# UCAgents: Unidirectional Convergence for Visual Evidence Anchored Multi-Agent Medical Decision-Making

Qianhan Feng[1] , Zhongzhen Huang[2], Yakun Zhu[2,3] , Xiaofan Zhang[2,3*], Qi Dou[1*]

[1]The Chinese University of Hong Kong, Hong Kong SAR, China
[2]Shanghai Jiao Tong University, Shanghai, China
[3]Shanghai Innovation Institute, Shanghai, China

`qianhan.feng@link.cuhk.edu.hk, qidou@cuhk.edu.hk`

## Abstract

*Vision-Language Models (VLMs) show promise in medical diagnosis, yet suffer from reasoning detachment, where linguistically fluent explanations drift from verifiable image evidence, undermining clinical trust. Recent multi-agent frameworks simulate Multidisciplinary Team (MDT) debates to mitigate single-model bias, but open-ended discussions amplify textual noise and computational cost while failing to anchor reasoning to visual evidence, the cornerstone of medical decision-making. We propose **UCAgents**, a hierarchical multi-agent framework enforcing unidirectional convergence through structured evidence auditing. Inspired by clinical workflows, UCAgents forbids position changes and limits agent interactions to targeted evidence verification, suppressing rhetorical drift while amplifying visual signal extraction. In UCAgents, a one-round inquiry discussion is introduced to uncover potential risks of visual-textual misalignment. This design jointly constrains visual ambiguity and textual noise, a dual-noise bottleneck that we formalize via information theory. Extensive experiments on four medical VQA benchmarks show UCAgents achieves superior accuracy (71.3% on PathVQA, +6.0% over state-of-the-art) with 87.7% lower token cost, the evaluation results further confirm that UCAgents strikes a balance between uncovering more visual evidence and avoiding confusing textual interference. These results demonstrate that UCAgents exhibits both diagnostic reliability and computational efficiency critical for real-world clinical deployment. Code is available at https://github.com/fqhank/UCAgents.*

## 1. Introduction

Recent advances in Vision-Language Models (VLMs) [3, 7, 8, 15, 17, 21, 35, 36] have enabled unified reasoning across visual and textual modalities, achieving remarkable
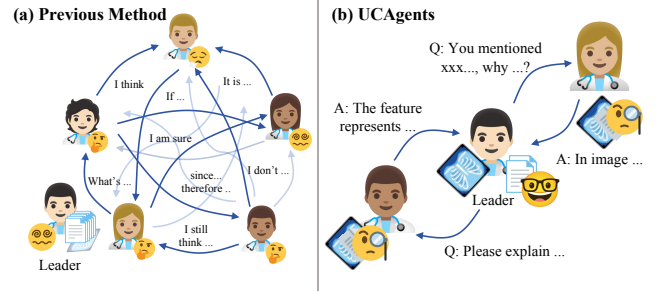


Figure 1. Unlike redundant discussion in previous multi-agent system, UCAgents uses one-round unidirectional inquiry to cut down textual noise and help focus on visual evidence.

generalization in open-domain tasks. However, when applied to medical Visual Question Answering (VQA), a key benchmark for reliable AI-assisted diagnosis, their reliability sharply declines. The limitation lies in the fragile coupling between visual evidence and diagnostic reasoning. Even a subtle misinterpretation can reverse a clinical conclusion, exposing the need for reasoning that remains tightly anchored to verifiable image evidence. Medical reasoning fundamentally differs from general multimodal tasks: each inference step must be grounded in observable, clinically meaningful visual cues. Yet, current VLMs often produce linguistically fluent but visually unsupported explanations, namely reasoning detachment, where textual fluency conceals evidential drift and undermines clinical trust.

To address this, recent works have explored *multi-agent collaboration* to simulate Multidisciplinary Team (MDT) discussions in clinical practice [10, 13]. While debate-based systems such as MDAgents [10] encourage diverse reasoning, they suffer from an information–noise paradox: as discussion rounds increase, total information expands but textual noise $N_t$ inflates, leading to rhetorical overfitting and unstable convergence. Conversely, single-agent systems avoid excessive $N_t$ but amplify visual noise $N_v$ due to lack

---

*Corresponding Author

of cross-validation. These dual sources of uncertainty form the dual-noise bottleneck, where suppressing one inevitably worsens the other.

We propose **UCAgents** (Unidirectional Convergence Agents), a hierarchical multi-agent framework that replaces open debate with structured, entropy-reducing collaboration. UCAgents enforces a hierarchical flow of reasoning through three tiers: (1) independent agents perform controlled initial divergence to quantify uncertainties; (2) a supervisory reviewer audits visual-textual alignment to remove false consensus; and (3) a leader agent conducts evidence-anchored one-round risk inquiry by clarifying ambiguous findings and selecting the most visually consistent hypotheses. This process mirrors the clinical workflow of screening, verification, and adjudication, progressively filtering noise and consolidating diagnostic consensus around image-grounded evidence.

We evaluate UCAgents on several medical benchmarks, covering pathology, radiology, and mixed imaging modalities. Across both open-source and proprietary models, UCAgents consistently achieves superior visual-evidence anchored performance. Our main contributions lie in:

- We identify and formalize the dual-noise bottleneck in medical VQA, revealing how visual ambiguity and textual drift jointly degrade diagnostic reliability.
- We propose a hierarchical, entropy-minimizing multi-agent framework that transforms open debates into unidirectional, evidence-anchored convergence.
- We demonstrate that UCAgents achieves consistent gains across four medical VQA benchmarks and various backbones, delivering clinically reliable reasoning with significantly improved accuracy and interpretability.

## 2. Related Work

### 2.1. Medical Multimodal Reasoning and Visual Evidence Grounding

Medical VQA has become a key benchmark for evaluating multimodal understanding in medical AI [1, 9, 11, 12, 14, 31, 37], requiring models to interpret fine-grained radiological features and integrate clinical queries into coherent reasoning [9, 11, 12]. Datasets such as PathVQA [9], VQA-RAD [12], and SLAKE [11] provide large-scale, clinically annotated image-question-answer pairs, lay the foundation for standardized evaluation across perception, interpretation, and diagnostic inference.

Recent advances in VLMs such as GPT-4 [21], LLaVA [15], Qwen-VL [3], and Gemini [8] have achieved impressive performance on open-domain multimodal reasoning. However, when applied to medical VQA, their visual encoders trained on natural images lack sensitivity to clinical features such as lesion texture, margin irregularity, or enhancement patterns [18]. They often produce lin-

guistically fluent but visually unsupported reasoning, while domain-specific medical foundation models [5, 22] capture fine-grained features but exhibit poor generalization across modalities or question types.

To improve reliability, researchers have explored three directions. Prompt-engineering strategies like Chain-of-Thought (CoT) and Self-Consistency (SC) [29, 30] enhance logical structure but remain language-centric. Knowledge-augmented models [25] integrate medical ontologies but rarely verify consistency between reasoning and visual evidence. Despite progress, existing systems fail to jointly constrain visual ambiguity and textual redundancy, leaving diagnostic reasoning loosely anchored to image evidence.

