

# Med-VRAgent: A Framework for Medical Visual Reasoning-Enhanced Agents

Guangfu Guo<sup>1,2\*</sup>      Xiaoqian Lu<sup>1,2\*</sup>      Yue Feng<sup>1</sup>  
fs23990@bristol.ac.uk      fd23007@bristol.ac.uk      y.feng.6@bham.ac.uk  
<sup>1</sup> University of Birmingham    <sup>2</sup> University of Bristol

## Abstract

Visual Language Models (VLMs) achieve promising results in medical reasoning but struggle with hallucinations, vague descriptions, inconsistent logic and poor localization. To address this, we propose a agent framework named Medical Visual Reasoning Agent (**Med-VRAgent**). The approach is based on Visual Guidance and Self-Reward paradigms and Monte Carlo Tree Search (MCTS). By combining the Visual Guidance with tree search, Med-VRAgent improves the medical visual reasoning capabilities of VLMs. We use the trajectories collected by Med-VRAgent as feedback to further improve the performance by fine-tuning the VLMs with the proximal policy optimization (PPO) objective. Experiments on multiple medical VQA benchmarks demonstrate that our method outperforms existing approaches. Our implementation is publicly available <https://github.com/KwongFuk/Med-VRAgent>.

## 1 Introduction

Visual Language Models (VLMs) enable context-aware medical reasoning and have shown strong performance in tasks like radiology report generation (Hartsock and Rasool, 2024; Tanno et al., 2025; Li et al., 2024). However, they remain prone to hallucinations, where outputs deviate from the visual input—posing risks in clinical settings (Chen et al., 2025; Jin et al., 2024). This issue is exacerbated by the factual unreliability of underlying large language models (LLMs) (Huang et al., 2025; Zhu et al., 2024b; Pal et al., 2023), highlighting the urgent need for effective mitigation strategies (Kim et al., 2025; Bai et al., 2025).

Researchers also have explored several enhancements, the Chain-of-Thought (CoT) has become a popular approach to enhance the logical reasoning capability (Wei et al., 2023). Visual prompting—using region-specific cues—have

\*Equal contribution. This work was completed during a research assistant at the University of Birmingham.

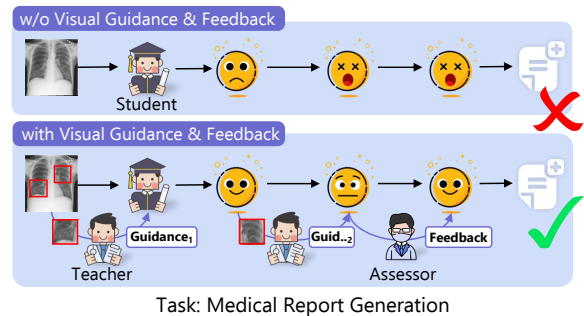


Figure 1: Top: A student struggles, feeling confused and making mistakes. Bottom: With guidance, the student overcomes the confusion and successfully completes the task.

been shown to improve model performance in fields such as radiology and pathology where precise localization is required. (Denner et al., 2025). Plan-then-Generate decouple reasoning into structured planning followed by execution (Zhou et al., 2023). Self-enhancement mechanisms, such as self-reflection, self-correction, and self-critique, and external feedback systems enable models to revise their own reasoning (Madaan et al., 2023; Wang et al., 2024). Additionally, Retrieval-Augmented Generation (RAG) incorporates external knowledge to support the reasoning process (Lewis et al., 2021).

While the above approaches are effective, some key challenges remain. (1) In high-stakes domains like radiology and pathology, VLMs often lack fine-grained image-text alignment, producing overly generic descriptions that miss critical local details, spatial structures, and abnormal patterns. (Vishwanath et al., 2025; Liévin et al., 2023). (2) Although complex medical prompting strategies have been proposed to address this issue, they are often domain-specific, labor-intensive, and require expert knowledge. (Boiko et al., 2023; Xia et al., 2024a). (3) Furthermore, current models, even with visual prompting, focus on a single ROI and struggle to integrate overall medical image structure and spatially distributed lesions, limiting performance in

cases with high spatial complexity. (Wang et al., 2025; Huang et al., 2024b). (4) Current frameworks offer limited feedback, usually evaluating only the final output, making error detection and correction during reasoning difficult. (5) Finally, retrieval enhancement methods often introduce irrelevant or noisy information, potentially distorting clinical reasoning. (Gao et al., 2024; Ji et al., 2023).

We propose a multimodal agent framework **Med-VRAgent**, to tackle challenges like error propagation, suboptimal planning, limited feedback, and the fragility of retrieval-based methods. Med-VRAgent consists of three core modules—**Teacher**, **Student**, and **Assessor**—and two key components: a **Visual Extraction Module**, and a **Retrieval-Augmented Reflection (RAR)**. The Visual Extraction Module identifies Regions of interest (ROIs) in medical images and uses Visual Token Edit to improve the agent’s regional perception. The Teacher provides ROI-specific visual guidance. The Student generate diagnostic outputs with ROI and teacher’s guidance. The Assessor offers fine-grained feedback for iterative refinement. We use RAR module to enhance factual grounding by incorporating external medical knowledge and introduce Monte Carlo Tree Search (MCTS) to explore high-quality reasoning paths using an adaptive expansion strategy while better balancing performance and efficiency. Our framework only needs to be trained once for both the teacher and the assessor, which can achieve good generalization ability and save computational resources.

Results across three benchmarks confirm the superior performance of Med-VRAgent, achieving new state-of-the-art (SOTA) results. It outperforms reasoning baselines (Visual CoT) on GMAI (Table 3), surpasses retrieval-augmented methods on IU-Xray (Demner-Fushman et al., 2016) (Table 4), and exceeds advanced fine-tuning strategies like MMedPO on VQA-RAD (Lau et al., 2018) and MIMIC-CXR (Johnson et al., 2019) (Table 2). These results highlight the effectiveness of our visual guidance-based medical multimodal agent framework.

In summary, our contributions are as follows:

- We propose a **Teacher-Student-Evaluator** framework for medical visual reasoning based on **Visual Guidance** and **Feedback**.
- We use **Visual Extraction** and **Visual Token Edit** to improve the visual capabilities of multimodal agents.

- We develop a **Retrieval-Augmented Reflection** module to further boost agent reasoning via External knowledge.
- Extensive experiments on multiple medical multimodal benchmarks demonstrate that our framework achieves SOTA performance.

## 2 Related Work

### 2.1 Foundation Large models

Large Language Models (LLMs) like GPT-3 (Brown et al., 2020), PaLM (Chowdhery et al., 2022), and LLaMA (Touvron et al., 2023) have shown strong capabilities in reasoning, generation, and understanding across natural language tasks, excelling in few-shot learning, in-context reasoning, and text generation. These models are central to the development of multi-modal systems. VLMs have demonstrated remarkable generalization across cross-modal tasks such as image captioning, retrieval, and visual question answering (VQA). Early models like CLIP (Radford et al., 2021) and Flamingo (Alayrac et al., 2022) use large-scale image-text pairs for contrastive or retrieval-based learning. Recent models like BLIP-2 (Li et al., 2023b) and MiniGPT-4 (Zhu et al., 2023) integrate LLMs with visual encoders to enhance reasoning and support open-ended question answering. These advances in Foundation Large Models (FLMs) lay the foundation for tasks that require deep cross-modal understanding.

### 2.2 Multi-step Reasoning in FLMs

Reasoning in Foundation Large Models (FLMs) has advanced with frameworks enhancing multi-step inference and decision-making. CoT (Wei et al., 2023) enables intermediate reasoning steps, improving performance on complex tasks. ToT (Yao et al., 2023) explores multiple reasoning paths using tree search strategies, boosting decision-making. The ReAct framework (Yao et al., 2022) combines reasoning with environment interaction, improving tool-augmented tasks. In multimodal reasoning, Visual Chain-of-Thought (Rose et al., 2024) extends CoT by integrating visual grounding to bridge logical gaps. The Reinforced Ranker-Reader ( $R^3$ ) architecture (Zhang et al., 2023) improves open-domain question answering by combining a ranker and reader with reinforcement learning, optimizing accuracy over retrieved documents.

