflective rea-
1057 soning with a limited volume of high-quality in-
1058 struction data.

D.2 Response Length.

We further report the response length of MedReflect, measured by the number of generated tokens (excluding the prompt). Table 3 summarizes the average length over all samples, as well as the average length of correct and incorrect responses. On MedQA and PubMedQA, correct and incorrect responses have almost the same average length. On MMLUPro-health, correct responses are slightly shorter on average. In contrast, on MMLUPro-biology and GPQA-Molecular Biology, incorrect responses are notably longer than correct ones. On the highly specialized GPQA-Genetics benchmark, correct responses are substantially longer than incorrect ones (904.7 vs. 761.7 tokens on average), suggesting a stronger length–correctness association on this dataset. Overall, response length appears dataset-dependent under MedReflect.

E Case Study

Cases for Multi-Choice Question The multi-choice cases are presented in Figure 7, where we compare the performance of MedReflect, HuatuoGPT-o1, and Deepseek-R1.

I We conducted a specific analysis of the performance of MedReflect and the general reasoning large model (DeepSeek-R1), using septicemia pathogen identification as a representative clinical scenario. Experimental results reveal that both DeepSeek-R1 and MedReflect initially exhibit similar diagnostic errors when confronted with complex clinical evidence. However, a critical divergence emerges in their subsequent reasoning pathways:

- DeepSeek-R1, constrained by its reasoning paradigm, terminates analysis upon reaching an incorrect conclusion. It demonstrates no self-correction, maintaining confidence in its initial output.
- MedReflect, leveraging its intrinsic reflective architecture, critically evaluates morphological inconsistencies – questioning the mismatch between observed coccobacilli features and *Neisseria gonorrhoeae*’s typical diplococcal structure. It further assesses clinical epidemiology to validate diagnostic hypotheses.

This structured self-reflection mechanism enables MedReflect to iteratively retrieve evidence from the knowledge base and dynamically adjust diagnostic weights, thereby achieving autonomous er-

ror correction and precise diagnosis. In contrast, DeepSeek-R1 despite engaging in constant deep thinking, lack a mechanism for reflection and adjustment of initial responses.

II We conducted a specific analysis of the reasoning differences between MedReflect and HuatuoGPT-o1, using mass conversion problems as a representative stoichiometric scenario.

- MedReflect first incorrectly identified ethane. Then, MedReflect exhibited dynamic error correction capability. When significant deviations from given options were detected in its computed results, MedReflect implicitly switched to the correct molecular formula for recalculation, ultimately arriving at the correct answer (G).
- HuatuoGPT-o1 manifests a fundamental flaw in its misconception of core principles, in which case it misinterprets the stoichiometric ratio. This foundational error causes the entire computational framework to deviate from theoretical benchmarks. Then, when the calculated result (122.4g) shows significant deviation from options, the model arbitrarily selects the closest match (I) without verifying the validity of stoichiometric assumptions, which indicates a lack of scientific verification or reflection.

Overall, MedReflect tends to follow the physician-like reasoning process and adjust its reasoning process better. In contrast, HuatuoGPT-o1 failed to promptly correct its initial error, compounding the mistake until it ultimately settled for an answer it deemed a close approximation.

Cases for Consulting Question Using a diagnostic case of pediatric respiratory infection presenting with lingual symptoms, we evaluated the performance of three models: DeepSeek, MedReflect, and HuatuoGPT-o1. The performance is shown in Figure 8.

- MedReflect demonstrated more rigorous clinical reasoning through its dynamic reflection mechanism. After generating an initial response, it actively challenged its preliminary assumptions. This process involved critically analyzing relevant pathological mechanisms and prioritizing non-infectious etiologies. Its treatment recommendations appropriately advocated for conservative interventions such as avoiding irritants and dietary modifications.