### 2.2. Multi-Agent Collaboration for Clinical Decision-Making

Multi-agent systems inspired by Multidisciplinary Team collaboration [6, 10, 34] have been proposed to enhance robustness through structured dialogue among specialized agents. Early frameworks such as CAMEL [13] and Auto-Gen [32] decompose tasks into fixed roles, improving interpretability but lacking agents explicitly responsible for verifying visual–textual alignment. Later works, MDAgents [10], ReConcile [4], and Reflexion [23], introduced iterative debates or self-reflection to refine consensus. While these systems simulate human-like deliberation, open-ended exchanges often amplify textual noise and rhetorical drift, causing reasoning to deviate from visual evidence. Dynamic optimization frameworks [16, 26] further adapt team size or communication depth based on task complexity, yet focus mainly on interaction efficiency rather than maintaining evidence-centered reasoning.

In summary, current multi-agent methods oscillate between rigid role assignments that limit perspective diversity and unconstrained debates that increase linguistic entropy. Few explicitly model how multimodal evidence should be verified and filtered during collaboration. Our framework, UCAgents, addresses this gap through hierarchical, unidirectional convergence, enforcing structured information flow that stabilizes reasoning and anchors collaboration to verifiable visual evidence.

## 3. Method: UCAgents

### 3.1. Problem Formulation and Motivation

Medical Visual Question Answering is a critical evidence-driven task where diagnostic reliability depends on how well reasoning anchors to medical images, the core evidence in multimodal inputs. We formalize this task through an information-theoretic lens to reveal performance bottlenecks in existing methods and motivate our structured multi-agent framework.

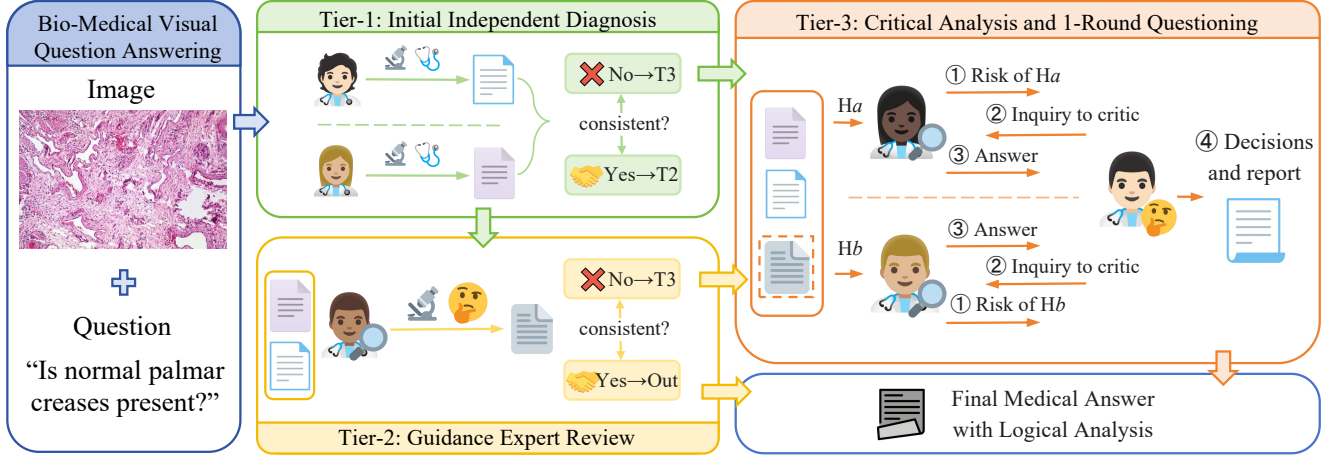**Formalizing Medical VQA: The Primacy of Visual-**

Figure 2. The overview of UCAgents. UCAgents system is composed of 3 dynamic Tiers: Initial Independent Diagnosis, Guidance Expert Review and Critical Analysis and Questioning. $H_a$, $H_b$: Divergent Candidate Hypotheses from previous tiers.

**Evidence Alignment**. Given a medical image $V$, clinical query $T$, the goal is to extract core visual evidence according to the textual constraint, and predict the correct answer $Y$ from the diagnostic hypothesis space $\mathcal{H}$. According to Fano's Inequality [28], the diagnostic error rate $P_e$ is bounded by:

$$P_e \geq \frac{H(Y) - I(Y;\mathcal{I}) - 1}{\log |\mathcal{H}|}, \quad (1)$$

where $H(Y)$ is the prior entropy of the answer space, and $\mathcal{I} = (V, T, M)$ represents total observable information, visual evidence $V$, textual query $T$, and auxiliary information $M$ (e.g., agent interactions). The chain rule of mutual information yields:

$$I(Y;V,T) = I(Y;V) + I(Y;T|V), \quad (2)$$

where $I(Y;V)$ captures core evidence association and $I(Y;T|V)$ captures task-guided visual reasoning. This leads to a critical insight: to minimize diagnostic errors, we must maximize $I(Y;\mathcal{I})$ by prioritizing $I(Y;V)$ while enhancing $I(Y;T|V)$, aligning with medical VQA's evidence-centric nature where diagnoses are ultimately grounded in images.

Unlike natural image tasks, Medical VQA is more difficult for VLMs because of the scarcity of high-quality medical data, stringent demands for accurate visual-textual alignment, and the high-stakes nature of medical decision-making. Single-agent VLMs for Medical VQA suppress $I(Y;\mathcal{I})$ through inadequate multimodal alignment:

- **Unaudited visual-textual misalignment**: Small perceptual errors (e.g., misidentifying benign lesions) propagate unchallenged, reducing $I(Y;V)$ and breaking the link between textual arguments and visual features.
- **Single-point bias**: Lack of cross-validation leads to model biases (e.g., overdiagnosing common conditions),

limiting $I(Y;T|V)$ as reasoning drifts from task-specific queries to biased priors.

Recent multi-agent frameworks inspired by MDT debates allow open, multi-round discussions. While designed to enhance $I(Y;\mathcal{I})$, they introduce detrimental entropy inflation that degrades cross-modal alignment:

- **Unconstrained decision entropy**: Agents freely revise positions and explore open-ended justifications, expanding conditional entropy $H(Y|A_i)$. For medical VQA where answers are uniquely anchored to $V$, this diverts attention from visual-textual alignment to rhetorical persuasion.
- **Mutual information interference**: Lengthy debates introduce redundant textual noise, making $I(Y;V,M) \leq I(Y;V)$. Instead of enhancing alignment, auxiliary information $M$ creates distraction: $I(Y;M|V) > 0$ but often correlates with incorrect justifications.
- **Group consensus bias**: Open debates enable "herd behavior": agents with correct thinking may be swayed by majority incorrect reasoning, suppressing $I(Y;V)$.
- **Prohibitive computational cost**: Open debates require extensive tokens, but entropy inflation means cost doesn't translate to proportional performance gains.

**Principles of UCAgents**. Previous failures point to a clear principle: multi-agent collaboration must be structured to serve multimodal alignment and evidence-centricity, not open-ended debates. To this end, we propose a new multi-agent system for Medical VQA, UCAgents, which is demonstrated in Fig. 2, achieves three complementary goals:

- **Enhance visual-textual alignment**: Leverage multi-agent collaboration to audit and refine visual feature extraction directly boosting $I(Y;V)$ and $I(Y;T|V)$.
- **Suppress irrelevant textual noise**: Constrain interactions to avoid open debates, ensuring exchanges target

validating visual evidence.
- **Eliminate group consensus bias**: Fix agent roles and enforce unidirectional convergence so reasoning relies solely on visual-textual evidence.

