

# Unlocking Latent Medical Reasoning in LLMs via Inference-Time Representation and Prefix Interventions

Anonymous ACL submission

## Abstract

Recent reasoning advances in large language models (LLMs) have broadened their applicability to medical tasks. Yet most prior work remains dependent on scarce, high-quality rationales and compute-intensive post-training, with limited exploration of how to leverage the medical capabilities acquired during pretraining. Consequently, a key challenge is how to elicit these latent capabilities in a data-efficient manner. To address this gap, our study introduces **RIPT**, a lightweight framework for data-efficient capability activation. RIPT explicitly decomposes the objective into two complementary components: reasoning enhancement and medical knowledge elicitation. For the former, we extract steering vectors from hidden activations on a small set of high-quality paired reasoning/direct responses to shape LLMs’ reasoning behavior. For the latter, we obtain prefix vectors via prefix tuning on simple medical QA pairs to elicit domain-specific knowledge. At inference, we freeze the backbone LLM and apply a hybrid intervention that jointly injects both steering and prefix vectors. Experiments under limited-resource settings show that RIPT consistently outperforms strong baselines, suggesting an efficient pathway for unlocking LLMs’ medical reasoning capabilities.

## 1 Introduction

Recent advances in reasoning LLMs, such as OpenAI’s o1 (Jaech et al., 2024), DeepSeek-R1 (Guo et al., 2025), and medical variants like Baichuan-M2 (Dou et al., 2025), have demonstrated strong performance. The Superficial Alignment Hypothesis (Zhou et al., 2023) suggests that such reasoning abilities are largely acquired during pre-training rather than alignment, implying that LLMs exposed to large-scale medical corpora may already possess latent medical reasoning capabilities. Nevertheless, how to elicit these capabilities in a data-efficient manner remains a practical challenge.

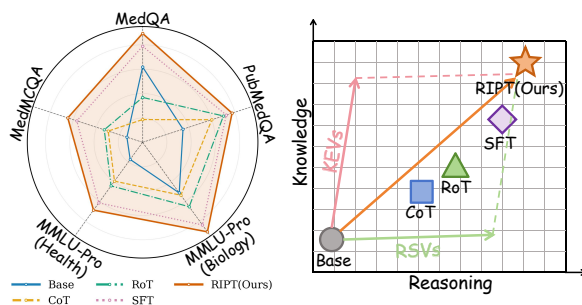


Figure 1: Overall performance of RIPT. **Left:** The performance on benchmarks after normalization. **Right:** The performance on Reasoning and Knowledge dimension. **KEVs** represent Knowledge Elicitation Vectors, while **RSVs** represent Reasoning Steering Vectors.

This challenge is particularly pronounced in medicine, where high-quality supervision is intrinsically scarce. Clinician-authored rationales are costly and subject to privacy and governance constraints (Malin et al., 2018), with data scarcity further exacerbated for long-tail and rare diseases (Kolekar et al., 2024). While post-training paradigms, such as supervised fine-tuning (SFT) (Singhal et al., 2023) and reinforcement learning (RL) (Gu et al., 2025), can be effective, they typically demand substantial computational resources and access to high-quality rationales. In contrast, prompting-based methods (Shi et al., 2024; Wei et al., 2022) are lightweight but often suffer from brittleness and instability (Chen et al., 2025). Together, these limitations motivate a central question: *Can latent medical reasoning be effectively unlocked at inference time in a data-efficient manner?*

Prior work has indicated that medical reasoning requires synergizing logical inference under uncertainty with specialized domain knowledge grounding (Berger et al., 2025). Inspired by this insight, our study proposes **Representation Injection with Prefix Tuning (RIPT)**, a novel framework that decouples this capability into two complementary

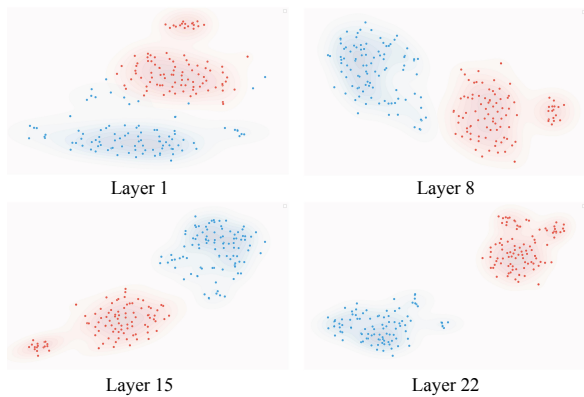


Figure 2: T-SNE visualization of reasoning (red) and direct (blue) response representations across different layers on Qwen2.5-7B-Instruct. The two clusters become more separable at deeper layers.

objectives: reasoning enhancement and medical knowledge elicitation. For reasoning enhancement, we adopt a representation engineering perspective (Zou et al., 2023). As visualized in Figure 2, reasoning patterns form separable clusters in the latent space, enabling us to extract **Reasoning Steering Vectors** from contrastive activations on a small reasoning-dense set to steer LLMs. For knowledge elicitation, we obtain **Knowledge Elicitation Vectors** via prefix tuning on simple medical QA pairs. Prior work has shown that prefix tuning can elicit domain-specific knowledge by modulating the attention mechanism while mitigating unintended degradation (Petrov et al., 2023; Lobo et al., 2025). At inference, both Reasoning Steering and Knowledge Elicitation Vectors are jointly injected into a frozen backbone, enabling RIPT to unlock latent capabilities in a data-efficient manner.

In our experiments, RIPT requires only 100 curated reasoning examples and 3,000 easily accessible QA pairs, demonstrating its data efficiency. As shown in Figure 1, RIPT yields consistent improvements (~2-3% accuracy gains) across benchmarks, with gains along both reasoning and knowledge that validate the effectiveness of our decomposed activation strategy. Detailed analyses further reveal that reasoning capabilities transfer across domains while knowledge elicitation requires adaptive, question-specific activation. Overall, our contributions can be summarized as follows:

- We propose RIPT, a data-efficient framework that effectively decouples reasoning enhancement and medical knowledge elicitation to unlock latent medical reasoning, without updating the backbone parameters.

- To our knowledge, RIPT is the first to unify activation steering and prefix vectors in an inference-time pipeline, injecting both Reasoning Steering Vectors (via contrastive representation engineering) and Knowledge Elicitation Vectors (via lightweight prefix tuning).
- Using only limited medical data, RIPT yields consistent accuracy gains across medical benchmarks. Furthermore, our analyses corroborate the complementary roles of reasoning enhancement and domain-specific knowledge elicitation.

## 2 Related Works

### 2.1 Medical Reasoning Enhancement

Methods for improving LLMs’ medical reasoning can be broadly categorized into prompting-based and post-training approaches. Prompting-based methods, such as CoT prompting (Liu et al., 2024a) and Med-PaLM’s ensemble refinement (Singhal et al., 2023), elicit reasoning through carefully designed prompts but often suffer from brittleness and sensitivity to prompt variations. Post-training paradigms align models via SFT (e.g., Med-PaLM2 (Singhal et al., 2025), UltraMedical (Zhang et al., 2024)) or RL with distilled rationales (e.g., HuatuoGPT-o1 (Chen et al., 2024), Baichuan-M2 (Dou et al., 2025)), yet typically require scarce high-quality annotations and substantial computational resources. Therefore, how to effectively unlock latent medical reasoning at inference time in a data-efficient manner remains an open challenge.