| | |
|---|---|
| <p>Q: A 39-year-old female emergency patient presented with fever, chills, left lower abdominal pain..... Laboratory tests showed low platelets, low fibrinogen and high D-dimer.most likely blood culture is: A: Coagulase-positive B: Encapsulated C: Spore-forming D: Lactose-fermenting E: Gamma-hemolytic</p> | |
| <p>MedReflect-32B</p> <p>D.</p> <p>1. Analyze background information The patient presents.....Given the clinical context and laboratory findings, the most likely pathogen causing this severe infection with sepsis and DIC is a Gram-negative bacterium, particularly one associated with pelvic infections.</p> <p>2. First decision making[Wrong] Therefore, the most appropriate answer is: Option B.....</p> <p>3. Self-reflect Wait, let me think about it.<Think> Which type of bacteria, based on its characteristics and typical growth medium, is more likely to be associated with severe systemic infections and sepsis in the context of a patient presenting with fever, chills, and evidence of pelvic inflammatory disease?.....</Think></p> <p>4. Modify and final decision making[True] <Modified>Option D: Lactose-fermenting, gram-negative rods forming pink colonies on MacConkey agar</Modified>Therefore, my answer is: option D</p> | <p>HuatuoGPT-o1-70B</p> <p>A.</p> <p>1. Analyze background information We've got a 39-year-old woman</p> <p>2. Deep Reasoning and Analyzing that's pointing towards some sort of clotting issue, like disseminated intravascular coagulation, or DIC. And the high D-dimer supports that too, since it indicates clotting activity. They did some kind of test with phenol on her blood, and they found this compound with a phosphorylated N-acetylglucosamine dimer and 6 fatty acids. Wait a minute, that sounds like lipoteichoic acid, which is part of the cell wall of gram-positive bacteria. So, if we're thinking about gram-positive bacteria,That's often linked to Staphylococcus aureus, which is a gram-positive cocci that's coagulase-positive.</p> <p>3. Decision Making Let's check the options. Option A mentions coagulase-positive, gram-positive cocci that grow on methicillin-containing agar. That's MRSA, which is a type of Staphylococcus aureus. This fits well with our findings so far. Therefore, my answer is A.</p> |
| <p>Q: Using this formula, if a reaction was provided with 84g of ethane and unlimited oxygen, how many grams of carbon dioxide would result? The unbalanced reaction of ethane gas to carbon dioxide and water is as follows $C_2H_4 + O_2 \rightarrow CO_2 + H_2O$ A: 78g, B: 210g, C: 312g, D: 528g, E: 156g, F: 44g, G: 264g, H: 180g, I: 132g, J: 396g</p> | |
| <p>MedReflect-7B</p> <p>G.</p> <p>1. Analyze background information To solve this problem, we need to first balance the chemical equation for the reaction of ethane (C_2H_4) with oxygen (O_2) to form carbon dioxide (CO_2) and water (H_2O):</p> <p>2. First Calculating making[Wrong] However, based on the given choices and the closest logical fit: Option I: 132g is the closest to the actual calculated value of 88 grams.</p> <p>3. Self-reflect Wait, let me think about it.<Think> How does the actual yield of carbon dioxide compare to the calculated amount when using 84 grams of ethane in the reaction with an excess of oxygen?.....</Think></p> <p>4. Modify and final decision making[True] <Modified>Option G-If we proceed with the correct conversion,,</Modified>Therefore, my answer is: option G</p> | <p>HuatuoGPT-o1-8B</p> <p>A.</p> <p>1. Analyze background information Alright, let's figure out how much carbon dioxide we can get from 84 grams of ethane.</p> <p>2. Deep Reasoning and Analyzing Now, with 84 grams of ethane, I can find out how many moles there are.....Next, I need to see what happens when ethane reacts with oxygen.....But wait, let's double-check that. Hmm, maybe I missed something.</p> <p>3. Decision Making Maybe there's a small rounding issue or some other factor I didn't consider. It's possible that the numbers don't exactly match due to these kinds of things. So, I'm going to go with option I, which is 132 grams.</p> |

Figure 7: Cases of Multi-Choice Questions

- 1156 • HuatuoGPT-o1 also engaged in deep reason-
1157 ing and reflection during its analysis. How-
1158 ever, its reflective process was comparatively
1159 superficial. Rather than conducting an in-
1160 depth exploration of disease etiology, it em-
1161 ployed overly simplistic reasoning, such as ex-
1162 cluding diagnoses based solely on "the child
1163 not complaining of pain".
- 1164 • DeepSeek-R1 systematically listed potential
1165 diagnoses including lingual papillitis, oral
1166 thrush, and viral exanthems. However, its
1167 diagnostic suggestions exhibited a scattered
1168 approach, failing to effectively integrate core
1169 clinical clues. Of course, as a powerful gen-
1170 eral reasoning model, it may be difficult for
1171 it to directly offer suggestions, so it provides
1172 many options for patients to choose.