### 3.2. Overview of UCAgents Framework

UCAgents is a dynamic 3-Tiers system as shown in Fig. 2, it establishes a hierachical unidirectional process: **diversify → verify → converge**. Each Tier systematically reduces irrelevant entropy and reinforces visual grounding through fixed roles and structured interactions. In UCAgents, no one single agent is allowed to change their stances, and the opinion information provided to each one is strictly limited, thus to reduce the distraction from the image. We introduce each Tier and the special mechanism of UCAgents in detail in the following sections.

### 3.3. Tier-1: Initial Independent Diagnosis

Tier-1 objectively assesses the given medical case through structured collaboration. MDAgents [10] uses a single agent to evaluate the difficulty of the case, thus to determine whether the case is given to a single expert agent to process or to a MDT group to discuss. However, this highly subjective and biased process often leads to inappropriate case allocations. In UCAgents, we assess the case by conducting parallel independent diagnosis: two identical expert agents $\mathcal{A}_{1-1}$ and $\mathcal{A}_{1-2}$ are deployed to answer the question based on input image respectively. There is no communication between agents, thus to ensure reasoning anchors to original input $V$ and $T$, and divergence arises solely from independent visual interpretation.

It is important to encourage diversity between the two agents in their thinking. Classical operations include image data augmentations, but this would inevitably increase the difficulty of interpreting medical image which is already quite noisy. Some other method let different agents focus on different perspective like global versus local part. However, this operation explicitly ignores the textual needs, while most medical case need to combine both local and global information to give the correct diagnosis. To generate diversity without compromising image integrity, we introduce asymmetry via temperature modulation ($\tau = 0.7$). This slightly higher temperature encourages exploration of subtle visual variations without drifting into unconstrained textual reasoning. Since agents share identical setups, divergence amplifies inherent image ambiguity (e.g., ill-defined lesion boundaries) rather than introducing false differences. Each agent outputs an initial report that includes:
- $H_{1-i}$: diagnostic hypothesis.
- $R_{1-i}$: brief textual justification tied to visual features.

The two independent initial reports, $H_{1-1}$ and $H_{1-2}$, are compared: $D = \mathbb{I}[H_{1-1} \neq H_{1-2}]$. Disagreement $D = 1$ signals high divergence and difficulty requiring further in-depth discussions from experienced experts. However, consensus $D = 0$ cannot indicate reliable diagnosis, but may hide shared biases needing verification. We allocate the case into different reviewing routes based on $D$:

$$\text{Route}(H_{1-1}, H_{1-2}) = \begin{cases} \text{Tier-3}, & D = 1, \\ \text{Tier-2}, & D = 0. \end{cases} \quad (3)$$

Tier-1's entropy partitioning operator $\mathcal{P}_1$ formalizes the mutual information division:

$$I_1(Y; V, T) = I(Y; V, T|D = 1) + I(Y; V, T|D = 0), \quad (4)$$

where $I(Y; V, T|D = 1)$ and $I(Y; V, T|D = 0)$ represent high-entropy and low-entropy cases respectively. The no-interaction design preserves $\mathcal{I} = (V, T)$ and avoids textual noise, while dual-agent reasoning enhances $I(Y; V)$ by double-mining visual evidence.

### 3.4. Tier-2: Consensus Purification via Visual-Textual Alignment Verification

Tier-2 tries to address the potential "false consensus", where Tier-1 agents agree due to shared biases rather than accurate visual-textual alignment. Its mission is to verify consensus authenticity, filter misaligned conclusions, and either finalize reliable diagnoses or escalate complex cases.

Tier-2 deploys a single Guidance Expert Agent $\mathcal{A}_2$, analogous to a supervising physician. $\mathcal{A}_2$ specializes in bidirectional validation - confirming whether consensus anchors to visual evidence and identifying potential flaws. It operates at moderate temperature ($\tau = 0.5$) for stability and rigorous evidence checking.

The agent's input materials leverages Tier-1's outputs: original medical image $V$, clinical query $T$, and consensus report from Tier-1 (including $H_1 = H_{1-1} = H_{1-2}$ and justifications $R_{1-1}, R_{1-2}$). Its structured workflow is formulated as:
1. **Comprehensive Visual Scan**: Identify core diagnostic features in $V$ (e.g., lesion shape, density), but without providing preconceived judgments or positions. This step solely helps Guidance Expert have an initial knowledge of the case.
2. **Evidence verification**: Parse $R_{1-1}$ and $R_{1-2}$ to verify each claim aligns with visual features (eliminating hallucinations) and no critical features are omitted.
3. **Logic verification**: Parse $R_{1-1}$ and $R_{1-2}$ to check whether the visual evidence is correctly interpreted, and whether the logical reasoning in the reports is rigorous and correct.
4. **Independent hypothesis generation**: Produce a new diagnosis $H_2$ and report $R_2$ based on extracted verification rather than independent reasoning.

The consensus purification operator $\mathcal{P}_2$ in Tier-2 is:

$$I_2(Y; V, T|D = 0) = I(Y; V, T|H_2 = H_1), \quad (5)$$

Guidance Expert gains access to more ideas and information while conducting reviews, but this information is neither disordered nor noisy. Instead, it is confined to different lines of reasoning that lead to the same conclusion. When the expert's conclusion $H_2$ aligns with that of Tier-1 $H_1$, it indicates the conclusion has high reliability. However, if they identify flaws after reviewing Tier-1's conclusion and propose an alternative solution, this signifies the case is quite high in complexity and requires further research:

$$\text{Route}(H_1, H_2) = \begin{cases} \text{Terminates,} & H_2 = H_1, \\ \text{Tier-3,} & H_2 \neq H_1. \end{cases} \quad (6)$$

Consensus is finalized only if hypotheses agree and alignment is verified, ensuring conclusions are correct for the right reasons. Disagreement or detected misalignment (e.g., unsubstantiated claims, omitted features) triggers escalation. This enhances mutual information: $I_2(Y; V, T | D = 0) > I_1(Y; V, T | D = 0)$ by filtering false consensus derived from shared biases.

### 3.5. Tier-3: Unidirectional Risk Auditing with Targeted Inquiry

Building upon the structured uncertainty estimation (Tier-1) and bias correction (Tier-2), Tier-3 conducts final adversarial risk auditing to ensure convergence toward evidence-grounded diagnosis. Facing this challenge, Tier-3 deploys specialized agents with fixed, non-overlapping roles:

- **2 Critical Analyst Agents** ($\mathcal{C}_1, \mathcal{C}_2$): Examine the potential risk of existing hypotheses.
- **1 Leader Agent** ($\mathcal{A}_L$): Acts as the MDT chairperson who facilitates discussion and makes the final arbitration.

By organizing a discussion with these agents, Tier-3 operates via a four-step process to establish the final diagnosis.

**Step 1: Unidirectional Risk Mining.** Both $\mathcal{C}_1$ and $\mathcal{C}_2$ are exclusively assigned to challenge 1 specific hypothesis. Unlike Tier-2's Guidance Expert who evaluates whether reasoning is sound, Critical Analysts adopt a **"devil's advocate" stance**: actively mining potential flaws and risks that could invalidate the hypothesis with moderate 0.5 temperature. Each analyst outputs a structured risk report:

$$\mathcal{C}(H_i, R_i, V, T) \rightarrow R_i^{risk}, \quad (7)$$

where $R_i^{risk}$ contains identified flaws in analysis.