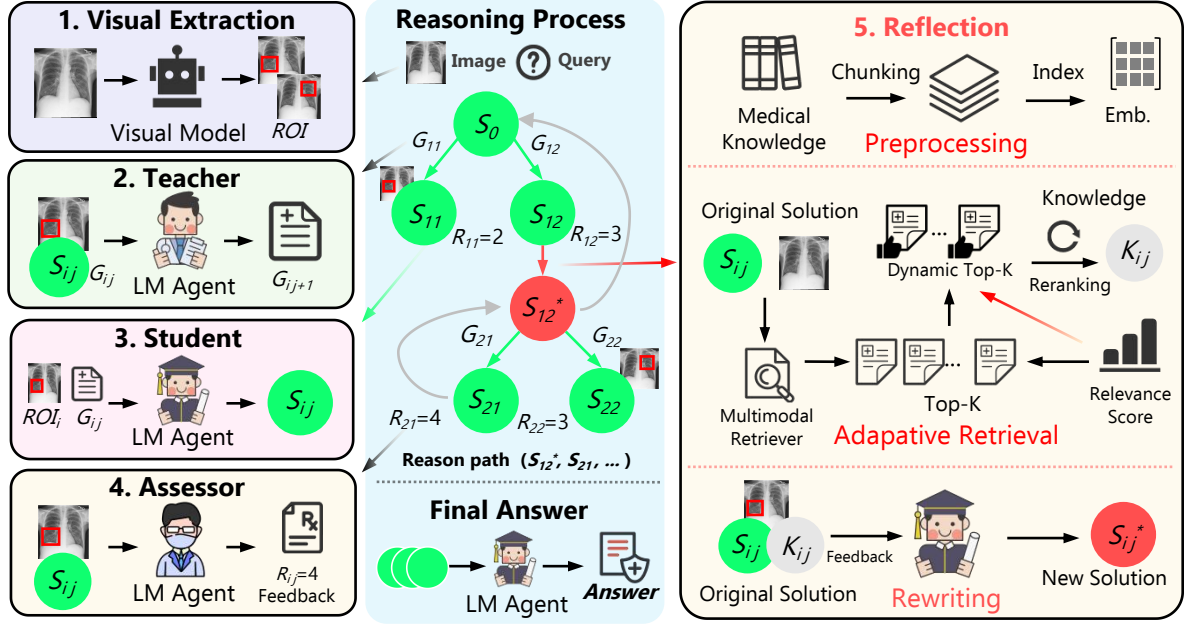


Figure 2: Overview of the **Med-VRAgent** framework. The system uses MCTS to generate solutions  $\mathcal{S}_{ij}$  based on Regions of Interest  $\mathcal{ROI}_{ij}$ , visual guides  $\mathcal{G}_{ij}$ , rewards  $\mathcal{R}_{ij}$ , and external knowledge  $\mathcal{K}_{ij}$ .  $\mathcal{S}_{ij}^*$  is the solution after reflection.

### 2.3 Medical-Specific Reasoning Frameworks

Medical-Specific Reasoning has advanced with specialized frameworks to enhance LLMs’ clinical reasoning. MedAgents (Tang et al., 2024) creates a multi-agent system where LLM-based medical experts collaborate on diagnostic tasks, improving zero-shot reasoning. MedReason (Wu et al., 2025) aligns LLM reasoning with medical graphs, enhancing decision-making accuracy and interpretability. FineMedLM-o1 (Yu et al., 2025) uses supervised fine-tuning and test-time training on curated dialogues for complex tasks like differential diagnosis. DeepSeek R1 (Moëll et al., 2025) benchmarks LLM outputs against expert behavior, revealing both advanced reasoning and domain-specific biases. These models highlight the value of tailored frameworks and medical knowledge in improving LLM clinical reasoning.

## 3 Methodology

To enhance medical visual reasoning, we propose Med-VRAgent, a novel reasoning scheme. It combines a Visual guidance and Reward-Feedback Paradigm in a search algorithm to optimize reasoning paths.

### 3.1 Med-VRAgent Reasoning Process

Fig 2 illustrates the Med-VRAgent process. We model the agent reasoning process as a tree search, where each node  $\mathcal{S}_{ij}$  represents a state defined as:

$$\mathcal{S}_{ij} = [\mathcal{Q}, \mathcal{I}, \mathcal{G}_{ij}, \mathcal{A}_{ij}, \mathcal{R}_{ij}, \mathcal{F}_{ij}, \mathcal{A}_{ij}^*, \mathcal{O}_{ij}, \mathcal{ROI}_i] \quad (1)$$

where  $\mathcal{Q}$  is the query,  $\mathcal{I}$  is the medical image,  $\mathcal{G}_{ij}$  is the visual guidance,  $\mathcal{A}_{ij}$  is the current step answer,  $\mathcal{R}_{ij}$  is the reward,  $\mathcal{F}_{ij}$  is feedback,  $\mathcal{A}_{ij}^*$  is the answer after reflection,  $\mathcal{O}_{ij}$  represents the observation information, including all ancestor and sibling node guidance and answers, and  $\mathcal{ROI}_i$  is the visual ROI.

Given an image  $\mathcal{I}$  and query  $\mathcal{Q}$ , the goal is for Student  $\mathcal{S}_{\text{model}}$  to generate step-by-step reasoning using ROIs  $\mathcal{ROI}_i$  from Vision Extraction  $\mathcal{V}_{\text{model}}$  and visual guidance  $\mathcal{G}_{ij}$  from Teacher  $\mathcal{T}_{\text{model}}^\theta$ . Assessor  $\mathcal{A}_{\text{model}}^\theta$  evaluates guidance and answers, providing reward  $\mathcal{R}_{ij}$  and feedback  $\mathcal{F}_{ij}$ . If answer quality is low, the reflection module uses external knowledge  $\mathcal{K}$  from retriever  $\mathcal{R}_{\text{model}}$  to refine it. MCTS searches for the optimal reasoning path for answering  $\mathcal{Q}$ .

### 3.2 Visual Extraction Module

**Visual ROIs Extraction** We use a lightweight VLM to extract medical entities  $\mathcal{E}$  relevant to the question and image. Following MedVP (Zhu et al., 2025), we adopt a fine-tuned Grounding DINO (Liu et al., 2024) as the visual extractor. Grounding DINO is an open-vocabulary detector that localizes entities from image  $\mathcal{I}$  and text prompts  $\mathcal{E} = \{e_1, e_2, \dots, e_N\}$ .

$$ROI = \{ROI_i\}_{i=1}^N = \text{G-DI}(I, E), \quad ROI_i = (b_i, s_i, l_i) \quad (2)$$

$ROI$  is the set of extracted regions, with each  $ROI_i = (b_i, s_i, l_i)$  representing the bounding box, confidence score, and matched entity label.

**Visual Token Edit** To refine the Agent’s focus on a ROI without retraining, we apply Visual Token Edit (VTE), a single edit to visual tokens in the first ( $K \leq 3$ ) self-attention layers. For each patch embedding  $\mathbf{v}_i \in \mathbb{R}^d$  and binary ROI mask  $m_i \in \{0, 1\}$ , we replace:

$$\mathbf{v}_i \longrightarrow \mathbf{v}_i^* = \mathbf{v}_i + \beta m_i \mathbf{b} \quad (3)$$

where  $\mathbf{b}$  is a fixed direction (e.g.,  $\mathbf{1}$  or  $\mathbf{v}_i$ ). Because the key and value projections are linear, Eq. (3) increases the  $\ell_2$  norm of ROI tokens and thus raises their soft-max attention weights *indirectly*, concentrating information flow on the referenced region while keeping background tokens intact.