### 2.2 Representation Engineering

Representation Engineering (RepE) enhances LLM controllability via direct intervention on internal representations (Zou et al., 2023; Nanda et al., 2023; Liu et al., 2024b), showing promise in instruction following (Stolfo et al., 2024), personality steering (Cao et al., 2024), and hallucination mitigation (Arditi et al., 2024). Although recent studies have applied RepE to mathematical reasoning (Tang et al., 2025; Li et al., 2025), its potential in the medical domain remains unexplored. Crucially, unlike mathematical reasoning which relies primarily on logical deduction, medical reasoning necessitates both logical inference and specialized domain knowledge. Existing methods, however, target reasoning enhancement alone, thereby overlooking this dual requirement.

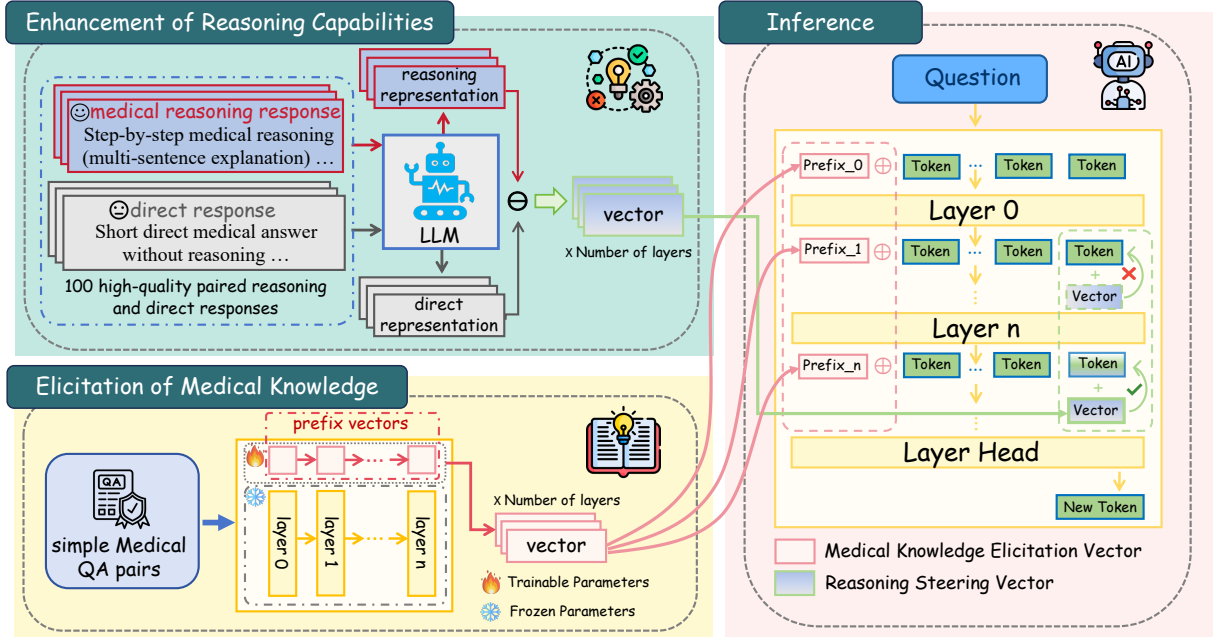


Figure 3: **Overview of the RIPT framework.** (a) **Enhancement of Reasoning Capabilities:** Reasoning steering vectors are extracted from contrastive activations between a small set of high-quality paired reasoning and direct responses. (b) **Elicitation of Medical Knowledge:** Knowledge elicitation vectors are learned via prefix tuning on simple, readily available medical QA pairs to activate internalized domain knowledge. (c) **Inference:** During generation, knowledge elicitation vectors are prepended to each layer, and reasoning steering vectors are injected into specific layers, synergizing logical inference with domain knowledge to unlock latent medical reasoning in LLMs while keeping the backbone frozen.

### 3 Methods

In this section, we propose RIPT, a lightweight framework designed to unlock latent medical reasoning in a data-efficient manner. Specifically, it decouples medical reasoning into two complementary objectives: reasoning enhancement and medical knowledge elicitation.

We first formalize the task and inference process (§3.1), then analyze the representations and detail the extraction and injection of steering vectors (§3.2). Finally, we present the prefix tuning approach for knowledge elicitation (§3.3).

#### 3.1 Task Formulation

We formulate medical reasoning as a conditional text generation task. Given a medical query  $q$ , the model generates a response  $a$  that synthesizes both logical reasoning and domain-specific medical knowledge, which can be expressed as modeling the conditional distribution  $p(a|q)$ .

We build our approach upon the auto-regressive Transformer architecture (Vaswani et al., 2017), which underlies state-of-the-art LLMs. The model processes the input via a stack of  $L$  decoder layers, each comprising a Multi-Head Self-Attention

(MHSA) mechanism and a Feed-Forward Network (FFN) connected via residual streams. The hidden state  $h_t^l \in \mathbb{R}^d$  at layer  $l$  for token  $t$  is computed via the residual stream as:

$$h_t^l = h_t^{l-1} + a_t^l + m_t^l \quad (1)$$

where  $a_t^l$  and  $m_t^l$  are the outputs of the MHSA and FFN at layer  $l$  for token  $t$ , respectively. The model then predicts the next token auto-regressively based on the final layer’s hidden state  $h_t^L$ .

#### 3.2 Enhancement of Reasoning Capabilities

In this section, we detail the enhancement of reasoning capabilities. Our method proceeds in three stages: (1) extracting contrastive representations from a small set of paired reasoning/direct responses; (2) analyzing the geometric separability of these representations in the latent space; and (3) deriving steering vectors to enhance reasoning behaviors of the frozen backbone during inference.

**Extraction of Representations.** To extract representations, we construct a small set of high-quality medical queries  $\mathcal{Q} = \{q_i\}_{i=1}^N$  together with paired direct and reasoning responses  $(d_i, r_i)$ . For each query  $q_i$ , we construct two input sequences:

197  $x_i^d = [q_i; d_i]$  and  $x_i^r = [q_i; r_i]$ , where  $[\cdot; \cdot]$  denotes  
 198 concatenation. We then perform forward passes on  
 199 these sequences and extract hidden states. Follow-  
 200 ing standard practice for auto-regressive models,  
 201 we take the representation from layer  $l$  at the final  
 202 response-token position:

$$203 \quad R^l(d_i) = \mathbf{h}^l(x_i^d); \quad R^l(r_i) = \mathbf{h}^l(x_i^r). \quad (2)$$

204 Here,  $\mathbf{h}^l(x)$  denotes the hidden state at layer  $l$  of  
 205 the last token in sequence  $x$ .

206 **Analysis of Representations.** To analyze the  
 207 latent structure of reasoning and direct repre-  
 208 sentations, we apply t-SNE (Maaten and Hin-  
 209 ton, 2008) to project the extracted hidden states  
 210  $\{R^l(d_i), R^l(r_i)\}$  onto a 2D space. As shown in  
 211 Figure 2, representations corresponding to reason-  
 212 ing responses form a compact cluster within a spe-  
 213 cific region of the latent space. Notably, this clus-  
 214 ter is clearly separated from the region occupied  
 215 by direct responses, with the separation becoming  
 216 more pronounced in deeper layers. This reveals  
 217 that reasoning patterns occupy a distinct region in  
 218 the latent space, enabling reasoning enhancement  
 219 through directional steering from direct to reason-  
 220 ing representations.