This mechanism has 2 advantages: (1) Opinions on each hypothesis are consolidated into one report respectively, avoiding herd behavior caused by discrepancies in the number of agents holding different views, (2) Targeted critical risk reports can offset the persuasive rhetoric of supporting opinions to a certain extent.

**Step 2&3: Leader-Directed Inquiry and Expert Response.** After receiving both risk reports, the Leader Agent

$\mathcal{A}_L$ reviews them to identify areas requiring deeper investigation. Unlike open debates where agents freely exchange arguments, the Leader issues a single targeted question $Q_i$ to each Critical Analyst to probe specific ambiguities or underdeveloped claims in their risk assessment. This inquiry operator is formalized as:

$$\hat{R}_i^{risk} = \mathcal{A}_L(R_i^{risk}, Q_i, V, T), \quad (8)$$

where $\hat{R}_i^{risk}$ is the supplemented response. This step mirrors clinical practice where the MDT chairperson asks pointed questions to clarify conflicting interpretations without allowing open-ended debate. The inquiry $Q_i$ acts like applying attention to the logic and encourage reinforcement. If the original $H_i$ is wrong, $\hat{R}_i^{risk}$ would be more robust. In the contrast, if $H_i$ is correct, then $\hat{R}_i^{risk}$ would reveal more flaws in $R_i^{risk}$. Uniform 0.1 temperature is used.

**Step 4: Final Arbitration.** Given the response from both $\mathcal{C}_1$ and $\mathcal{C}_2$, Leader agent $\mathcal{A}_L$ is required to evaluate the overall risks by comparing the reports and aggregating all the information, and make the final diagnosis:

$$Y^* = \mathcal{A}_L(\hat{R}_a^{risk}, \hat{R}_b^{risk}, R, V, T), \quad (9)$$

and the Leader is also allowed to make a diagnosis out of $H_a, H_b$, if enough evidence suggests that both are risky and identified evidences lead to a new solution.

The key innovation of Tier-3 lies in its unidirectional risk auditing mechanism: instead of having agents debate which diagnosis is correct, which may lead to rhetorical persuasion and entropy inflation, we task agents with actively searching for reasons why one specific hypothesis might be *wrong*, and one-round inquiry-response communication helps the unidirectional convergence of the critical logics. This adversarial approach forces deep scrutiny of visual evidence while avoiding the pitfalls of open debate.

The unidirectional auditing maximizes mutual information $I(Y; V, T)$ while suppressing irrelevant entropy from textual debates $M$. The adversarial risk reports force agents to extract counter-evidence from $V$ that might invalidate each hypothesis, effectively amplifying signal from visual features that distinguish correct from incorrect diagnoses. By constraining communication to a single targeted inquiry per agent, we prevent the entropy inflation $H(Y | M)$ that plagues open debate systems, where lengthy arguments introduce noise orthogonal to visual evidence. Overall, The framework maximizes $I(Y; \mathcal{I})$ by prioritizing $I(Y; V)$ and $I(Y; T | V)$ while suppressing textual noise from $M$.

## 4. Experiments

### 4.1. Experimental Setup

We conduct experiments on multiple medical VQA datasets and across various vision-language backbones. All experiments on open-source models are implemented under the

Ollama framework [20] on a single NVIDIA H800 GPU, with results averaged over three independent trials. For fairness, we adopt a unified temperature configuration across all datasets and models. All agents are guided by standardized prompt templates, which can be found in appendix A to ensure consistent behavior across settings, without exemplar demonstrations. This guarantees reproducibility while minimizing prompt engineering bias, enabling a clean comparison between different multi-agent paradigms. All datasets used are open-source, and the research has obtained the necessary licenses and certifications.

## 4.2. Results based on GPT-4

We first evaluate UCAgents using **GPT-4** [21] as the base model. GPT-4 represents one of the advanced vision–language models, capable of joint image and text understanding across diverse tasks. We follow MDAgents to conduct experiments on two representative medical VQA benchmarks: **PathVQA** [9] and **MIMIC-CXR-VQA** [2]. PathVQA focuses on pathology slides involving fine-grained questions on tissue structure, cell morphology, and diagnostic patterns; MIMIC-CXR-VQA contains clinical chest X-rays paired with textual questions covering thoracic diseases such as pneumonia, effusion, and cardiomegaly. These two datasets collectively span both microscopic and radiographic diagnostic scenarios.

As shown in Table 1, UCAgents achieves 71.3% on PathVQA and 60.3% on MIMIC-CXR-VQA, consistently outperforming all baselines. Compared to the best single-agent baseline, UCAgents yields a 9.9% improvement. Among multi-agent methods, UCAgents surpasses the MDT-based MDAgents by 6.0%, confirming the effectiveness of unidirectional convergence mechanism in maintaining visual-textual alignment and diagnostic stability. These results demonstrate that even with a strong base model like GPT-4, structured and directed collaboration remains critical for reliable medical reasoning, validating UCAgents as a robust enhancement over existing multi-agent paradigms.

## 4.3. Results based on Open-Source Models

We further extend experiments to open-source VLMs including **Qwen2.5VL** [3] series and **LLaVA** [15]. In addition to the PathVQA and MIMIC-CXR-VQA, we include two complementary benchmarks, **VQA-RAD** [12] and **SLAKE-VQA** [11]. VQA-RAD integrates radiology images across CT, MRI, and X-ray modalities, while SLAKE-VQA focuses on anatomical reasoning and multilingual comprehension. Both datasets include a mixture of *closed-form* and *open-ended* questions. For the open-ended subset, we follow a standardized setup: GPT-4 is used to generate 3–4 plausible answer options including the correct one, ensuring balanced difficulty and fair comparison. We thus report both *full-set/closed-set* scores per dataset.

Table 1. Accuracy Results (%) based on GPT-4 model.

| Category | Method | Path-VQA | MIMIC-CXR |
|---|---|---|---|
| Single Agent | Zero-Shot | $57.9_{\pm1.6}$ | $40.0_{\pm5.3}$ |
| | Few-Shot | $57.5_{\pm4.5}$ | $35.3_{\pm5.0}$ |
| | +CoT[30] | $58.6_{\pm3.1}$ | $36.2_{\pm5.2}$ |
| | +CoT-SC[29] | $61.2_{\pm2.1}$ | $51.7_{\pm4.0}$ |
| | ER[24] | $61.4_{\pm4.1}$ | $50.0_{\pm0.0}$ |
| | MedPrompt[19] | $59.2_{\pm5.7}$ | $53.4_{\pm4.3}$ |
| Multi Agents | Reconcile[4] | $57.5_{\pm3.3}$ | $33.3_{\pm3.4}$ |
| | AutoGen[33] | $43.0_{\pm8.9}$ | $43.3_{\pm8.2}$ |
| | DyLAN[16] | $41.3_{\pm1.2}$ | $38.7_{\pm1.2}$ |
| | MedAgents[27] | $45.4_{\pm8.1}$ | $43.3_{\pm7.0}$ |
| | Meta-Prompt[26] | $55.3_{\pm2.3}$ | $42.0_{\pm4.0}$ |
| | MDAgents[10] | $65.3_{\pm3.9}$ | $55.9_{\pm9.1}$ |
| | **UCAgents** | $\mathbf{71.3}_{\pm1.3}$ | $\mathbf{60.3}_{\pm0.5}$ |

As shown in Table 2, UCAgents consistently outperforms both single-agent and SOTA multi-agent method MDAgents across all settings. Notably, on Qwen2.5VL-3B, UCAgents achieves an average gain of 11.4%, while on the larger 72B model, it still provides a solid 3.5% improvement. On SLAKE-VQA, UCAgents reaches up to 80.0%. Meanwhile, on VQA-RAD, UCAgents delivers +8.5% and +10.6% gains on the total and closed sets respectively, further validating its consistency across question types. What's more, with Qwen2.5VL-3B, MDAgents classifies almost all cases into '*Basic*' category, thus performing equally as single agent.