The gain  $\beta > 0$  is chosen on-the-fly to prevent over- or under-boosting:

$$\beta = s_i \kappa \phi\left(\frac{\bar{a}_{\text{bg}}}{\bar{a}_{\text{ROI}}} - 1\right), \quad \kappa \in [0, 1], \quad (4)$$

where  $\bar{a}_{\text{ROI}}$  and  $\bar{a}_{\text{bg}}$  are the average pre-softmax attention logits of ROI and background patches obtained from a provisional forward pass, and  $s_i$  is the detector confidence for the ROI,  $\phi(\cdot)$  is any element-wise activation that is non-negative and monotonically non-decreasing. When the model already attends to the ROI ( $\bar{a}_{\text{ROI}} \geq \bar{a}_{\text{bg}}$ ), Eq. (4) yields  $\beta = 0$ , leaving the representation unchanged. Setting  $\kappa = 0$  disables VTE entirely, making the mechanism safe, computationally negligible, and fully reversible.

### 3.3 Teacher-Student-Assessor Mechanism

**Teacher Agent.** In natural language tasks, the exponential growth of tag combinations severely limits vanilla MCTS. To improve efficiency, we incorporate a prompt-driven Teacher  $\mathcal{T}_{\text{model}}^\theta$  that expands the policy space via heuristics. See the appendix ?? for prompt. At each node,  $\mathcal{T}_{\text{model}}^\theta$  gathers prior guidance–answer pairs  $(\mathcal{G}_{1..i}, \mathcal{A}_{1..i})$  and feedback  $\mathcal{F}$ , then generates the next-step guidance:

$$\mathcal{G}_{ij+1} = \mathcal{T}_{\text{model}}(\mathcal{ROI}_i, \mathcal{G}_{i1..j}, \mathcal{A}_{i1..j}, \mathcal{F}_i) \quad (5)$$

**Student Agent.** The Student  $\mathcal{S}_{\text{model}}$  leverages a vision-language backbone to perform step-wise reasoning. At each stage of problem, it receives the

Teacher  $\mathcal{T}_{\text{model}}^\theta$ -generated guidance  $\mathcal{G}_{ij}$  and the corresponding image  $\mathcal{ROI}_i$ , and produces an intermediate answer  $\mathcal{A}_{ij}$ . After search, the best reasoning path selected by MCTS is used to compose the final answer. This process is formally defined as:

$$\mathcal{A}_{ij} = \mathcal{S}_{\text{model}}(\mathcal{ROI}_i, \mathcal{G}_{ij}), \quad (6)$$

**Assessor Agent.** In the MCTS, it is essential to quantitatively evaluate each reasoning step and provide high-quality feedback to guide the search process. To this end, we adopt a *LLM-as-a-Judge* (Gu et al., 2025) approach, we introduce an Assessor model  $\mathcal{A}_{\text{model}}^\theta$ , implemented using a VLM, and grounded in the *Self-Rewarding* paradigm (Yuan et al., 2025) The Assessor  $\mathcal{A}_{\text{model}}^\theta$  employs a 5-point scoring system to evaluate task progress, where the score reflects both the quality and contribution of each intermediate answer. The Assessor  $\mathcal{A}_{\text{model}}^\theta$  receives the image ROI  $\mathcal{ROI}_i$ , the current guidance  $\mathcal{G}_{ij}$ , and the student’s answer  $\mathcal{A}_{ij}$ . It then produces both a descriptive feedback  $\mathcal{F}_{ij}$  and a quantitative rating  $\mathcal{R}_{ij}$ , See Appendix 7 for prompt. The process is formalized as:

$$\mathcal{F}_{ij}, \mathcal{R}_{ij} = \mathcal{A}_{\text{model}}(\mathcal{ROI}_i, \mathcal{G}_{ij}, \mathcal{A}_{ij}). \quad (7)$$

### 3.4 Retrieval-Augmented Reflection

The reflection phase is designed to enhance ROI analysis tasks that the Student  $\mathcal{S}_{\text{model}}$  fails to complete under Teacher  $\mathcal{T}_{\text{model}}^\theta$  guidance. We use IU-Xray, MIMIC-CXR, VQA-RAD and other datasets as knowledge sources. The reflection process consists of two stages:

**Retrieval Phase.** We adopt the domain-aware retriever from MMed-RAG (Xia et al., 2025), which uses ResNet-50 (He et al., 2015) and BioClinicalBERT (Alsentzer et al., 2019) as the image and text encoders, respectively. During reflection, the retriever takes as input the guidance  $\mathcal{G}_{ij}$ , image  $\mathcal{I}$ , and answer  $\mathcal{A}_{ij}$ . It first retrieves a Top- $\mathcal{K}$  candidate set  $\mathcal{K}_1$  from the external knowledge base using FAISS (Johnson et al., 2017). A cross-attention-based relevance scoring model cross-encoder/ms-marco-MiniLM-L-6-v2 (Reimers and Gurevych, 2019) then refines these candidates into a subset  $\mathcal{K}_2$ , which is finally reranked to produce the final knowledge set  $\mathcal{K}_{ij}$ . This multi-stage knowledge retrieval process is formally expressed as:

$$\mathcal{K}_{ij} = \text{Rerank}\left(\text{Relevance}(\text{RetrieveTop-}K(\mathcal{I}, \mathcal{G}_{ij}, \mathcal{A}_{ij}))\right) \quad (8)$$

**Rewriting Phase.** When reflection is needed, the student  $\mathcal{S}_{\text{model}}$  receives the original answer  $\mathcal{A}_{ij}$ , guidance  $\mathcal{G}_{ij}$ , the input  $\mathcal{ROL}$ , feedback  $\mathcal{F}_{ij}$ , and retrieved knowledge  $\mathcal{K}_{ij}$ . It then synthesizes these inputs to produce a refined answer  $\mathcal{A}_{ij}^*$ . This rewriting process can be formalized as:

$$\mathcal{A}_{ij}^* = \mathcal{S}_{\text{model}}(\mathcal{ROL}_i, \mathcal{G}_{ij}, \mathcal{A}_{ij}, \mathcal{F}_{ij}, \mathcal{K}_{ij}) \quad (9)$$

### 3.5 Monte Carlo Tree Search Process

Monte Carlo Tree Search (MCTS) operates through four main phases—selection, expansion, evaluation, and backpropagation—repeating until satisfactory reasoning results are produced or computational limits are reached. In the **Selection** phase, the algorithm starts at the root node (initial state  $\mathcal{S}_0$ ) and recursively selects child nodes using the Upper Confidence Bound (UCB) formula, which balances exploration and exploitation:

$$UCB(s) = \frac{R(s)}{N(s)} + C \cdot \sqrt{\frac{2 \cdot \ln N(p)}{N(s)}} \quad (10)$$

where  $R(s)$  is the reward,  $N(s)$  the visit count of node  $s$ ,  $N(p)$  the visit count of its parent  $p$ , and  $C$  is a constant. The **Expansion** phase involves selecting an unprocessed ROI along the current path and expanding it by sampling  $\mathcal{N}$  guidance suggestions from the Teacher  $\mathcal{T}_{\text{model}}^\theta$ . This step incorporates a heuristics mechanism, where feedback from Assessor  $\mathcal{A}_{\text{model}}^\theta$  and all observations—including guidance, answer from ancestor and sibling nodes are provided to the Teacher  $\mathcal{T}_{\text{model}}^\theta$ . In the **Evaluation** phase, each new child node is assessed using feedback from the Assessor  $\mathcal{A}_{\text{model}}^\theta$ . Finally, in the **Backpropagation** phase, the reward  $\mathcal{R}(s')$  is used to update the average reward and visit counts for node  $\mathcal{S}'$  and its ancestors.