221 **Reasoning Steering Vectors.** The observed separa-  
 222 tion implies the existence of a directional shift in  
 223 the latent space from direct to reasoning generation.  
 224 We leverage this geometric separability to extract  
 225 steering vectors that capture this shift. To abstract  
 226 away instance-specific variations and capture the  
 227 general reasoning direction, we derive the **Reason-**  
 228 **ing Steering Vectors**  $\mathbf{p}_r^l$  at layer  $l$  by averaging the  
 229 contrastive differences across all queries in  $\mathcal{Q}$ :

$$230 \quad \mathbf{p}_r^l = \frac{1}{|\mathcal{Q}|} \sum_{i=1}^{|\mathcal{Q}|} (R^l(r_i) - R^l(d_i)).$$

231 We inject  $\mathbf{p}_r^l$  into the hidden state of the last input  
 232 token during inference, as this position encodes  
 233 the full context and directly influences the gener-  
 234 ation of subsequent reasoning tokens. We inject  
 235 the steering vector with a controllable coefficient  
 236  $\lambda_p$  and then re-normalize to preserve the original  
 237 magnitude:

$$238 \quad \tilde{\mathbf{h}}_{\text{last}}^l = \mathbf{h}_{\text{last}}^l + \lambda_p \cdot \mathbf{p}_r^l, \quad \tilde{\mathbf{h}}_{\text{last}}^l \leftarrow \tilde{\mathbf{h}}_{\text{last}}^l \cdot \frac{\|\mathbf{h}_{\text{last}}^l\|_2}{\|\tilde{\mathbf{h}}_{\text{last}}^l\|_2}.$$

### 239 3.3 Elicitation of Medical Knowledge

240 Prior work indicates that prefix tuning functions  
 241 by modulating the attention mechanism to guide  
 242 models to attend to specific internal representa-  
 243 tions (Petrov et al., 2023). Motivated by this, we  
 244 employ prefix tuning as a targeted activation mech-  
 245 anism—not to inject new information, but to elicit  
 246 the model’s internalized medical knowledge. By  
 247 keeping the backbone parameters frozen, this ap-  
 248 proach avoids the substantial high-quality data re-  
 249 quirements and potential degradation of general  
 250 reasoning capabilities often associated with full-  
 251 scale medical fine-tuning (Lobo et al., 2025).

252 Specifically, we introduce a set of **Knowledge**  
 253 **Elicitation Vectors**  $\mathbf{p}_k \in \mathbb{R}^{k \times d}$ , where  $k$  is the  
 254 prefix length and  $d$  is the hidden dimension. These  
 255 layer-specific vectors are prepended to the input  
 256 representations at each layer and optimized on sim-  
 257 ple, readily available medical QA pairs while keep-  
 258 ing all base model parameters frozen. The training  
 259 objective minimizes the language modeling loss:

$$260 \quad \mathcal{L}(\mathbf{p}_k) = \mathbb{E}_{(q,a) \sim \mathcal{D}} [\mathcal{L}_{\text{LM}}(f([\mathbf{p}_k; q]; \theta), a)] \quad (3)$$

261 where  $\mathcal{D}$  is the medical QA corpus,  $(q, a)$  is a  
 262 query-answer pair, and  $[\mathbf{p}_k; q]$  indicates concatena-  
 263 tion of the prefix with input embeddings. The func-  
 264 tion  $f(\cdot; \theta)$  represents the frozen pretrained model  
 265 with parameters  $\theta$ . During inference, prepending  
 266  $\mathbf{p}_k$  to the input at each layer elicits domain knowl-  
 267 edge without modifying the model’s parameters.

### 268 3.4 Overall Framework

269 **RIPT** employs a hybrid intervention strategy that  
 270 synergizes reasoning steering vectors and knowl-  
 271 edge elicitation vectors to unlock latent medical  
 272 reasoning in a data-efficient manner. The complete  
 273 inference process can be formalized as:

$$274 \quad a = f_{\theta}([\mathbf{p}_k; q]; \{\mathbf{h}^l + \lambda_p \cdot \mathbf{p}_r^l\}_{l \in \mathcal{L}}) \quad (4)$$

275 where knowledge elicitation vectors  $\mathbf{p}_k$  are  
 276 prepended to the input query  $q$ , and reasoning steer-  
 277 ing vectors  $\mathbf{p}_r^l$  are injected into hidden states at  
 278 target layers  $\mathcal{L}$  with controllable strength  $\lambda_p$ , while  
 279 the base model parameters  $\theta$  remain frozen.

## 280 4 Experiments

### 281 4.1 Experimental Setup

282 **Data Construction.** We construct two datasets  
 283 for RIPT’s components. For reasoning steering vec-  
 284 tors, we leverage HuatuoGPT-o1 (Chen et al., 2024)

Model	Method	MedQA	PubMedQA	MMLU-Pro		MedMCQA	Overall
				Biology	Health		
Qwen2.5-7B -Instruct	Vanilla	63.7	72.8	73.1	55.9	55.7	64.2
	+CoT	61.1 (-2.6)	74.8 (+2.0)	73.2 (+0.1)	57.0 (+1.1)	56.5 (+0.8)	64.5 (+0.3)
	+RoT	62.2 (-1.5)	75.6 (+2.8)	73.9 (+0.8)	57.2 (+1.3)	56.7 (+1.0)	65.1 (+0.9)
	+SFT	<u>64.7</u> (+1.0)	<u>75.8</u> (+3.0)	<u>75.0</u> (+1.9)	<u>58.1</u> (+2.2)	<u>57.9</u> (+2.2)	<u>66.3</u> (+2.1)
	<b>+RIPT</b>	<b>65.4</b> (+1.7)	<b>76.2</b> (+3.4)	<b>75.5</b> (+2.4)	<b>58.4</b> (+2.5)	<b>58.2</b> (+2.5)	<b>66.7</b> (+2.5)
Llama3.1-8B -Instruct	Vanilla	66.7	73.8	66.5	54.7	58.9	64.1
	+CoT	67.2 (+0.5)	78.6 (+4.8)	67.4 (+0.9)	56.5 (+1.8)	60.1 (+1.2)	66.0 (+1.9)
	+RoT	67.7 (+1.0)	77.8 (+4.0)	67.9 (+1.4)	56.9 (+2.2)	60.3 (+1.4)	66.1 (+2.0)
	+SFT	<u>68.4</u> (+1.7)	<u>77.4</u> (+3.6)	<u>69.2</u> (+2.7)	<b>59.3</b> (+4.6)	<u>61.2</u> (+2.3)	<u>67.1</u> (+3.0)
	<b>+RIPT</b>	<b>68.7</b> (+2.0)	<b>78.8</b> (+5.0)	<b>69.9</b> (+3.4)	<u>58.7</u> (+4.0)	<b>61.5</b> (+2.6)	<b>67.5</b> (+3.4)
UltraMedical -8B	Vanilla	<u>70.6</u>	76.8	63.2	55.0	55.4	64.2
	+CoT	67.8 (-2.8)	73.8 (-3.0)	64.0 (+0.8)	53.3 (-1.7)	53.6 (-1.8)	62.5 (-1.7)
	+RoT	69.8 (-0.8)	<u>77.0</u> (+0.2)	<u>64.7</u> (+1.5)	<u>55.5</u> (+0.5)	54.9 (-0.5)	<u>64.4</u> (+0.2)
	<b>+RIPT</b>	<b>72.1</b> (+1.5)	<b>79.2</b> (+2.4)	<b>67.0</b> (+3.8)	<b>57.7</b> (+2.7)	<b>55.9</b> (+0.5)	<b>66.4</b> (+2.2)