These results validate the proposed unidirectional convergence strategy's generalization across architectures and modalities. Critically, UCAgents turns lightweight open-source VLMs into clinically reliable assistants while boosting lightweight models to match larger models' diagnostic performance and enhancing powerful models' stability by reducing dangerous misdiagnoses.

## 4.4. Ablation Study

We evaluate several key components proposed in UCAgents, and the results are reported in Table 3, which is conducted based on LLaVA-7B model on VQA-RAD dataset. The supervisor agent in Tier-2 provides further review on potential false consensus from Tier-1, and provides a 3.54% accuracy gain to the final performance. In Tier-3, which is the core mechanism of UCAgents, there are 3 important methods ensuring the success of discussion: (1) one-round inquiry, (2) 2 independent Critical Analyst Agents and (3) unidirectional critics instead of support. In Table 3, removing the leader's questioning step reduces overall performance by 15.60%, and an extreme drop of 27.23% on the accuracy of Tier-3 expert review can also be witnessed. Also, using one single agent (instead of two independent

Table 2. Performance on open-source models, with accuracy reported in percentage (%). For VQA-RAD and SLAKE-VQA, the performances are reported in the format of *full-set/closed-set*. The consumed tokens are reported in the format of Input/Output token numbers. 1: Qwen2.5VL-3B based MDAgents is not reported since the system classifies all cases into the '*Basic*' category. 2: The running time of Qwen2.5VL-72B based MDAgents on Path-VQA and MIMIC-CXR-VQA largely surpasses affordable range.

| Method | Path-VQA | MIMIC-CXR-VQA | VQA-RAD | SLAKE-VQA | AVG. | API Calls | Tokens(K) |
|---|---|---|---|---|---|---|---|
| **Qwen2.5VL-3B** | $37.00_{\pm16.26}$ | $44.05_{\pm1.24}$ | $40.16_{\pm7.98}/43.03_{\pm8.62}$ | $52.05_{\pm0.30}/55.78_{\pm3.59}$ | 45.35 | 1.00 | 1.02/0.06 |
| w/MDAgents[1] | - | - | - | - | - | - | - |
| w/UCAgents | $\mathbf{61.46}_{\pm1.67}$ | $\mathbf{53.09}_{\pm2.80}$ | $\mathbf{52.64}_{\pm1.46}/\mathbf{54.09}_{\pm1.82}$ | $\mathbf{58.64}_{\pm2.25}/\mathbf{60.58}_{\pm1.00}$ | **56.75** | 5.38 | 6.49/0.40 |
| **Qwen2.5VL-7B** | $55.15_{\pm10.21}$ | $52.61_{\pm2.28}$ | $53.59_{\pm0.11}/57.37_{\pm1.13}$ | $59.66_{\pm0.78}/60.71_{\pm0.21}$ | 56.15 | 1.00 | 0.99/0.06 |
| w/MDAgents | $58.13_{\pm3.42}$ | $47.23_{\pm2.63}$ | $53.44_{\pm1.60}/59.05_{\pm2.23}$ | $57.60_{\pm8.60}/\mathbf{63.02}_{\pm1.79}$ | 56.41 | 10.15 | 21.02/1.74 |
| w/UCAgents | $\mathbf{61.90}_{\pm0.62}$ | $\mathbf{54.45}_{\pm2.55}$ | $\mathbf{56.12}_{\pm0.00}/\mathbf{61.35}_{\pm0.00}$ | $\mathbf{61.45}_{\pm2.93}/62.48_{\pm0.11}$ | **59.63** | 4.49 | 5.97/0.45 |
| **Qwen2.5VL-32B** | $63.12_{\pm0.08}$ | $53.30_{\pm0.62}$ | $62.37_{\pm1.91}/65.12_{\pm0.31}$ | $71.53_{\pm1.88}/69.72_{\pm0.59}$ | 64.19 | 1.00 | 0.95/0.08 |
| w/MDAgents | $61.80_{\pm3.13}$ | $51.68_{\pm0.28}$ | $60.59_{\pm2.06}/65.46_{\pm1.90}$ | $73.13_{\pm1.48}/69.24_{\pm2.69}$ | 63.65 | 8.43 | 23.76/2.23 |
| w/UCAgents | $\mathbf{63.96}_{\pm0.28}$ | $\mathbf{54.14}_{\pm0.40}$ | $\mathbf{66.34}_{\pm1.73}/\mathbf{67.60}_{\pm0.57}$ | $\mathbf{73.63}_{\pm1.18}/\mathbf{73.60}_{\pm1.71}$ | **66.55** | 4.10 | 5.52/0.37 |
| **Qwen2.5VL-72B** | $63.43_{\pm3.60}$ | $56.36_{\pm2.44}$ | $65.49_{\pm1.78}/74.91_{\pm1.69}$ | $73.22_{\pm0.16}/74.61_{\pm1.25}$ | 68.00 | 1.00 | 0.97/0.06 |
| w/MDAgents[2] | - | - | $64.28_{\pm2.28}/70.31_{\pm2.68}$ | $\mathbf{75.22}_{\pm1.34}/72.82_{\pm2.38}$ | - | 6.87 | 14.75/1.61 |
| w/UCAgents | $\mathbf{67.61}_{\pm0.57}$ | $\mathbf{62.55}_{\pm1.24}$ | $\mathbf{68.42}_{\pm0.16}/\mathbf{75.67}_{\pm1.65}$ | $74.69_{\pm0.89}/\mathbf{80.04}_{\pm1.25}$ | **71.50** | 3.71 | 4.78/0.32 |
| **LLaVA-7B** | $47.78_{\pm1.23}$ | $44.05_{\pm1.87}$ | $44.93_{\pm1.67}/47.30_{\pm4.67}$ | $47.08_{\pm2.47}/46.69_{\pm0.69}$ | 46.31 | 1.00 | 0.84/0.07 |
| w/MDAgents | $50.49_{\pm5.25}$ | $44.23_{\pm3.43}$ | $47.72_{\pm5.30}/47.87_{\pm1.60}$ | $43.84_{\pm1.89}/46.05_{\pm1.20}$ | 46.70 | 13.39 | 33.57/2.86 |
| w/UCAgents | $\mathbf{54.56}_{\pm1.22}$ | $\mathbf{54.63}_{\pm1.53}$ | $\mathbf{61.03}_{\pm3.14}/\mathbf{54.16}_{\pm0.89}$ | $\mathbf{60.03}_{\pm1.58}/\mathbf{50.04}_{\pm0.95}$ | **55.74** | 4.48 | 5.78/0.45 |

agents) to criticize both hypotheses reduces accuracy by 4.58%, since the bias is inevitable. If the agents are instructed to provide supportive reports on hypotheses instead of finding risks, the discussion quality also suffers, which demonstrates that unidirectional critic is important for mitigating the eloquence effect. To further validate UCAgents' effectiveness, we analyze the performance of different diagnostic routes. According to the GPT-4 based experiment on Path-VQA in Table 4, the validated consensus from supervisor in Tier-2 achieves an accuracy of 73.47%. For complex cases where disagreements arose at the Tier-1 level, the accuracy rate improved by 11.64% after consultation with Tier-3 experts, compared to using a single agent to make a diagnosis directly (49.48%). Identifying consensus risks in Tier-2 reviews often indicates an extremely high level of complexity in the cases, as evidenced by the poor performance of using a single agent to make decisions on these cases (43.03%). However, through the efficient expert consultation of Tier-3, these cases can be handled with a high accuracy of 62.42%.