To improve search performance and efficiency in MCTS, we apply some strategy.

**Early Stopping.** Expansion is terminated when the node score exceeds 4 or when KL divergence and semantic similarity suggest the Student  $\mathcal{S}_{\text{model}}$  and Teacher  $\mathcal{T}_{\text{model}}^\theta$  outputs align with the previous node. This allows the agent to shift to other ROIs.

**Alpha-Beta Pruning.** During selection and expansion, Alpha (min guaranteed by maximization) and Beta (max guaranteed by minimization) bounds are maintained. Subtrees are pruned when node scores fall outside this range, avoiding unnecessary evaluations.

**Reflection.** If early stopping is triggered repeatedly or the expansion limit is reached without

achieving a score of 4, the reflection module is activated. In this case, the Student  $\mathcal{S}_{\text{model}}$  retrieves external knowledge to continue reasoning.

### 3.6 Training Strategy and Optimization

To enhance the Teacher  $\mathcal{T}_{\text{model}}^\theta$  and Assessor  $\mathcal{A}_{\text{model}}^\theta$ , we fine-tune both VLMs using proximal policy optimization (PPO) with feedback trajectories collected by Med-VRAgent. PPO optimizes the policy by maximizing expected rewards while constraining updates to avoid performance degradation. The objective is:

$$\mathcal{L}_{\text{PPO}}(\theta) = \mathbb{E} \left[ \min \left( r_\theta \hat{A}_t, \text{clip}(r_\theta, 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (11)$$

where

$$r_\theta = \frac{\pi_\theta(A_{1..i}|O_{1..i})}{\pi_{\theta_{\text{old}}}(A_{1..i}|O_{1..i})} \quad (12)$$

Here,  $A_{1..i}$  and  $O_{1..i}$  denote sampled actions (guidance) and observations, respectively, while  $\hat{A}_t$  is the advantage estimate and  $\epsilon$  is the clipping threshold. We collect trajectories

$$\mathcal{T}_{\text{Med-VRAgent}} = (A_{1..i}, O_{1..i}, R_{1..i}) \quad (13)$$

from Med-VRAgent to estimate advantages and update the policy parameters  $\theta$ . The clipping in  $\mathcal{L}_{\text{PPO}}(\theta)$  ensures conservative, stable updates.

## 4 Experiments

### 4.1 Experimental Datasets

Dataset	Modality	Size	Task Type
IU-Xray	X-ray	590	Report Generation
MIMIC-CXR	Chest X-ray	500	Report Generation
VQA-RAD	X-ray, CT	451	Visual Question Answering
GMAI-MMbench	38 modalities	4 task	Visual Question Answering

Table 1: The medical visual datasets used in this experiment

We evaluate Med-VRAgent on various medical visual-linguistic datasets covering report generation and VQA tasks. As shown in Table 1, for report generation, we use the IU-Xray (Demner-Fushman et al., 2016) dataset containing 590 test samples and the MIMIC-CXR (Johnson et al., 2019) dataset test with 500 test samples. For VQA, we use VQA-RAD containing test 451 QA pairs based on X-rays and CT images and GMAI-MMbench (Chen et al., 2024) we use 4 clinical tasks.

For the Med-VQA task, for open questions, we report recall in the Open column. For closed questions, we report precision in the Closed column.

For the report generation task, we use BLEU (Papineni et al., 2002) Score, ROUGE-L (Lin, 2004), and METEOR as metrics (Banerjee and Lavie, 2005). BLEU score represents the average of BLEU-1/2/3/4.

## 4.2 Compared Methods

We evaluate the performance of various methods across different approaches.

For training methods, we employ the LLaVA-Med (Li et al., 2023a) model and assess its performance on the VQA-RAD and MIMIC-CXR datasets. The training approaches compared include SFT, Self-Rewarding (Yuan et al., 2025), Direct Preference Optimization (DPO) (Rafailov et al., 2024), STLLaVA-Med, and MMedPO (Zhu et al., 2024a).

For reasoning methods, we use the DeepSeek-VL-7B (Lu et al., 2024) and MiniCPM-V2 (Yao et al., 2024) models, evaluating their performance on the GMAI-MMBench. The reasoning approaches compared include CoT (Wei et al., 2023), ToT (Yao et al., 2023), and Visual CoT (Shao et al., 2024).

Finally, for Decoding-based and Retrieval-Augmented methods, we use the LLaVA-Med v1.5 model and evaluate its performance on the IU-Xray dataset. The Decoding-based methods include Greedy Decoding, BeamSearch (Xie et al., 2023), DoLa (Chuang et al., 2024), OPERA (Huang et al., 2024a), VCD (Leng et al., 2023). The RAG approaches compared include MedDr (He et al., 2024), FactMM-RAG (Sun et al., 2025), RULE (Xia et al., 2024b), and MMed-RAG (Xia et al., 2025). Please see the appendix 6 for details.

## 4.3 Model Implementation

We applied Med-VRAgent to LLaVA-Med v1.5, DeepSeek-VL-7B, and MiniCPM-V2. To ensure fair comparison, we follow the same experimental settings as prior work, using a decoding temperature of 0.7. We use DeepSeek-VL-7B as the Teacher  $\mathcal{T}_{\text{model}}^{\theta}$  and Assessor  $\mathcal{A}_{\text{model}}^{\theta}$  and perform PPO fine-tuning.

For PPO fine-tuning, We follow the official training scripts and use the "peft" and "trl" Python packages to implement LoRA and PPO. The fine-tuning process is completed within 7–8 hours on 4 Nvidia A6000 GPUs. The "lora\_target\_modules" are set to ["q\_proj", "v\_proj"], with lora\_r set to 16, lora\_alpha set to 32, and lora\_dropout set to 0.05. The micro\_batch\_size is 1, the batch\_size is

Methods	VQA-RAD		MIMIC-CXR		
	Open	Closed	BLEU	ROUGE-L	METEOR
LLaVA-Med v1.5	29.24	63.97	10.25	9.38	7.71
SFT	31.38	64.26	12.39	10.21	8.75
Self-Rewarding	32.69	65.89	12.15	10.05	8.77
DPO	32.88	64.33	12.37	10.38	9.10
STLLaVA-Med	33.72	64.70	12.21	10.12	8.98
MMedPO	34.03	67.64	13.28	13.22	<b>10.20</b>
<b>Med-VRAgent (Ours)</b>	<b>35.70</b>	<b>68.72</b>	<b>13.90</b>	<b>13.53</b>	9.58

Table 2: Comparison of Med-VRAgent with fine-tuning methods, including SFT, Self-Rewarding, DPO, STLLaVA-Med, and MMedPO, evaluated on VQA-RAD (Open/Closed Accuracy) and MIMIC-CXR (BLEU, ROUGE-L, METEOR) datasets, based on LLaVA-Med v1.5. The best result for each model is bolded.

Methods	AR	BVR	B	CR	Average
DeepSeek-VL-7B	38.43	47.03	42.31	37.03	41.20
CoT	39.24	46.60	43.26	38.18	41.57
ToT	40.23	46.07	44.42	39.58	42.08
Visual CoT	41.57	46.76	44.13	41.59	43.51
<b>Med-VRAgent (Ours)</b>	<b>44.81</b>	<b>51.82</b>	<b>47.52</b>	<b>42.79</b>	<b>46.74</b>
MiniCPM-V2	40.74	43.01	36.46	37.57	39.45
CoT	41.69	43.90	37.69	38.74	40.51
ToT	42.14	44.32	38.27	39.29	41.01
Visual CoT	43.20	44.70	39.12	41.28	42.08
<b>Med-VRAgent (Ours)</b>	<b>44.81</b>	<b>47.32</b>	<b>40.18</b>	<b>41.34</b>	<b>43.41</b>

Table 3: Comparison of Med-VRAgent with reasoning methods, including CoT, Tree-of-Thought (ToT), and Visual CoT, evaluated on the GMAI (Accuracy) dataset, based on DeepSeek-VL-7B and MiniCPM-V2. GMAI include AR (Attribute Recognition), BVR (Blood Vessels Recognition), B (Bone), and CR (Cell Recognition). The best result for each model is bolded, and average values are in blue.