1173 MedReflect's iterative reasoning model, grounded
1174 in clinical clue integration, differs significantly
1175 from DeepSeek's divergent deep reasoning. It more
1176 closely aligns with medical differential diagnostic
1177 logic. Simultaneously, its reflective process demon-
1178 strates greater rigor and professionalism compared
1179 to HuatuoGPT-o1's relatively simplistic and cur-
1180 sory approach.

| | | |
|---|--|---|
| <p>Q: My grandson is 6 years old. He has had a cold for several weeks and still has a slight fever. The cold and fever have subsided, but he still has a runny nose and cough. Now, there are small white bumps on his tongue. What is the cause? What medicine can he take?</p> | | |
| <p>MedReflect-7B</p> <p>1. First diagnosis decision making[Wrong] White bumps on the tongue of your 6-year-old grandson can result from oral thrush due to an infection, bumps, or cold sores caused by a virus.</p> <p>2. Self-reflect Wait, let me think about it.<Think> How does the actual yield of carbon dioxide compare to the calculated amount when using 84 grams of ethane in the reaction with an excess of oxygen?.....</Think></p> <p>3. Modify and final decision making[True] <Modified>White bumps on the tongue of your 6-year-old grandson can result from inflamed papillae due to an infection, canker sores, or cold sores caused by a virus.</p> <p>4. First treatment decision making[Wrong] You should get If the bumps are painful, you can apply a small amount of antifungal cream to the bumps.</p> <p>5. Self-reflect Wait, let me think about it. <Think> Given that the condition involves white bumps on the tongue possibly due to inflamed papillae what type of treatment could provide relief from the potential pain or discomfort associated with these bumps?</p> <p>6. Modify and final treatment decision making[True] <Modified> If the bumps are painful, you can apply a small amount of Watermelon Frost Spray to the bumps. </Modified> Stick to soft and bland foods.</p> | <p>HuatuoGPT-o1-8B</p> <p>1. Analyze background information Now, he's left with this annoying runny nose and a cough.</p> <p>2. The first consideration of the cause What could cause white bumps on the tongue? I remember thrush can cause that. It's that yeast infection in the mouth, right? Yeah, oral thrush.</p> <p>3. The first consideration of the cause Let's consider other possibilities. Could it be something viral? Like herpes simplex virus? That can cause little sores in the mouth too. But usually, those are more painful. My poor kid isn't complaining about pain, so maybe not that.Geographic tongue?</p> <p>4. Final decision making Alright, back to thrush. If it is thrush, we'd need to treat it with an antifungal. Nystatin suspension is usually what's used for kids. You just swab it in the mouth or mix it with some water. That should help clear it up if it's thrush..... So, I'll make an appointment with the pediatrician. They'll know exactly what to do and can confirm if it's thrush or something else. We'll go from there with their advice. That way, my grandson can feel better soon.</p> | <p>Deepseek-R1</p> <p>1. Analyze background information Lingering cold symptoms like a runny nose and cough after several weeks can certainly be worrying,</p> <p>2. Analysis and diagnosis of possible causes Those small white bumps could be a few different things: Temporary papillitis of the tongue (small white or red bumps) after a viral cold; Thrush (creamy white patches) caused by antibiotic use; Lumps or thickened tongue coating caused by postnasal drip and leakage irritation; Hand-foot-mouth disease, which is less common (usually accompanied by other symptoms).</p> <p>3. Suggestion making They need to determine the exact cause to recommend the right treatment. I cannot prescribe medication without examining him. If it's thrush, they might prescribe an antifungal mouth rinse.....</p> |