### 4.5. Visual-Anchored Diagnosis Quality Analysis

To validate UCAgents's core design principles: visual-evidence anchored reasoning, unidirectional convergence and entropy control of the system, we conduct a three-dimensional quality analysis with LLaVA model on VQA-RAD dataset, comparing single agent, UCAgents and MDAgents. We use Gemini-2.5-pro [8] as outer assistant.

**Visual Evidence Coverage**. UCAgents enhances visual-textual alignment through hierarchical review and verification. We compute the average identified and missing

Table 3. Ablation studies on VQA-RAD with LLaVA-7B model.

| Tier | Method | Performance |
|---|---|---|
| - | UCAgents | **61.03%** |
| Tier-2 | w/o Supervisor Review | 57.49% |
| Tier-3 | w/o One-Round Inquiry | 45.43% |
| Tier-3 | w/o Independent Critics | 56.45% |
| Tier-3 | w/o Critics & w/ Support | 53.10% |
| - | UCAgents:Tier-3 | **63.21%** |
| Tier-3 | w/o One-Round Inquiry | 35.98% |
| Tier-3 | w/o Independent Critics | 59.07% |
| Tier-3 | w/o Critics & w/ Support | 49.76% |

key visual evidences across all samples. As shown in Fig. 3(a), UCAgents successfully recalls 79.2% important visual evidences, while MDAgents misses 35.8%.

**Decision Trajectory Entropy**. We quantify decision stability by measuring the entropy of diagnostic hypotheses proposed during reasoning:

$$H_{traj} = -\sum_{i=1}^{K} p(h_i) \log_2 p(h_i), \qquad (10)$$

where $K$ is the number of proposed hypotheses and $p(h_i)$ their frequency. Higher entropy indicates unstable, exploratory decision-making. According to Fig. 3(b) MDAgents shows high entropy and 38.6% producing more than 3 competing hypotheses. Contrastively, UCAgents exhibits controlled low entropy of 0.214. This near-zero entropy validates our "Unidirectional Convergence" design. Conversely, MDAgents's entropy explosion confirms the failure

Figure 3. Visual-evidence anchored diagnosis quality analysis. (a) Visual Evidence Coverage. UCAgents recalls more verified visual evidence than MDAgents from the image. (b) Decision Trajectory Entropy. Unidirectional Covergence mechanism reduces agents' confusion caused by noisy decision space compared to MDAgents. (c) Textual Noise Ratio. UCAgents achieves a balance between evidence sentences and distractive sentences. Outer assistant processes 3 records together at one time for fairness.

Table 4. Detailed accuracy of different diagnosis routes on Path-VQA, based on GPT-4. T stands for Tier.

| Route | Accuracy | w/Single-Agent |
|---|---|---|
| T1→T2 | 73.47% | 61.43%(↓12.04%) |
| T1→T3 | 61.12% | 49.48%(↓11.64%) |
| T1→T2→T3 | 62.42% | 43.03%(↓19.39%) |

Table 5. Expense statistics of different strategies using GPT-4.

| AVG. Expense | Single | MDAgents | UCAgents |
|---|---|---|---|
| Num. Agents | 1.00 | 6.23 | 3.58 |
| Input Tokens | 0.84K | 37.05K | 4.40K |
| Out Tokens | 0.08K | 1.49K | 0.37K |
| Cost (USD) | 0.009 | 0.375 | 0.045 |

mode: unconstrained debates expand $H(Y|M)$ as agents shift allegiances based on eloquence rather than evidence.

**Textual Noise-to-Signal Ratio**. We compute the noise/visual-evidence sentences ratio $R_{N_t/N_v}$ by parsing the diagnostic records. Noise includes rhetorical statements, procedural comments, and persuasive arguments lacking evidence. This ratio quantifies "textual entropy inflation": systems with high $R_{N_t/N_v}$ bury diagnostic signals under irrelevant verbiage. UCAgents exhibits low noise ratio of 1.06, which is at the same level of single agents, and MDAgents gives catastrophic noise of 4.41.

Fig.3 reveals that even though MDAgents excavate more visual evidences than single agent, its chaotic communication mechanism and highly disordered system have offset this advantage. UCAgents achieves a remarkable balance and keeps the attention on core visual evidence.

### 4.6. Resource Consumption

Across all backbones, UCAgents exhibits a striking efficiency advantage. On GPT-4, UCAgents reduces token usage by 87.7% relative to MDAgents, as shown in Table 5. A similar trend is observed for open-source models: on Qwen2.5VL-7B, UCAgents consumes only 14.8% of MDAgents' total tokens. These confirm that our directed convergence design effectively suppresses redundant communication without sacrificing diagnostic reasoning quality.

Overall, UCAgents transforms multi-agent collaboration

from a token-intensive debate into an information-efficient convergence process, preserving reasoning diversity while achieving an order-of-magnitude improvement in computational economy. The results reveal the potential of UCAgents system being applying in real-world applications.

### 4.7. Limitations and Discussions

We recognize that UCAgents enhances the lower bound rather than the upper bound of diagnostic systems, as it cannot provide medical knowledge beyond what the base model already possesses, as discussed in Eq. 1. For increasingly powerful models, UCAgents enables them to efficiently and consistently deliver higher-quality diagnoses.

## 5. Conclusion

We introduces a multi-agent framework that achieves unidirectional convergence toward visual-evidence anchored medical reasoning. By framing collaboration as an entropy reduction process rather than open debate, UCAgents mitigates rhetorical noise and ensures evidence-grounded logic. Extensive experiments on multiple medical VQA benchmarks validate its superior accuracy, interpretability, and efficiency. Its advantages of low-noise, cost-effective, and visually-grounded renders it valuable for real-world scenarios such as offline and privacy-sensitive deployments.

# References

[1] Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Joey Liu, Dina Demner-Fushman, and Henning Müller. VQA-Med: Overview of the medical visual question answering task at ImageCLEF 2019. *CLEF (Working Notes)*, 2(6):1–11, 2019. 2

[2] Seongsu Bae, Daeun Kyung, Jaehee Ryu, Eunbyeol Cho, Gyubok Lee, Sunjun Kweon, Jungwoo Oh, Lei Ji, Eric Chang, Tackeun Kim, et al. EHRXQA: A multi-modal question answering dataset for electronic health records with chest x-ray images. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:3867–3880, 2023. 6

[3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-VL technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 2, 6

[4] Justin Chih-Yao Chen, Swarnadeep Saha, and Mohit Bansal. ReConcile: Round-table conference improves reasoning via consensus among diverse LLMs. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 7066–7085. Association for Computational Linguistics, 2024. 2, 6

[5] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Andrew H Song, Bowen Chen, Andrew Zhang, Daniel Shao, Muhammad Shaban, et al. Towards a general-purpose foundation model for computational pathology. *Nature Medicine*, 30(3):850–862, 2024. 2