8, and num\_epochs is 1. For optimization, we set the learning\_rate to 1.41e-5, the reward baseline to 3.75, and the random seed to 0.

## 4.4 Overall Performance

**Evaluating Training Strategy.** As shown in Table 2, we evaluated Med-VRAgent on medical VQA tasks using the LLaVA-Med v1.5 model, comparing it with five baselines: Zero-shot, SFT, Self-Rewarding, DPO, and STLLaVA-Med. On VQA-RAD, Med-VRAgent achieved 35.70 (open) and 68.72 (closed); on MIMIC-CXR, it scored 3.90 (BLEU), 13.53 (ROUGE-L), and 9.58 (METEOR), outperforming other methods in generalization and generation quality.

**Evaluating Reasoning Strategy.** As shown in Table 3, we tested Med-VRAgent’s reasoning strategy, hypothesizing that improved visual guidance and feedback and higher-quality auxiliary information enhance performance. On the DeepSeek-VL-7B and MiniCPM-V2 models, Med-VRAgent outperformed others, achieving top scores in BVR (51.82

Methods	BLEU	ROUGE-L	METEOR
LLaVA-Med v1.5	9.64	12.26	8.21
Greedy	11.47	15.38	12.69
Beam Search	12.10	16.21	13.17
DoLa	11.79	15.82	12.72
OPERA	10.66	14.70	12.01
VCD	10.42	14.14	11.59
MedDr	12.37	16.45	13.50
FactMM-RAG	14.70	18.05	15.92
RULE	27.53	23.16	27.99
MMed-RAG	31.38	25.59	32.43
<b>Med-VRAgent</b>	<b>33.45</b>	<b>26.81</b>	<b>33.12</b>

Table 4: Comparison of Med-VRAgent with RAG methods, including FactMM-RAG, MMed-RAG etc, on the IU-Xray (BLEU, ROUGE-L, METEOR) dataset, based on LLaVA-Med v1.5 model. The best score for each metric is highlighted in bold.

and 47.32) and average (46.74 and 43.41). Compared to Zero-shot, CoT, and ToT, it excelled in abnormality recognition, visual reasoning, and relational understanding, confirming the effectiveness of the Med-VRAgent in complex medical VQA.

**Performance Comparison of Decoding-based and RAG-based Methods.** As shown in Table 4, on the IU-Xray dataset, LLaVA-Med v1.5 performed poorly (BLEU=9.64), with modest improvements from Greedy and Beam Search (BLEU=12.10). MMed-RAG showed significant improvement (BLEU=31.38), while Med-VRAgent achieved the best results (BLEU=33.45, ROUGE-L=26.81, METEOR=33.12), demonstrating that Med-VRAgent enhances medical report generation quality.

## 5 Discussion

This section presents three experiments examining Med-VRAgent’s performance in medical visual reasoning tasks. The first investigates the importance of each component. The second explores the impact of MCTS width and depth on model accuracy. The third experiment evaluates the adaptive retrieval strategy (ARS) in the Reflection component compared with the traditional fixed Top-K method.

### 5.1 Analysis of Med-VRAgent’s Components

We conduct an ablation study on **Med-VRAgent** to assess the contribution of its key components to medical visual reasoning. As shown in Fig 3, removing any component leads to performance degradation, highlighting the critical role of each module in reasoning progression, relevance, coherence,

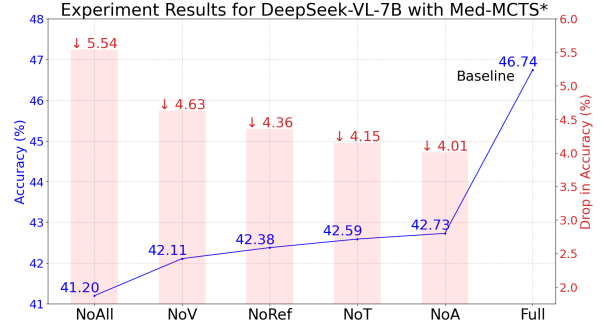


Figure 3: Ablation Experiment 1 Results (accuracy; %) for DeepSeek-VL-7B with Med-VRAgent on dataset GMAI-MMBench. Noall means removing all components, NoV means removing visual extraction, NOA means removing Assessor, and NoT means removing Teacher.

Accuracy Heatmap for Different Width and Depth

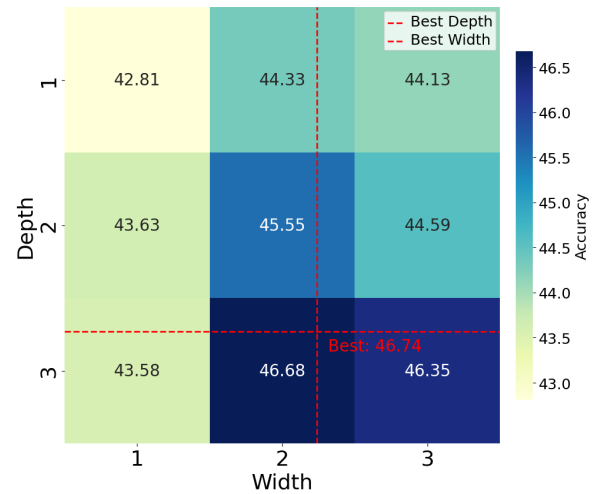


Figure 4: Ablation Experiment 2 Results (accuracy; %) for DeepSeek-VL-7B with Med-VRAgent on dataset GMAI-MMBench. Best is an adaptive exploration strategy, the average width and depth are 1.74 and 2.23 (red line), and other combinations are fixed width and depth.

and adaptability. The visual extraction component has the greatest impact. Specifically, omitting any module increases the error rate in LLMs, affecting reasoning quality.

### 5.2 Width and Depth Optimization

We studied the impact of different MCTS fixed widths and depths on performance. The results are shown Fig 4. By adjusting the fixed width and depth in the search strategy, we found that accuracy could be improved. Search benefits decrease as width and depth rise, likely due to VLM’s limited processing capacity. The best fixed combination (width 2, depth 3) achieved the highest accuracy of 46.68%. The adaptive strategy (width 1.74, depth 2.23) achieved an even higher accuracy of 46.74%. This result demonstrates that our adaptive strategy can maintain a balance between exploration and


	<p><b>Zero-shot:</b> ... clear lung field, .... The lungs appear to have normal bronchovascular markings, ... There is <b>no visible pleural effusion or pneumothorax, and the diaphragm appears intact.</b>... The lung parenchyma is free of significant pathological changes. The costophrenic angles are well-defined, and no blunting or fluid collection is visible. The chest X-ray appears normal .... <b>No signs of pneumonia, lung masses, or cardiovascular abnormalities are present.</b></p>	<p><b>GroundTruth:</b> ... The left lung is relatively well aerated and clear. The right hemithorax is markedly opacified with volume loss, circumferential pleural thickening and <b>pleural fluid</b> with near complete opacification of the right lung with right basal pleural catheter noted. ... Cardiac contours are somewhat obscured but unremarkable. ... Bibasilar opacities, larger on the left side, <b>could be due to atelectasis</b> but superimposed infection cannot be excluded. If any, there is a small right pleural effusion. There is elevation of the right hemidiaphragm. There is mild vascular congestion.</p>	<p><b>Med-VRAgent:</b> Increased density is observed in the left lower lung field, .... The grayscale value in this area is higher than the contralateral side, with reduced translucency, ... Blurred architecture of the left lower lobe..... <b>atelectasis ... should be considered;</b> ... Blurring of the right heart border may indicate involvement of the right middle or lower lobe, such as exudation or <b>pleural effusion.</b> Increased markings and reduced translucency ... possible inflammatory changes. .... No significant cardiomegaly or mediastinal shift is noted.</p>
---	--	---	--

Figure 5: Med-VRAgent Medical Report Generation Case Study

Experiment	Filter	Rerank	BLEU ↑	ROUGE-L ↑	METEOR ↑
Fixed Top-K	✗	✗	13.66	13.10	8.94
Rerank Only	✗	✓	13.75	13.20	9.22
Dynamic Top-K	✓	✗	13.80	13.75	9.12
Adaptive Retrieval	✓	✓	<b>13.90</b>	<b>14.10</b>	<b>9.58</b>

Table 5: Ablation study on the MIMIC-CXR dataset using the LLaVA-Med v1.5 model. Each retrieval strategy varies in its use of Filter and Rerank.

exploitation.