Table 1: **Main results on medical benchmarks.** Accuracy (%) of three models with different methods across four medical benchmarks: MedQA, PubMedQA, MMLU-Pro (Biology and Health), and MedMCQA. Numbers in parentheses indicate improvement over vanilla baseline. The best and second-best results are **bolded** and underlined.

reasoning dataset, which provides medical queries with paired reasoning and direct responses. We employ GPT-4o (Hurst et al., 2024) and Gemini-2.5-flash (Comanici et al., 2025) to score each sample’s reasoning quality under five dimensions (detailed scoring rules in Appendix A). Samples with high scores are retained, from which we randomly select 100 examples. For knowledge elicitation, we use 3,000 samples from the training set of MedMCQA (Pal et al., 2022), a knowledge-intensive multiple-choice dataset, reformatted into direct QA pairs.

**Models.** We evaluate RIPT on three representative models: Qwen2.5-7B-Instruct (Team et al., 2024) and Llama-3.1-8B-Instruct (Dubey et al., 2024), two widely-used open-source general models, and UltraMedical-8B (Zhang et al., 2024), a medical model fine-tuned on domain-specific data.

**Baselines.** We compare RIPT with representative baselines across three paradigms: CoT prompting (Wei et al., 2022) as the prompting-based method, Representation of Thought (RoT) (Hu et al., 2024) as the representation engineering approach, and SFT as the post-training method. Detailed configurations are provided in Appendix B.2.

**Benchmarks.** We evaluate RIPT on four medical benchmarks. **MedQA** (Jin et al., 2021) consists of

USMLE-style questions from medical licensing exams. **MedMCQA** (Pal et al., 2022) covers diverse medical topics across 21 specialties from entrance examinations. **PubMedQA** (Jin et al., 2019) requires answering questions based on biomedical research abstracts. **MMLU-Pro** (Wang et al., 2024) (Biology and Health) contains challenging questions with up to ten choices. Together, these benchmarks evaluate different aspects of medical competence, including clinical reasoning, domain knowledge, and scientific literature comprehension.

## 4.2 Experimental Results

**Comparison with Baselines.** As shown in Table 1, RIPT consistently outperforms all baseline methods. Compared to CoT prompting, RIPT achieves an average improvement of 2.2%, demonstrating that representation-level intervention provides more stable and effective control over reasoning processes than prompting alone. RIPT also outperforms RoT by an average of 1.6%, validating the effectiveness of our decoupled strategy that separately targets reasoning enhancement and knowledge elicitation, rather than relying on a single intervention mechanism. Furthermore, RIPT achieves competitive or superior performance compared to SFT across most benchmarks, while requiring less than 0.3% trainable parameters. This highlights

RIPT’s parameter efficiency without sacrificing performance, confirming that latent medical reasoning capabilities can be effectively unlocked through lightweight representation-level interventions.

**Generalization Across Benchmarks.** RIPT exhibits strong generalization across benchmarks with diverse formats and complexity levels. It delivers consistent gains on standard medical licensing exams (MedQA, MedMCQA), which simulate clinical decision-making scenarios. Furthermore, these improvements extend to context-rich and challenging specialized tasks, such as PubMedQA and the MMLU-Pro subsets. This consistent improvement across benchmarks varying in format, difficulty, and required skills suggests that RIPT enhances fundamental medical reasoning capabilities rather than exploiting dataset-specific artifacts.

**Effectiveness Across Models.** RIPT demonstrates robust efficacy across models with different specializations. For general-purpose LLMs, RIPT yields substantial gains, indicating that these models possess internalized medical capabilities that can be effectively unlocked via representation-level intervention. RIPT also benefits the medical-specialized UltraMedical-8B, confirming its applicability beyond general domains. Notably, prompting-based methods hurt performance on certain benchmarks for UltraMedical-8B, confirming that such approaches tend to be brittle and unstable. In contrast, RIPT yields consistent improvements across all benchmarks, demonstrating that our decoupled strategy offers more precise and stable control over model behavior.

## 5 Analysis and Discussion

### 5.1 Ablation Study

To quantify the contribution of each component in our hybrid intervention strategy, we conduct ablation studies by selectively removing reasoning steering vectors or knowledge elicitation vectors. Table 2 presents results on MedQA and MMLU-Pro (Biology and Health). The results indicate that latent medical reasoning capabilities exist, which can be separately activated. Each component independently improves performance: reasoning steering enhances multi-step inference, while knowledge elicitation activates domain-specific knowledge. The combined strategy outperforms either component alone, confirming that these two capabilities are complementary.

Model	Vectors		MedQA	MMLU-Pro	
	Reas.	Know.		Bio.	Hlth.
Qwen2.5-7B -Instruct	-	-	63.7	73.1	55.9
	✓	-	64.1	74.9	57.8
	-	✓	64.9	74.4	56.8
	✓	✓	<b>65.4</b>	<b>75.5</b>	<b>58.4</b>
UltraMedical- 8B	-	-	70.6	63.2	55.0
	✓	-	71.2	66.2	57.1
	-	✓	71.5	63.0	55.3
	✓	✓	<b>72.1</b>	<b>67.0</b>	<b>57.7</b>

Table 2: Ablation study on MedQA and MMLU-Pro datasets. **Reas.:** Reasoning Steering Vectors; **Know.:** Medical Knowledge Elicitation Vectors.

The relative impact of the two components differs across models. For Qwen2.5-7B-Instruct, both components yield substantial improvements, indicating that general models possess latent reasoning capabilities and internalized knowledge that can be effectively activated. In contrast, UltraMedical-8B benefits more from reasoning steering than from knowledge elicitation, suggesting that domain-specific training has already strengthened medical knowledge access, leaving reasoning enhancement as the primary bottleneck. These findings validate our decoupled strategy, supporting the view that RIPT unlocks existing capabilities rather than injecting new ones through extensive post-training.

### 5.2 Decoupled Quality Assessment

To separately assess reasoning and knowledge quality, we conduct a fine-grained evaluation on 200 sampled MedQA instances using GPT-4o (Hurst et al., 2024) as an automatic judge. Each response is scored on both dimensions using a 0-5 scale, with detailed rubrics provided in Appendix C.1.

As shown in Table 3, different methods exhibit distinct strengths. CoT primarily boosts reasoning while offering marginal improvement in knowledge, indicating that prompting can enhance reasoning processes but fails to elicit domain knowledge. SFT delivers strong reasoning gains but only moderate knowledge improvement, suggesting it mainly refines reasoning style rather than substantially activating internalized knowledge. In contrast, RIPT achieves the superior scores on both dimensions, delivering the most significant knowledge enhancement while maintaining reasoning performance. This confirms the efficacy of our decoupled approach in effectively unlocking both latent reasoning and medical knowledge.