Figure 8: Case of Consulting Question

| ID | Detail |
|------------|--|
| Prompt I | <p>You are a medical expert skilled in multiple choice question answering with thorough reasoning. <Query>:[Question]</Query></p> <p>If you need to specify an option, always use option A/B/C/D. Now, conclude your answer [Therefore, my answer is: option A/B/C/D.] after thorough reasoning:</p> |
| Prompt II | <p>You are a text processor skill in splitting text according to the option. Split the text into prefix and suffix, the prefix is the content the first time the option is explained in the text. Do not alter the original text. <Text>[Text to be segmented]</Text></p> <p><Option>[Incorrect option letter]</Option></p> <p><Prefix></p> |
| Prompt III | <p>You are a professional doctor. I will provide you with a question and part of the answer. Please read the answer and extract the relevant entities that appear in the answer. Extract professional core terms with high clinical medical characteristics from the perspective of professional doctors, such as etiology, disease, diagnosis, medical examination and testing, medicine, medical therapy and other medical entities. Find all the medical entities that you think suitable in the answer. If you think there are no suitable entities, just ok to return null. Please provide it to me in JSON format, such as [entity type:xxx,entity name:xxx]. No explanation is needed, just give me JSON in English.</p> |
| Prompt IV | <p>Here is a medical query from your patient: <Query>:[question]</Query></p> <p>Here is the answer template for the medical query: <Answer template>:[temp text]</Answer template></p> <p>Please complete the mask section based on the template I provided. The case format for the mask section is as follows: < mask, type: TYPE>, The TYPE is a hint I gave you, implying the entity type in the mask section. Each part of the mask can only be filled with one answer, which can be a word or phrase. You need to find all the mask entities in this sentence and give me the answer you generated, with <mask, type: TYPE, answer: YOURANSWER> is returned without explanation. The TYPE must be the same as the TYPE I gave you. Just give me the return format. Your must give me a answer. And it cannot be the same as the TYPE. Now you can answer here:</p> |
| Prompt V | <p>Here is a medical query from your patient: <Query>:[query]</Query></p> <p>Here is your response for the medical query: <Response>:[sentences]</Response></p> <p>(<correct> [your correct preceding content] </correct>,<incorrect>[your wrong answer] </incorrect>)</p> <p>The initial answer [retry answer] in your response is incorrect, so you need to ask a reflective question based on your correct preceding content (if any) and your wrong answer.</p> <p>Now please provide a brief question(Strictly follow this format:<Reflective Question>your response</Reflective Question>):</p> |

Table 5: Prompts for Data Generation (Part 1)

| ID | Detail |
|------------|---|
| Prompt VI | <p>Here is a medical query from your patient: <Query>:[query]</Query> Here is your response for the medical query: <Response>:[sentences]</Response> Here is your own reflective question for your response: <Reflective Question>:[Reflective Question]</Reflective Question> (<correct> [your correct preceding content] </correct>,<incorrect>[your wrong answer] </incorrect>) Now, please provide a concise answer for the reflective question(Strictly follow this format:<Reflective Answer>your response</Reflective Answer>):</p> |
| Prompt VII | <p>Here is a medical query from your patient: <Query>:[query]</Query> Here is your own reflection on your initial answer: <Self-Reflection>:[reflect question][reflect answer]</Self-Reflection> Here is the response you need to complete(Complete each blank): <Response>:[sentences]</Response> Now, according on this reflection, your completed answer is(Strictly follow this format:<Answer>:your refine entity,eg:[weight reduction,cancer]</Answer>):</p> |

Table 6: Prompts for Data Generation (Part 2 - Continued)

| | |
|--------|---|
| Prompt | <p>You are an experienced doctor. Please answer the following medical questions. Conclude your response with ‘Therefore, my answer is **your option letter**’. For example, if the your answer is option A, you should conclude your response with ‘Therefore, my answer is **A**’.</p> <p><Query>:[Question]</Query> <Internal Thinking>[reflection_info]< /Internal Thinking> [Warning: The <Internal Thinking> represents your internal thoughts about the <Question>, it’s not always correct, it’s only for your reference.]</p> |
|--------|---|

Table 9: Prompts for the experiments in Section 5.6.2. Depending on the experimental setting, reflection_info may include both the reflection question and answer, only the reflection question, or only the reflection answer.