### 5.3 Evaluation of Adaptive Retrieval Strategy

The Table 5 presents an ablation study on the MIMIC-CXR dataset using the LLaVA-Med v1.5 model, evaluating different retrieval strategies. The experiments compare the impact of enabling filtering and Rerank mechanisms on the quality of generated outputs. The results indicate that using either Filter or Rerank alone leads to modest performance improvements. For instance, compared to the Fixed Top-K baseline, the Rerank Only strategy shows slight gains across all metrics. The best performance is Adaptive Retrieval, which combines both Filter and Rerank. It obtains the highest scores across all metrics.

### 5.4 Performance and Efficiency Analysis

In this experiment, we compared the performance of four methods (CoT, ToT, Med-VRAgent (Fix), Med-VRAgent (Ours)) on the GMAI-MMBench dataset. Fix is a fixed width of 2 and depth of 3. The results show that Med-VRAgent (Ours) performs best in terms of accuracy, reaching 46.74%. In addition, Med-VRAgent (Ours) has an advantage over Med-VRAgent (Fix) in inference time, which is 36.7 seconds, significantly lower than the fixed strategy of 45.7 seconds. Although the ToT method is slightly higher than CoT in accuracy (42.08% vs. 41.52%), its inference time is longer, reaching 31.3 seconds. The CoT method is the most efficient in inference time, only 18.3 seconds, but its accuracy is lower. Overall, Med-VRAgent (Ours) has achieved a good balance between accuracy and inference

Method	Accuracy (%)	Inference Time (s)
CoT	41.52	<b>18.3</b>
ToT	42.08	31.3
Med-VRAgent (Fix)	46.68	45.7
Med-VRAgent (Ours)	<b>46.74</b>	36.7

Table 6: DeepSeek-VL-7B compares the inference accuracy and average time of CoT, ToT and Med-VRAgent (fixed and adaptive policies) on the GMAI-MMBench dataset.

time, showing its comprehensive advantages over fixed strategies and other methods. This shows that adaptive strategies can optimize inference time while improving accuracy, and have better application potential.

### 5.5 Case Study

As shown in the Fig 5, the case comes from the Deepseek-VL and the MIMI-CXR dataset. Med-VRAgent outperforms the Zero-shot in generating clinically accurate and factually grounded chest X-ray reports. While the Zero-shot model incorrectly states clear lungs and no pleural abnormalities, Med-VRAgent correctly identifies increased density, reduced translucency, and possible pleural effusion in the left lung, closely matching the expert GroundTruth. It avoids major hallucinations and captures subtle findings like blurred architecture and right heart border changes, suggesting infection or inflammation. Med-VRAgent also includes diagnostic considerations such as atelectasis, reflecting expert-level reasoning.

## 6 Conclusion

This study introduces Med-VRAgent, a novel medical visual reasoning framework that enhances multimodal large models' performance in medical image understanding. It incorporates a teacher-student-evaluator mechanism, visual guidance and self-feedback paradigm, and a multi-step reasoning strategy based on MCTS. Med-VRAgent achieved top performance across several medical multimodal benchmark datasets, demonstrating proficiency in



image-text alignment, spatial structure understanding, and lesion recognition. Future research will focus on improving search efficiency, using advanced multimodal models, and expanding deployment in real clinical settings.

## Limitations

Although Med-VRAgent has achieved significant improvements in medical visual reasoning, it still has limitations. Despite optimization, tree search is still resource-intensive. Due to node expansion strategies and computational resource constraints, it may not be possible to fully search all possible reasoning paths. It may not be directly transferable to other domains and additional domain adaptation is required. Visual guidance may have limited effect in complex images or low-quality images. Inaccurate reasoning may still occur when faced with fine-grained errors or very complex cases. Performance and reliability in actual clinical settings have not been fully verified.

## Ethical Considerations

Ethical considerations are central to our research. In this study, we ensure adherence to ethical principles by exclusively using publicly available datasets and employing models that are open-source or widely accepted within the research community. We emphasize transparency in all stages of our work and prioritize the responsible application of technology, particularly in the sensitive domain of medical reasoning, to ensure that our contributions promote fairness, reliability, and societal benefit.

## References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, and et al. 2022. [Flamingo: a visual language model for few-shot learning](#).
- Emily Alsentzer, John R. Murphy, Willie Boag, and et al. 2019. [Publicly available clinical bert embeddings](#).
- Zechen Bai, Pichao Wang, Tianjun Xiao, and et al. 2025. [Hallucination of multimodal large language models: A survey](#).
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Daniil A. Boiko, Robert MacKnight, and Gabe Gomes. 2023. [Emergent autonomous scientific research capabilities of large language models](#).
- Tom B. Brown, Benjamin Mann, Nick Ryder, and et al. 2020. [Language models are few-shot learners](#).
- Jiawei Chen, Dingkan Yang, Tong Wu, et al. 2025. [Detecting and evaluating medical hallucinations in large vision language models](#). In *Proceedings of ICLR 2025*. Med-HallMark Benchmark.
- Pengcheng Chen, Jin Ye, Guoan Wang, and et al. 2024. [Gmai-mmbench: A comprehensive multimodal evaluation benchmark towards general medical ai](#). *arXiv preprint arXiv:2408.03361*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, and et al. 2022. [Palm: Scaling language modeling with pathways](#).
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, and et al. 2024. [Dola: Decoding by contrasting layers improves factuality in large language models](#).
- Dina Demner-Fushman, Madhur D. Kohli, Michael B. Rosenman, and et al. 2016. [Preparing a collection of radiology examinations for distribution and retrieval](#). *Journal of the American Medical Informatics Association*, 23(2):304–310. Epub 2015 Jul 1.
- Stefan Denner, Markus Bujotzek, Dimitrios Bounias, and et al. 2025. [Visual prompt engineering for vision language models in radiology](#).
- Yunfan Gao, Yun Xiong, Xinyu Gao, and et al. 2024. [Retrieval-augmented generation for large language models: A survey](#).
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, and et al. 2025. [A survey on llm-as-a-judge](#).
- Iryna Hartsock and Ghulam Rasool. 2024. [Vision-language models for medical report generation and visual question answering: a review](#). *Frontiers in Artificial Intelligence*, 7.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. [Deep residual learning for image recognition](#).
- Sunan He, Yuxiang Nie, Hongmei Wang, and et al. 2024. [GSCO: Towards generalizable ai in medicine via generalist-specialist collaboration](#).
- Lei Huang, Weijiang Yu, Weitao Ma, and et al. 2025. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ACM Transactions on Information Systems*, 43(2):1–55.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, and et al. 2024a. [Opera: Alleviating hallucination in multimodal large language models via over-trust penalty and retrospection-allocation](#).