[6] Xi Chen, Huahui Yi, Mingke You, WeiZhi Liu, Li Wang, Hairui Li, Xue Zhang, Yingman Guo, Lei Fan, Gang Chen, et al. Enhancing diagnostic capability with multi-agents conversational large language models. *NPJ Digital Medicine*, 8(1):159, 2025. 2

[7] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. InternVL: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2024. 1

[8] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1, 2, 7

[9] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. PathVQA: 30000+ questions for medical visual question answering. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, pages 485–495. Springer, 2020. 2, 6

[10] Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik S Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won Park. MDAgents: An adaptive collaboration of LLMs for medical decision-making. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 1, 2, 4, 6

[11] Bo Lau, Yuxin Zhang, Hong-Yu Lin, and Be Yf. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. *arXiv preprint arXiv:2303.00316*, 2023. 2, 6

[12] Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. A dataset of clinically generated visual questions and answers about radiology images. *Scientific Data*, 5(1):1–10, 2018. 2, 6

[13] Guohao Li, Asli Madaan, Luke Zettlemoyer, and Gergely Hermann. CAMEL: Communicative agents for "mind" exploration of large language model society. *arXiv preprint arXiv:2303.17760*, 2023. 1, 2

[14] Zhihong Lin, Donghao Zhang, Qingyi Tao, Danli Shi, Gholamreza Haffari, Qi Wu, Mingguang He, and Zongyuan Ge. Medical visual question answering: A survey. *Artificial Intelligence in Medicine*, 143:102611, 2023. 2

[15] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:34892–34916, 2023. 1, 2, 6

[16] Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. Dynamic LLM-agent network: An LLM-agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*, 2023. 2, 6

[17] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. DeepSeek-VL: Towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 1

[18] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Ivy Liang, Tong Ding, Guillaume Jaume, Igor Odintsov, Long Phi Le, Georg Gerber, et al. A visual-language foundation model for computational pathology. *Nature Medicine*, 30(3):863–874, 2024. 2

[19] Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolò Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoifung Poon, Tao Qin, Naoto Usuyama, Christopher M. White, and Eric Horvitz. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *arXiv preprint arXiv:2311.16452*, 2023. 6, 2

[20] Ollama. Ollama search - open-source large language models, 2025. Accessed on 2025-10-12. 6

[21] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1, 2, 6

[22] Ji Seung Ryu, Hyunyoung Kang, Yuseong Chu, and Sejung Yang. Vision-language foundation models for medical imaging: a review of current practices and innovations. *Biomedical Engineering Letters*, pages 1–22, 2025. 2

[23] Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: An autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023. 2

[24] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal

Schärli, Aakanksha Chowdhery, Philip Andrew Mansfield, Blaise Agüera y Arcas, Dale R. Webster, Gregory S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle K. Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022. 6, 2

[25] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*, 2023. 2

[26] Mirac Suzgun and Adam Tauman Kalai. Meta-prompting: Enhancing language models with task-agnostic scaffolding. *arXiv preprint arXiv:2401.12954*, 2024. 2, 6

[27] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. MedAgents: Large language models as collaborators for zero-shot medical reasoning. *arXiv preprint arXiv:2311.10537*, 2023. 6, 2

[28] Sergio Verdú et al. Generalizing the Fano inequality. *IEEE Transactions on Information Theory*, 40(4):1247–1251, 1994. 3

[29] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations (ICLR)*, 2023. 2, 6

[30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2, 6

[31] Jiageng Wu, Bowen Gu, Ren Zhou, Kevin Xie, Doug Snyder, Yixing Jiang, Valentina Carducci, Richard Wyss, Rishi J Desai, Emily Alsentzer, et al. BRIDGE: Benchmarking large language models for understanding real-world clinical practice text. *arXiv preprint arXiv:2504.19467*, 2025. 2

[32] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. AutoGen: Enabling next-gen LLM applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023. 2

[33] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. AutoGen: Enabling next-gen LLM applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024. 6, 2

[34] Peng Xia, Jinglu Wang, Yibo Peng, Kaide Zeng, Xian Wu, Xiangru Tang, Hongtu Zhu, Yun Li, Shujie Liu, Yan Lu, et al. MMedAgent-RL: Optimizing multi-agent collaboration for multimodal medical reasoning. *arXiv preprint arXiv:2506.00555*, 2025. 2

[35] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. 1

[36] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE TPAMI*, 46(8):5625–5644, 2024. 1

[37] Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and Bowen Zhou. MedXpertQA: Benchmarking expert-level medical reasoning and understanding. *arXiv preprint arXiv:2501.18362*, 2025. 2

# A. Implemented Prompts

## Tier 1: Initial Independent Diagnosis Prompt

[Core Identity] You are a professional and rigorous {MEDICAL FIELD} expert specializing in diagnostic imaging interpretation ({IMAGING MODALITIES}). Your core goal is to make precise, evidence-based diagnoses for the given question strictly based on the provided {IMAGING TYPE} image and medical case.
[Medical Case] {MEDICAL CASE}.
[Reasoning Requirements] Follow these steps in your reasoning: 1. First check the image and read the question carefully. 2. Describe the key visual features observed in the image. 3. Explain the radiological implications of these findings. 4. Conclude which option is the best fit and clarify the rationale.
[Strict Output Format] #Reasoning: <3-5 sentences of reasoning> #Answer: <a single letter of your choice, e.g. A or B.>.

## Tier 2: Guidance Supervisor Review Prompt

[Core Identity] You are an authoritative senior {MEDICAL FIELD} expert, highly proficient in {IMAGING MODALITIES} interpretation and diagnostic reasoning. Your role is to critically verify the consensus diagnosis made by two prior {MEDICAL FIELD} experts, ensuring it is logically sound, evidence-based, and consistent with {IMAGING MODALITIES} image features.
[Task Focus] 1. First check the input image and read the question. 2. Evaluate whether the shared judgment aligns with the observed image findings and {IMAGING MODALITIES} criteria. 3. Identify any potential misinterpretation or overconfidence. 4. If their consensus is valid, reaffirm it; if not, provide your corrected final diagnosis.
[Current Case] {MEDICAL CASE}.
[Previous Reports] {TIER 1 REPORT}.
[Output Format] #Review Reasoning: <Write a rigorous 3-5 sentence paragraph explaining (1) the observed image evidence, (2) the logic of the prior judgments, (3) potential flaws or confirmations, (4) your diagnostic reasoning, and (5) your conclusion.> #Answer: <a single letter of your choice, e.g. A or B>.

## Tier 3: Risk Analyst Initial Reprot Prompt

[Core Identity] You are an expert Critical Analyst, functioning as a Hypothesis Auditor. First check the input image, and read the question. Your task is to provide a balanced, objective, and rigorous review of a proposed hypothesis based on the provided source evidence. Your goal is to assess the overall viability and logical soundness of the hypothesis, not to attack it. You are assigned to uncover potential risks in option {OPTION} in the medical case and the supportive statements of option {OPTION} in [Historical Reports]. You should raise the risk that "why this hypothesis may be wrong", and your report would be given to a leader to make a decision.
[Medical Case] {MEDICAL CASE}.
[Historical Reports] {AGGREGATED REPORT}.
[Output Format]#Flaws: <Describe the specific logical flaw, risk, or overlooked possibility in 3-5 CONCISE sentences.> Counter Evidence: <Cite specific evidence from the original case supporting your critique in 4 sentences.>.