Method	Reasoning		Knowledge	
	Score	$\Delta$	Score	$\Delta$
Qwen2.5-7B-Instruct	3.84	–	3.67	–
+ CoT	4.17	+0.33	3.71	+0.04
+ RoT	4.23	+0.39	3.75	+0.08
+ SFT	4.32	+0.48	3.80	+0.13
<b>RIPT (Ours)</b>	<b>4.36</b>	<b>+0.52</b>	<b>3.95</b>	<b>+0.28</b>

Table 3: LLM-as-judge evaluation on reasoning and knowledge dimensions. Scores range from 0-5.

Method	MedQA		MedMCQA	
	Acc.	$\Delta$	Acc.	$\Delta$
Qwen2.5-7B-Instruct	63.7	–	55.7	–
+ LoRA	58.0	-5.7	53.3	-2.4
+ Domain Steering	64.1	+0.4	54.9	-0.8
+ Prefix Tuning	64.9	+1.2	57.1	+1.4

Table 4: Comparison of different methods for knowledge elicitation under data-constrained setting on MedQA and MedMCQA.

### 5.3 Cross-Domain Reasoning Transfer

To investigate whether the reasoning enhancement component captures domain-agnostic capabilities, we conduct a cross-domain transfer experiment. Specifically, we extract reasoning steering vectors from 100 mathematical reasoning examples in STILL-2 (Min et al., 2024), using the same configuration as utilized for medical reasoning, and apply them to medical benchmarks.

As shown in Figure 4, steering vectors derived purely from mathematical tasks improve performance on medical benchmarks, demonstrating clear cross-domain transferability. The fact that vectors completely lacking medical semantics can still enhance medical reasoning validates our hypothesis: reasoning is a domain-agnostic capability that is distinct from parametric knowledge.

Nevertheless, medical-derived vectors maintain a performance advantage over math-derived ones on both benchmarks. This indicates that while general reasoning skills are transferable, in-domain vectors more effectively capture task-specific reasoning patterns required for optimal performance.

### 5.4 Rationale for Prefix Tuning in Knowledge Elicitation

To determine the most suitable strategy for data-efficient knowledge elicitation, we compare three

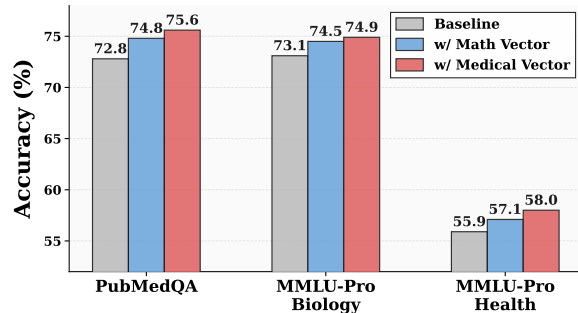


Figure 4: Cross-domain transfer of reasoning steering vectors. Math Vector: Derived from mathematical reasoning/direct QA pairs. Medical Vector: Derived from medical reasoning/direct QA pairs.

lightweight approaches: LoRA (Hu et al., 2022), Domain Steering (Tang et al., 2025) and Prefix Tuning (Li and Liang, 2021). LoRA and prefix tuning are trained on simple medical QA pairs, whereas domain steering derives fixed vectors from medical QA representations.

As shown in Table 4, LoRA induces significant performance degradation. We attribute this to LoRA directly modifying the model’s weights, causing it to memorize the simple patterns in our QA pairs and compromise general capabilities. In contrast, prefix tuning mitigates this risk by keeping the backbone frozen, and instead guides the model’s attention toward medical knowledge, making it effective for our data-constrained setting.

Domain steering yields marginal or even negative gains across benchmarks. Figure 6 visualizes the last-layer hidden states of the first decoded token. While reasoning steering exhibits a consistent direction, knowledge elicitation displays a more scattered radial pattern. This indicates that knowledge elicitation is inherently question-specific, making a single fixed vector insufficient to generalize across diverse medical queries.

### 5.5 Analysis of Hyperparameters

To investigate the optimal configuration for reasoning steering vectors, we analyze the effect of injection layer and strength on two models. A subset of 200 MedMCQA examples is sampled as a validation set for this analysis.

As shown in Figure 7 (top), the middle-to-late layers yield the best results, while the early layers show limited improvement and the final layers cause performance degradation. This pattern holds across both models, suggesting that reasoning capabilities are closely associated with these

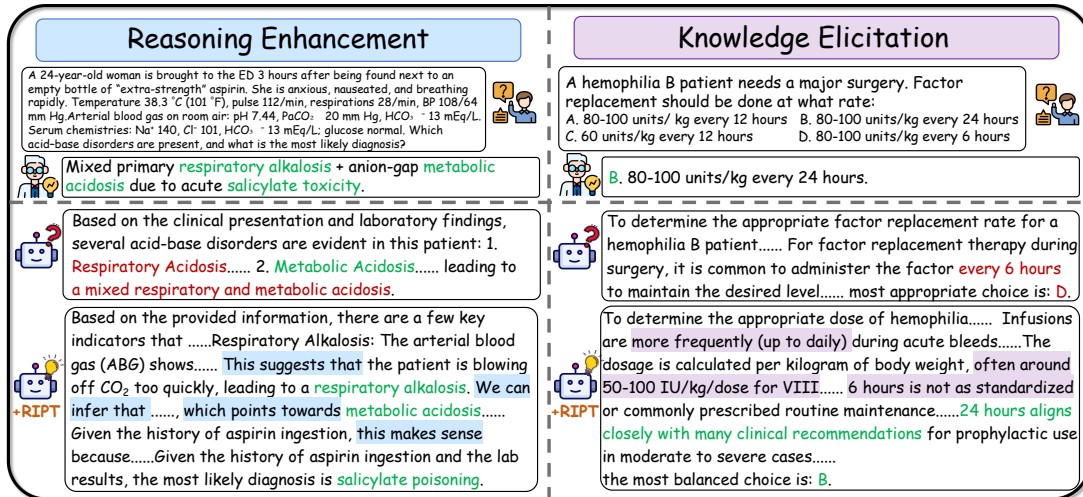


Figure 5: Case study of Qwen2.5-7B-Instruct with RIPT.

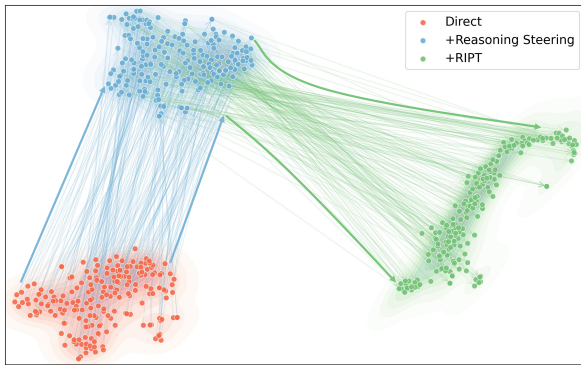


Figure 6: T-SNE visualization of hidden state shifts during RIPT. Reasoning steering produces consistent directional shifts (blue arrows), while knowledge elicitation introduces question-specific diversity (green arrows).