- Yuhao Huang, Xin Yang, Lian Liu, et al. 2024b. [Segment anything model for medical images?](#) *Medical Image Analysis*, 92:103061.
- Ziwei Ji, Nayeon Lee, Rita Frieske, and et al. 2023. [Survey of hallucination in natural language generation.](#) *ACM Computing Surveys*, 55(12):1–38.
- Qiao Jin, Fangyuan Chen, Yiliang Zhou, et al. 2024. [Hidden flaws behind expert-level accuracy of multimodal gpt-4 vision in medicine.](#) *npj Digital Medicine*, 7:190.
- Alistair E. W. Johnson, Tom J. Pollard, Nathaniel R. Greenbaum, and et al. 2019. [Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs.](#)
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. [Billion-scale similarity search with gpus.](#)
- Yubin Kim, Hyewon Jeong, Shan Chen, et al. 2025. [Medical hallucination in foundation models and their impact on healthcare.](#) *medRxiv preprint*.
- Jason Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. 2018. [A dataset of clinically generated visual questions and answers about radiology images.](#) *Scientific Data*, 5:180251.
- Sicong Leng, Hang Zhang, Guanzheng Chen, and et al. 2023. [Mitigating object hallucinations in large vision-language models through visual contrastive decoding.](#)
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, and et al. 2021. [Retrieval-augmented generation for knowledge-intensive nlp tasks.](#)
- Chunyuan Li, Cliff Wong, Sheng Zhang, and et al. 2023a. [Llava-med: Training a large language-and-vision assistant for biomedicine in one day.](#)
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. [Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models.](#)
- Yingshu Li, Yunyi Liu, Zhanyu Wang, and et al. 2024. [A systematic evaluation of gpt-4v’s multimodal capability for medical image analysis.](#)
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries.](#) In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, and et al. 2024. [Grounding dino: Marrying dino with grounded pre-training for open-set object detection.](#)
- Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. 2023. [Can large language models reason about medical questions?](#)
- Haoyu Lu, Wen Liu, Bo Zhang, and et al. 2024. [Deepseek-vl: Towards real-world vision-language understanding.](#)
- Aman Madaan, Niket Tandon, Prakhar Gupta, and et al. 2023. [Self-refine: Iterative refinement with self-feedback.](#)
- Björn Moëll, Fredrik Sand Aronsson, and Shahid Akbar. 2025. [Medical reasoning in llms: an in-depth analysis of deepseek r1.](#) *Frontiers in Artificial Intelligence*, 8:1616145.
- Sophie Ostmeier, Justin Xu, Zhihong Chen, and et al. 2024. [Green: Generative radiology report evaluation and error notation.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2024*, page 374–390. Association for Computational Linguistics.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2023. [Med-halt: Medical domain hallucination test for large language models.](#) In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 314–334.
- Kishore Papineni, Salim Roukos, Todd Ward, and et al. 2002. [Bleu: a method for automatic evaluation of machine translation.](#) In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Chris Hallacy, and et al. 2021. [Learning transferable visual models from natural language supervision.](#)
- Rafael Rafailov, Archit Sharma, Eric Mitchell, and et al. 2024. [Direct preference optimization: Your language model is secretly a reward model.](#)
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks.](#)
- Daniel Rose, Vaishnavi Himakunthala, Andy Ouyang, and et al. 2024. [Visual chain of thought: Bridging logical gaps with multimodal infillings.](#)
- Hao Shao, Shengju Qian, Han Xiao, and et al. 2024. [Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning.](#)
- Liwen Sun, James Zhao, Megan Han, and Chenyan Xiong. 2025. [Fact-aware multimodal retrieval augmentation for accurate medical radiology report generation.](#)
- Xiangru Tang, Anni Zou, Zhuosheng Zhang, and et al. 2024. [Medagents: Large language models as collaborators for zero-shot medical reasoning.](#)
- Ryutaro Tanno, T. Barrett, David G. Andrew Sellergren, et al. 2025. [Collaboration between clinicians and vision-language models in radiology report generation.](#) *Nature Medicine*, 31:599–608.
- Hugo Touvron, Louis Martin, Kevin Stone, and et al. 2023. [Llama 2: Open foundation and fine-tuned chat models.](#)

- Krithik Vishwanath, Anton Alyakin, Daniel Alexander Alber, and et al. 2025. [Medical large language models are easily distracted.](#)
- Chaojie Wang, Yanchen Deng, Zhiyi Lyu, and et al. 2024. [Q\\*: Improving multi-step reasoning for llms with deliberative planning.](#)
- Yubin Wang, Xinyang Jiang, De Cheng, and et al. 2025. [Exploring interpretability for visual prompt tuning with hierarchical concepts.](#)
- Jason Wei, Xuezhi Wang, Dale Schuurmans, and et al. 2023. [Chain-of-thought prompting elicits reasoning in large language models.](#)
- Juncheng Wu, Wenlong Deng, Xingxuan Li, and et al. 2025. [Medreason: Eliciting factual medical reasoning steps in llms via knowledge graphs.](#)
- Peng Xia, Ze Chen, Juanxi Tian, and et al. 2024a. [Cares: A comprehensive benchmark of trustworthiness in medical vision language models.](#)
- Peng Xia, Kangyu Zhu, Haoran Li, and et al. 2024b. [Rule: Reliable multimodal rag for factuality in medical vision language models.](#)
- Peng Xia, Kangyu Zhu, Haoran Li, and et al. 2025. [Mmed-rag: Versatile multimodal rag system for medical vision language models.](#)
- Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, and et al. 2023. [Self-evaluation guided beam search for reasoning.](#)
- Shinn Yao, Jeffrey Zhao, Dian Yu, and et al. 2022. [React: Synergizing reasoning and acting in language models.](#)
- Shunyu Yao, Dian Yu, Jeffrey Zhao, and et al. 2023. [Tree of thoughts: Deliberate problem solving with large language models.](#)
- Yuan Yao, Tianyu Yu, Ao Zhang, and et al. 2024. [Minicpm-v: A gpt-4v level mllm on your phone.](#)
- Hongzhou Yu, Tianhao Cheng, Yingwen Wang, and et al. 2025. [Finemedlm-o1: Enhancing medical knowledge reasoning ability of llm from supervised fine-tuning to test-time training.](#)
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, and et al. 2025. [Self-rewarding language models.](#)
- Wenhui Zhang, Jianlong Fu, Yijuan Lu, and et al. 2023. [Reinforced ranker-reader for open-domain question answering.](#)
- Denny Zhou, Nathanael Schärli, Le Hou, and et al. 2023. [Least-to-most prompting enables complex reasoning in large language models.](#)
- Deyao Zhu, Jun Chen, Xiaoqian Shen, and et al. 2023. [Minigt-4: Enhancing vision-language understanding with advanced large language models.](#)
- Kangyu Zhu, Ziyuan Qin, Huahui Yi, and et al. 2025. [Guiding medical vision-language models with explicit visual prompts: Framework design and comprehensive exploration of prompt variations.](#)
- Kangyu Zhu, Peng Xia, Yun Li, and et al. 2024a. [Mmedpo: Aligning medical vision-language models with clinical-aware multimodal preference optimization.](#)
- Zhiying Zhu, Yiming Yang, and Zhiqing Sun. 2024b. [Halueval-wild: Evaluating hallucinations of language models in the wild.](#)

## A Prompt

### Prompt

Role Setting: You are a medical expert providing guidance on medical image analysis to help students improve their understanding.

Task Description: Focus on the red-boxed area in the image, using previous guidance and student feedback to offer optimized suggestions for enhancing their analysis skills.

Guidance Content:

Analyze Key Area:

Identify the red-boxed region for closer analysis.

Observe structural features, shape changes, color contrasts, and any abnormalities.

Reference Feedback and Suggestions:

Evaluate the student's previous analysis.

Point out missed details or inadequate analysis, and offer visual techniques.