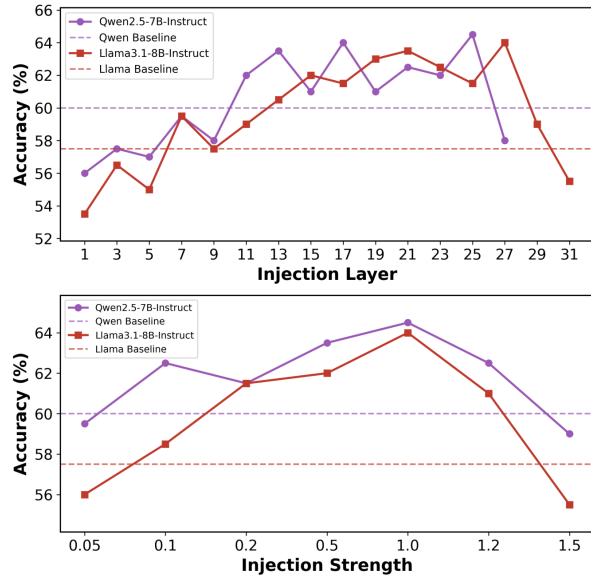


Figure 7: Effect of injection layer (top) and injection strength (bottom) on a subset of MedMCQA.

intermediate layers across different architectures. Accordingly, we select layer 25 for Qwen2.5-7B-Instruct and layer 27 for Llama-3.1-8B-Instruct.

Figure 7 (bottom) shows that performance peaks at  $\lambda_p = 1.0$ , as insufficient strength fails to elicit reasoning capabilities while excessive strength disrupts the representations. Overall, RIPT consistently outperforms the baseline across a wide range of layers and strengths, demonstrating robustness to hyperparameter choices.

## 5.6 Case Study

Figure 5 demonstrates how RIPT unlocks latent medical reasoning capabilities through our decoupled approach: enhanced reasoning enables structured multi-step inference that the base model fails to produce, while knowledge elicitation activates domain-specific clinical knowledge already internalized in the model.

## 6 Conclusion

In this paper, we propose RIPT, a lightweight framework designed to unlock LLMs' latent medical reasoning capabilities acquired during pre-training. By decoupling medical reasoning into reasoning enhancement and knowledge elicitation, our approach jointly injects two vectors at inference time, achieving data-efficient capability enhancement without updating backbone parameters. Extensive experiments demonstrate the strong performance of RIPT, validating the effectiveness of our hybrid intervention strategy. Overall, our work provides a precise and stable pathway to unlock latent capabilities, contributing to the adaptation of general LLMs to specialized domains.

## 7 Limitations

One limitation of our work is that RIPT requires access to internal model representations, making it infeasible for closed-source LLMs. Furthermore, due to computational constraints, we only conduct experiments on three representative models and four medical benchmarks. Our evaluation focuses on public benchmarks rather than clinical deployment due to cost constraints.

## References

Andy Ardit, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in language models is mediated by a single direction. *Advances in Neural Information Processing Systems*, 37:136037–136083.

Armin Berger, Sarthak Khanna, David Berghaus, and Rafet Sifa. 2025. Reasoning llms in the medical domain: A literature survey. *arXiv preprint arXiv:2508.19097*.

Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin, Lu Lin, Fenglong Ma, and Jinghui Chen. 2024. Personalized steering of large language models: Versatile steering vectors through bi-directional preference optimization. *Advances in Neural Information Processing Systems*, 37:49519–49551.

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

Shan Chen, Mingye Gao, Kuleen Sasse, Thomas Hartvigsen, Brian Anthony, Lizhou Fan, Hugo Aerts, Jack Gallifant, and Danielle S Bitterman. 2025. When helpfulness backfires: LLMs and the risk of false medical information due to sycophantic behavior. *npj Digital Medicine*, 8(1):605.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Chengfeng Dou, Chong Liu, Fan Yang, Fei Li, Jiyuan Jia, Mingyang Chen, Qiang Ju, Shuai Wang, Shunya Dang, Tianpeng Li, and 1 others. 2025. Baichuanm2: Scaling medical capability with large verifier system. *arXiv preprint arXiv:2509.02208*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Boyang Gu, Hongjian Zhou, Bradley Max Segal, Jinge Wu, Zeyu Cao, Hantao Zhong, Lei Clifton, Fenglin Liu, and David A Clifton. 2025. Clinical-r1: Empowering large language models for faithful and comprehensive reasoning with clinical objective relative policy optimization. *arXiv preprint arXiv:2512.00601*.

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

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.

Lijie Hu, Liang Liu, Shu Yang, Xin Chen, Zhen Tan, Muhammad Asif Ali, Mengdi Li, and Di Wang. 2024. Understanding reasoning in chain-of-thought from the hopfieldian view. *arXiv preprint arXiv:2410.03595*.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 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.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, pages 2567–2577.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.

Shreesh S Kolekar, Monali Gulhane, Swapna Kamble, Pragati Patil Bedekar, Hemchandra V Nerlekar, and Dinesh Goyal. 2024. AI systems for diagnosing rare diseases: Challenges and solutions. In *Proceedings of the 6th International Conference on Information Management & Machine Intelligence*, pages 1–11.

Qiming Li, Xiaocheng Feng, Yixuan Ma, Zekai Ye, Ruihan Chen, Xiachong Feng, and Bing Qin. 2025.



741 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer,  
742 Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping  
743 Yu, Lili Yu, and 1 others. 2023. Lima: Less is more  
744 for alignment. *Advances in Neural Information Pro-*  
745 *cessing Systems*, 36:55006–55021.

746 Andy Zou, Long Phan, Sarah Chen, James Campbell,  
747 Phillip Guo, Richard Ren, Alexander Pan, Xuwang  
748 Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,  
749 and 1 others. 2023. Representation engineering: A  
750 top-down approach to ai transparency. *arXiv preprint*  
751 *arXiv:2310.01405*.

752	<b>A Data Construction</b>		
753	<b>Reasoning Data Selection Pipeline.</b>	To curate a	
754		high-quality subset from the HuatuoGPT-o1 (Chen	
755		et al., 2024) reasoning dataset, we implement a	
756		rigorous filtering pipeline using GPT-4o (Hurst	
757		et al., 2024) and Gemini-2.5-Flash (Comanici et al.,	
758		2025) for cross-validation. We employ rubric-	
759		-based evaluation prompts (refer to Appendix D.1)	
760		to assess the quality of each <i>Complex_CoT</i> ratio-	
761		-nale across five dimensions: <i>Faithfulness</i> , <i>Medi-</i>	
762		<i>cal_Correctness</i> , <i>Reasoning_Clarity</i> , <i>Differential</i> ,	
763		and <i>Actionability</i> . We retain only instances that	
764		achieve a composite score of $\geq 23/25$ from both	
765		models. From this filtered pool, we randomly sam-	
766		-ple 100 pairs to constitute the final dataset.	
767	<b>Knowledge Elicitation Data Construction.</b>	We	
768		derive the knowledge elicitation training data by	
769		transforming a subset of MedMCQA (Pal et al.,	
770		2022) from multiple-choice format into direct QA	
771		pairs. Specifically, we process 10,000 examples	
772		from the training split using GPT-4o (Hurst et al.,	
773		2024) to assess whether each question can be re-	
774		-formulated as a standalone factual query without	
775		relying on the provided options. For convertible	
776		instances, the model generates a QA-style rewrite	
777		with the original ground-truth option as the answer.	
778		From all successfully converted instances, we uni-	
779		-formly sample 3,000 QA pairs to form the final	
780		knowledge elicitation corpus.	
781	<b>B Baselines and Training Configuration</b>		
782	<b>B.1 Benchmark Details</b>		
783		We select four diverse benchmarks that test differ-	
784		-ent facets of medical competency, ranging from	
785		knowledge recall to multi-step clinical reasoning.	
786	<b>MedQA (USMLE) (Jin et al., 2021):</b>	Derived	
787		from the USMLE, this dataset assesses clinical	
788		decision-making through complex patient scenar-	
789		-ios. We utilize the English test set, which requires	
790		applying medical knowledge to perform differen-	
791		-tial diagnoses.	
792	<b>MedMCQA (Pal et al., 2022):</b>	A large-scale	
793		dataset covering 21 distinct medical subjects (e.g.,	
794		Surgery, Anatomy) from entrance examinations. Its	
795		wide topical coverage serves as a robust test of the	
796		model’s <i>domain knowledge breadth</i> .	
797	<b>PubMedQA (Jin et al., 2019):</b>	A biomedical	
798		question-answering task based on research ab-	
799		-stracts. Unlike factoid QA, it requires answering	
		"Yes/No/Maybe" based strictly on the provided con-	800
		-text, emphasizing <i>logical inference</i> and evidence-	801
		-based reasoning over memorized knowledge.	802
	<b>MMLU-Pro (Health &amp; Biology) (Wang et al.,</b>		803
	<b>2024):</b>	A rigorous benchmark designed to chal-	804
		-lenge reasoning robustness by increasing the option	805
		-count (up to 10 choices) and removing retrieval	806
		-shortcuts. We focus on the Health and Biology	807
		-subsets to evaluate <i>complex reasoning</i> under high-	808
		-difficulty conditions.	809
	<b>B.2 Baseline implementation details</b>		810
		In this section, we provide a comprehensive de-	811
		-scription of the baseline methods employed in our	812
		-comparative analysis. These baselines are strategi-	813
		-cally selected to represent three distinct paradigms	814
		-of enhancing model reasoning: <i>Prompting-based</i>	815
		- <i>(CoT (Kojima et al., 2022))</i> , <i>Inference-time Repre-</i>	816
		- <i>sentation Intervention (RoT (Hu et al., 2024))</i> , and	817
		- <i>Post Training (SFT)</i> .	818
	<ul style="list-style-type: none"><li>• <b>CoT:</b></li></ul>	This method evaluates the intrinsic reason-	819
		-ing capability of the pre-trained model without	820
		-any parameter updates or internal interventions.	821
		We employ the standard Zero-shot CoT prompt-	822
		-ing strategy by appending the trigger phrase	823
		-" <i>Answer the following question step by step.</i> "	824
		-to the input query. This prompts the model to	825
		-generate intermediate reasoning steps prior to	826
		-deriving the final answer.	827
	<ul style="list-style-type: none"><li>• <b>RoT:</b></li></ul>	This method extracts contrastive represen-	828
		-tations based on whether a CoT prompt or a	829
		-non-CoT prompt is included in the input, and	830
		-then injects them into the model’s latent space.	831
	<ul style="list-style-type: none"><li>• <b>SFT:</b></li></ul>	This baseline represents the standard	832
		-gradient-based learning approach. We fine-tune	833
		-the backbone model using the curated high-	834
		-quality medical reasoning dataset described in	835
		-Appendix A. We summarize the SFT hyperpa-	836
		-rameters in Table 6.	837
	<b>B.3 Computational Resources</b>		838
		All experiments were conducted on NVIDIA	839
		-A100 GPUs (80GB). RIPT is computationally	840
		-lightweight, requiring only a single GPU for both	841
		-steering vector extraction and prefix tuning. Infer-	842
		-ence adds negligible overhead as it only involves	843
		-vector addition and prefix concatenation without	844
		-additional forward passes.	845

Hyperparameters	Value
Optimizer	AdamW
Learning rate	$5 \times 10^{-4}$
Batch size (per GPU)	2
Epochs	1
Prefix length (virtual tokens)	16

Table 5: Prefix-tuning hyperparameters.

Hyperparameters	Value
Optimizer	AdamW
Learning rate	$2 \times 10^{-4}$
Batch size per GPU	2
Epochs	5
Gradient accumulation	8
Effective batch size	$2 \times 8 = 16$

Table 6: SFT hyperparameters.

Hyperparameters	Value
Learning Rate	$5 \times 10^{-5}$
Batch Size	16
Epochs	5
LoRA Rank	4
LoRA Alpha	16
LoRA Dropout	0.15
Target Modules	q/k/v/o proj

Table 7: LoRA fine-tuning configuration.

Setting	Qwen	LLaMA
Layers	28	32
Token Position	Last token (Final Layer)	
Perplexity	30	
Iterations	1000	
Initialization	PCA	
Output Dim.	2	

Table 8: t-SNE hyperparameters for representation visualization.

## B.4 Prefix Tuning Configuration

Table 5 details the specific hyperparameters used for prefix tuning in our framework.

## C Additional Analyses

### C.1 Decoupled Quality Assessment

We assess model performance along two dimensions: *Reasoning* and *Knowledge*, each scored on a 0-5 scale. We employ GPT-4o (Hurst et al., 2024) as an LLM-based judge (Zheng et al., 2023) with rubric-constrained prompts (see Appendix D.2 and D.3). The judge is instructed to score based solely on explicit response content, without inferring unstated merits.

For an evaluation dataset consisting of  $N$  samples, we report the mean Reasoning score ( $\bar{s}^{(R)}$ ) and the mean Knowledge score ( $\bar{s}^{(K)}$ ), defined as:

$$\bar{s}^{(R)} = \frac{1}{N} \sum_{i=1}^N s_i^{(R)}, \quad \text{and} \quad \bar{s}^{(K)} = \frac{1}{N} \sum_{i=1}^N s_i^{(K)}, \quad (5)$$

where  $s_i^{(R)}, s_i^{(K)} \in \{0, 1, \dots, 5\}$  denote the integer scores assigned to the  $i$ -th sample. The judge’s outputs are parsed strictly as JSON objects; in cases of parsing failure, a default score of 0 is assigned, and the failure rate is reported.

### C.2 LoRA Configuration

To compare our representation-level intervention with parameter-efficient fine-tuning methods, we implement LoRA (Hu et al., 2022). LoRA modifies

the model behavior by optimizing low-rank decomposition matrices within the Transformer layers. The LoRA hyperparameters are in Table 7.

### C.3 Domain Steering Implementation

Following (Tang et al., 2025), we compute a static **Domain Steering** vector by averaging hidden states of the final token across medical QA data. Unlike RIPT’s dynamic approach, this method extracts a single global direction. During inference, this fixed vector is injected into hidden states at the target layer with coefficient  $\lambda_p$  to steer activations toward the medical subspace.