Optimize Analysis Directions:

Guide the student based on the image type (e.g., CT, X-ray, ultrasound).

Suggest perspectives like cross-sections or tissue density changes.

Important Notes:

Your goal is to help students master image analysis, not to do it for them.

Focus on a logical, systematic approach for comprehensive image interpretation.

Previous Guidance: </Guidance>

Student's Answer: </Answer>

Feedback Information: </Feedback >

Use this format for guidance: </Guidance> Guidance here </Guidance>

Figure 6: ROI-Guided Teaching Prompt

Prompt

You are a medical expert. Please review the image and visual analysis guidance and rate the student-generated answers using the additional 5-point rating system described below. The rating will be cumulative based on the following criteria:

5-point rating scale:

1. Relevant information: If the medical vision answer provides some information that is relevant to the user's query, even if the information is incomplete or contains some incompletely relevant content, 1 point can be awarded. 2. Partially solve the problem: If the answer solves most of the user's question, but does not fully answer the user's question or does not directly answer the core query, 2 points can be awarded. 3. Essential elements: If the answer answers the basic elements of the user's question from a medical vision perspective, although it may lack detail or completeness in some aspects, but is still helpful to the user, 3 points can be awarded. 4. Direct and comprehensive solution to the problem: If the answer directly and comprehensively solves the user's question, although there may be some room for improvement in clarity, conciseness or visual focus, 4 points can be awarded. 5. Tailored, professional and profound: If the answer is tailored to the user's question, provides an in-depth and professional answer through medical vision, avoids irrelevant information, and produces high-quality, engaging and insightful content, 5 points should be awarded.

Information: <guidance> Teacher's guidance </guidance> <answer> Student's answer </answer>

Evaluation steps:

Total rating: Please briefly explain your rating in 100 words or less.

Suggestions for teachers: Provide suggestions for teachers to build better guidance in 100 words or less.

Revision suggestions for students: Provide revision suggestions for students in 100 words or less.

Rating conclusion:

<score>Integer score</score>

<feedback1>Feedback to teachers</feedback1>

<feedback2>Revision suggestions for students</feedback2>

Figure 7: ROI-Guided Evaluation Prompt

## B Med-VRAgent algorithm

---

### Algorithm 1 Med-VRAgent

---

**Input:** Question  $Q$ , Image  $I$ , Visual extractor  $\mathcal{V}$ , Teacher  $T$ , Assessor  $A$ , Student  $S$ , Retriever  $R$ , max\_depth  $D_{\text{epth}}$ , max\_branch\_number  $b$ , max\_simulation\_number  $Sim$

**Output:** Best solution path  $\pi^*$ ; Final answer  $A_{\text{final}}$

```

ROI  $\leftarrow \mathcal{V}(I, Q)$ ; // Region-of-interest detection
 $\mathcal{T} \leftarrow \text{Initializetree}(Q, I)$ 
for  $t = 1$  to  $Sim$  do
   $C \leftarrow \text{root}(\mathcal{T})$ 
  ---Selection---
  while  $C$  is not leaf node do // C is not a leaf
     $C \leftarrow \text{argmax}_s \text{UCB}(s)$ 
    if  $\text{depth}(C) \geq D_{\text{epth}}$  then // max depth reached
       $\perp$  break
    if  $C$  has less than  $b$  children nodes then // node not fully expanded
       $\perp$  break
    if  $\text{depth}(C) \geq D_{\text{epth}}$  then // skip if depth limit
       $\perp$  continue
  ---Expansion&Evaluation---
   $O_g^{\text{anc}} \leftarrow \bigcup_{k \in \text{ancestor}(C)} G_k$ ; // Teacher's guidance from ancestor
   $O_a^{\text{anc}} \leftarrow \bigcup_{k \in \text{ancestor}(C)} A_k$ ; // Student's answers from ancestor
   $O_g^{\text{sib}} \leftarrow \bigcup_{k \in \text{siblings}(C)} G_k$ ; // Teacher's guidance from siblings
   $O_a^{\text{sib}} \leftarrow \bigcup_{k \in \text{siblings}(C)} A_k$ ; // Student's answers from siblings
   $O_f^{\text{sib}} \leftarrow \bigcup_{k \in \text{siblings}(C)} F_k$ ; // Assessor's feedback from siblings
   $O \leftarrow (O_g^{\text{anc}}, O_a^{\text{anc}}, O_g^{\text{sib}}, O_a^{\text{sib}}, O_f^{\text{sib}})$ 
   $roi \leftarrow \text{SelectOnProb}(\text{P\_softmax}(Conf_{roi}))$ 
   $G \leftarrow T(roi, O)$ ; // Generate guidance (§3.3)
   $A \leftarrow S(roi, G)$ ; // Student answer (§3.3)
   $(R, F) \leftarrow A(roi, G, A)$ ; // Score & feedback (§3.3)
  if  $R == 5$  then // Stop early if  $A \cup O_a^{\text{anc}}$  receives full 5-point score
     $\perp$  break
   $C' \leftarrow \text{CREATENEWCHILD}(G, A, R, F, O)$ ; // create a new child node for  $C$ 
   $\text{ADDCHILD}(C, C')$ ; // add  $C'$  to the children of  $C$ 
  ---Backpropagation---
   $\text{BACKPROPAGATE}(C)$ ; // Update visit-count reward
 $\pi^* \leftarrow \text{BESTPATH}(\mathcal{T})$ ; // Highest cumulative reward
---Reflection---
for node in  $\pi^*$  do
  if  $R < 4$  then
     $\mathcal{K} \leftarrow \text{Rerank}(\text{Relevance}(\text{Top} - K(A, G, I)))$ ; // Retrieval (§3.4)
     $A^* \leftarrow S(roi, G, A, \mathcal{K})$ ; // Rewrite
     $\text{UPDATENODE}(A^*)$ 
 $A_{\text{final}} \leftarrow \text{COMPOSEANSWER}(\pi^*)$  return  $\pi^*, A_{\text{final}}$ 

```

---

## C Ablation Studies

### C.1 Visual Token Edit Ablation Results

Method	BOX	Edit	Accuracy
No VTE	no	no	42.11
Only BOX No Edit	yes	no	43.22
Only Edit No BOX	no	yes	45.56
VTE	yes	yes	46.74

Table 7: Visual Token Edit Ablation Results for DeepSeek-VL-7B (Student) with Med-VRAgent on GMAI-MMBench.

BOX refers to bounding box prompts on ROI, and Edit refers to attention enhancement on ROI. The results show that using both simultaneously yields the best performance, while omitting either leads to performance drops.

### C.2 Teacher Guidance Ablation Results

Method	Guidance	Answer	Feedback	Accuracy
1	no	no	no	42.03
2	yes	no	no	42.31
3	yes	yes	no	43.01
4	yes	no	yes	42.64
5	yes	yes	yes	43.41

Table 8: Teacher Guidance Ablation Results for MiniCPM-V2 with Med-VRAgent on GMAI-MMBench.

The three middle columns denote the information available to the teacher. For example, Method 1 indicates the teacher sees nothing and merely samples multiple times. Rows 2–5 progressively allow the teacher to access prior guidance, the student’s answers, and feedback from the assessor, validating the effectiveness of heuristic-based teacher guidance.

### C.3 GREEN Evaluation Results

Method	GREEN
Zero-Shot	0.21
RULE	0.29
MMed-RAG	0.31
Med-VRAgent (ours)	0.34

Table 9: GREEN scores of LLaVA-Medv1.5 on IU-Xray.

We performed a preliminary evaluation using the GREEN (Ostmeier et al., 2024) approach. GREEN uses LLMs as evaluators and provides scores that are consistent with expert preferences and human-interpretable for clinically significant errors (0 to 1, higher is better). We sampled 100 examples from the IU-Xray dataset. We will add more new experimental results later.