### C.4 t-SNE Visualization Settings

We use t-SNE (Maaten and Hinton, 2008) to visualize layer-wise hidden states. Dimensionality reduction is performed independently for each layer to project representations into 2D space. Hyperparameter settings are detailed in Table 8.

890

## D Prompt Templates

891

### D.1 Quality Filtering Prompt

**System Role:** You are a reviewer for a medical reasoning dataset.

**Task:** Evaluate the quality of the Complex\_CoT using only the Question and the Complex\_CoT. Do not introduce facts beyond the stem. Do not add extra explanations. Output JSON exactly as requested, with no additional text.

**Evaluation Criteria:**

Please score the Complex\_CoT against the following minimal criteria and decide whether to keep it:

- **Faithfulness:** strictly based on the stem; no fabricated key data (0-5)
- **Medical\_Correctness:** correct concepts/mechanisms/logic; no obvious errors (0-5)
- **Reasoning\_Clarity:** clear causal/step-by-step chain (0-5)
- **Differential:** considers key differentials and life-threatening but reversible conditions with exclusion rationale (0-5)
- **Actionability:** provides a minimal verification or follow-up plan under uncertainty (0-5)

**Scoring Guide:** 0 = missing/incorrect; 3 = adequate; 5 = excellent. Total = sum of five items (0-25).

**Keep Rule:** Total  $\geq$  23 and no obvious medical error  $\rightarrow$  keep = true; otherwise keep = false.

**Output Format:** Output ONLY this JSON (no extra text):

```
{
  "keep": true,
  "overall": 0,
  "scores": {
    "Faithfulness": 0,
    "Medical_Correctness": 0,
    "Reasoning_Clarity": 0,
    "Differential": 0,
    "Actionability": 0
  },
  "strength": "One-sentence key strength",
  "improvement": "One-sentence top improvement"
}
```

892

893

### D.2 Reasoning Judge Prompt

**System Role:** You are a medical education expert specializing in the assessment of medical reasoning skills.

**Task:** Evaluate the model's response solely from the perspective of reasoning quality, not the medical correctness of the final conclusion.

**Important Notes:**

- Do NOT evaluate whether the final diagnosis is correct.
- Focus on whether the reasoning process aligns with clinical thinking.
- The score MUST be based only on reasoning that is explicitly stated in the response; do NOT make charitable or implicit inferences.

**Scoring Criteria (0-5):**

**5 (Excellent):**

- The reasoning is logically rigorous, with clear and complete steps.
- Key symptoms and signs are explicitly linked to underlying pathophysiological mechanisms.

894

- Differential diagnoses are systematically considered, with clear explanations for ruling out alternatives.
- The reasoning pathway reflects authentic clinical decision-making.
- The response repeatedly explains “why,” not just “what.”

#### 4 (Good):

- The reasoning is generally clear and follows a coherent analytical approach.
- Major symptoms are connected to diagnostic hypotheses.
- Differential diagnoses are mentioned, but depth or completeness is slightly insufficient.
- The reasoning steps are mostly complete, though some explanations lack detail.

#### 3 (Moderate):

- Some reasoning is present but lacks systematic structure.
- A conclusion is given, but the reasoning chain is incomplete or partially interrupted.
- Causal analysis is superficial, with limited explanation of “why.”
- Differential diagnosis is minimal or largely formulaic.

#### 2 (Poor):

- The reasoning shows clear logical jumps and omits key intermediate steps.
- The link between symptoms and conclusions is weak.
- There are obvious logical flaws or unexplained assumptions.
- Little to no evidence of differential diagnostic thinking.

#### 1 (Very Poor):

- The reasoning is disorganized or internally inconsistent.
- There is almost no logical connection between symptoms and conclusions.
- Basic clinical reasoning frameworks are absent.

#### 0 (Invalid):

- No reasoning process is present.
- The response is irrelevant, refuses to answer, or provides only a conclusion without reasoning.

**Output Format:** Output strictly in the following JSON format. Do NOT add any additional text or fields:

```
{
  "score": <integer from 0 to 5>,
  "reasoning": "<Explain the strengths and weaknesses of
               the reasoning quality, explicitly
               referencing content in the response>"
}
```

895

### D.3 Knowledge Judge Prompt

896

**System Role:** You are a medical knowledge assessment expert specializing in evaluating the accuracy and completeness of medical knowledge.

**Task:** Evaluate the model’s response solely from the perspective of medical knowledge quality, and NOT to assess reasoning processes, presentation style, or language fluency.

**Important Notes:**

- Do NOT evaluate whether the reasoning logic is sound.
- Do NOT evaluate clarity, organization, or persuasiveness of the response.

897

- The score MUST be based only on medical knowledge that is explicitly stated in the response; do NOT make charitable assumptions or fill in missing information.

**Scoring Criteria (0-5):**

**5 (Excellent):**

- All medical facts stated in the response are fully accurate, with no explicit or implicit errors.
- Medical terminology is used correctly, precisely, and in an appropriate clinical context.
- Knowledge coverage is comprehensive and includes all key points relevant to the question.
- The content is consistent with current mainstream clinical guidelines or evidence-based medicine (as can be determined directly from the response).
- Critical details are accurately described (e.g., dosages, reference ranges, indications/contraindications).

**4 (Good):**

- The main medical facts are accurate.
- Terminology is generally appropriate and does not cause substantive misunderstanding.
- Most key knowledge points are covered, with some omissions.
- The content aligns with commonly used clinical standards but does not reflect the latest advances or guideline-level details.
- Details are generally accurate but incomplete.

**3 (Moderate):**

- Core medical knowledge points are correct, but coverage is limited.
- Some terminology is imprecise or ambiguous.
- Several important pieces of information are missing, reducing knowledge completeness.
- The content is broadly consistent with general medical knowledge but lacks depth or systematic coverage.

**2 (Poor):**

- The response contains clear medical knowledge errors.
- Key information is missing or inaccurately described.
- Terminology is used inconsistently or incorrectly, potentially leading to misunderstanding.
- The content deviates substantially from commonly accepted clinical standards.

**1 (Very Poor):**

- Multiple major medical errors are present.
- Knowledge is severely incomplete, clearly outdated, or highly misleading.
- The response has little to no educational or clinical reference value.

**0 (Invalid):**

- The response is entirely incorrect or clearly irrelevant to the question.
- The response includes dangerous or inappropriate medical advice that could cause patient harm.

**Output Format:** Output strictly in the following JSON format. Do NOT add any additional text or fields:

```
{
  "score": <integer from 0 to 5>,
  "reasoning": "<Explain the strengths and weaknesses of
               the medical knowledge in terms of accuracy
               and completeness, explicitly referencing
               the response content>"
}
```

899 **E Ethical Considerations and Limitations**

900 **E.1 Data Usage and Compliance**

901 We utilize publicly available datasets and models  
902 in strict adherence to their respective licenses and  
903 comply with the terms of service of all third-party  
904 APIs used. All resources are employed solely for  
905 academic research and benchmarking purposes.

906 **E.2 Risks and Safety**

907 Although RIPT enhances reasoning capabilities,  
908 it does not guarantee factual correctness and may  
909 generate fluent but erroneous medical explanations.  
910 Such outputs must not be treated as medical advice.  
911 We emphasize that RIPT is not intended for clinical  
912 deployment or real-world decision-making without  
913 expert oversight.