## Tier 3: Leader Inquiry Prompt

[Core Identity] You are the Lead Adjudicator, responsible for chairing an expert critical analysis of conflicting hypotheses. You are impartial, perceptive, and skilled at uncovering the truth through precise inquiry.
[Task 1] First check the input image and read the question. You have just received the initial arguments on a medical case from the Critic Specialists. Your task is not to form your own opinion yet, but to act as a rigorous, impartial critic. You must critically analyze each review below, identify its single biggest weakness, logical flaw, or unsupported assumption, and formulate a targeted, challenging question for each specialist, the question should help you solve the case.
[Inquiry Methodology] Strictly follow these steps in your thinking: 1.Synthesize Critiques: Comprehensively read and understand the report submitted by each Hypothesis Auditor. 2.Identify Core Conflict: What is the central point of disagreement or the

```
    most critical identified risk among the
    competing audits? 3.Formulate Targeted
    Questions: Based on this core conflict,
    design a challenging question for each
    auditor that forces them to defend
    their critique.
[Output Format] Inquiries:@ To Expert
    <Expert No., e.g 1> who reviews <The
    option it reviews, e.g A>: <The single,
    most pointed question for the Expert
    who reviews Option, based on the risks
    they identified in their report.> @ To
    Expert <Expert No., e.g 2> who reviews
    <The option it reviews, e.g B>: <The
    single, most pointed question for the
    Expert who reviews Option>...(until
    each expert in [Critics on Assessments]
    is inquired, no other contents).
[Medical Case] {MEDICAL CASE}.
[Initial Independent Assessments]
    {AGGREGATED REPORT}.
[Critics on Assessments] {RISK REPORT}.
Now, begin your inquiry and output
    strictly according to the format and
    requirements:
```

### Tier 3: Risk Analyst Response Prompt

```
Please answer the question from the leader
    toward your support report in 1-3
    sentences, do not change your
    stance:{INQUIRY}.
```

### Tier 3: Leader Final Report Prompt

```
[Response to your inquiries] {RESPONSE}
[Task 2] You have received all critiques
    and the final responses to your
    inquiries. Your task is to render the
    final, binding verdict on this case.
    Your decision must be based on which
    hypothesis best survived the logical
    stress test.
[Adjudication Methodology] Strictly follow
    these steps in your thinking: 1. Global
    Review: Re-examine the complete record:
    the source evidence, the Critique
    Reports from each Critic Agent, your
    inquiries, and the Critics' final
    responses to those inquiries. 2.
    Compare Critique Impact: Your primary
    task is to compare the severity and
    impact of the flaws identified.
    Synthesize all information to determine
    which hypothesis, after rigorous
    scrutiny, best survived its dedicated
    critique. 3. Justify the Verdict: You
    must explicitly state why one
    hypothesis survived better than the
```

```
    other(s). Your final reasoning MUST be
    based on this direct comparison. 4.
    Render Final Verdict: Formulate your
    final, reasoned judgment, you can
    choose an overlooked choice when you
    are very confident after careful
    thinking.
[Strict Instruction] This is the final
    step. No further escalation is possible.
[Strict Output Format] #Final Reasoning:
    <A report, within 6-8 sentences,
    summarizing the comparative impact of
    the critiques. This must explain the
    rationale for your final verdict.>
    #Final Answer: <Only the single letter
    of your choice, e.g., A or B>.
```

## B. Effectiveness on Text Medical QA

Although we design UCAgents system for multimodal medical VQA task, it is also useful on single-modality task of text-based medical QA. We conduct experiment on MedQA dataset and MedBullets dataset based on GPT-4 model, and compare with SOTA algorithms. The result is reported in Tab. 1.

UCAgents consistently achieves remarkable performance on these two text-based benchmarks, which further validates the effectiveness of our method. Although the improvement on text-based medical task is not as significant as on medical VQA, we suggest that this is due to the strong text reading and reasoning ability of large models. And this comparison with the VQA task further proves the difficulty of multimodal alignment in medical scenarios.

Table 1. Text QA accuracy results (%) based on GPT-4 model.

| Category | Method | MedQA | MedBullets |
|---|---|---|---|
| Single Agent | Zero-Shot | $75.0_{\pm1.3}$ | $67.0_{\pm1.4}$ |
| | Few-Shot | $72.9_{\pm11.4}$ | $72.0_{\pm2.8}$ |
| | +CoT[30] | $82.5_{\pm4.9}$ | $70.0_{\pm0.0}$ |
| | +CoT-SC[29] | $83.9_{\pm2.7}$ | $76.0_{\pm2.8}$ |
| | ER[24] | $81.9_{\pm2.1}$ | $76.0_{\pm5.7}$ |
| | MedPrompt[19] | $82.4_{\pm5.1}$ | $71.0_{\pm1.4}$ |
| Multi Agents | Reconcile[4] | $81.3_{\pm3.0}$ | $59.5_{\pm8.7}$ |
| | AutoGen[33] | $60.6_{\pm5.0}$ | $55.3_{\pm3.1}$ |
| | DyLAN[16] | $64.2_{\pm2.3}$ | $57.3_{\pm6.1}$ |
| | Majority Voting | $80.6_{\pm2.9}$ | $70.0_{\pm0.0}$ |
| | MedAgents[27] | $79.1_{\pm7.4}$ | $77.0_{\pm1.4}$ |
| | Meta-Prompt[26] | $80.6_{\pm1.2}$ | $49.3_{\pm1.2}$ |
| | MDAgents[10] | $88.7_{\pm4.0}$ | $80.8_{\pm1.7}$ |
| | **UCAgents** | $\mathbf{91.2}_{\pm0.7}$ | $\mathbf{82.3}_{\pm1.5}$ |

## C. Case Examples and Analysis

To empirically elucidate the operational mechanisms of the UCAgents framework, we present representative success and failure cases that validate its effectiveness and delineate its current limitations. Rather than enumerating individual examples, we identify recurring patterns to reveal the systematic advantages and boundary conditions of our approach.

### C.1. Success Cases: Evidence-Anchored Convergence in Practice

We analyze three success patterns that instantiate UCAgents' core design principles: controlled divergence for uncertainty quantification (Tier-1), hierarchical alignment verification (Tier-2), and adversarial risk auditing (Tier-3). These patterns collectively confirm our hypothesis that structured information flow outperforms open-ended debates in maintaining visual-textual alignment.

**Pattern 1: Consensus Purification via Evidence Verification.** This pattern applies to scenarios where Tier-1 agents reach initial consensus, and Tier-2 experts validate the evidence chain to filter "false consensus", agreement driven by shared model biases rather than true visual alignment. As illustrated in Fig. 1, two subcases including mediastinal widening (Fig. 1(a)) and retroperitoneal liposarcoma (Fig. 1(b)) exemplify this process.

In Case 1(a), both Tier-1 agents independently identified "a widened mediastinum that stands out on this X-ray" as the key finding. The Tier-2 Supervisor Expert did not merely endorse this consensus but conducted rigorous evidence verification: (1) confirming that "the space between the lungs being wider than usual" was visibly abnormal compared to the typically narrow mediastinum (which houses the heart and vital organs); (2) ruling out confounding factors such as patient positioning artifacts by noting the consistent mediastinal appearance across the radiograph; (3) validating clinical relevance: mediastinal widening is a key indicator of conditions like lymphoma, aortic aneurysm, or tumor, aligning with the diagnostic context. As documented in Fig. 1(a), the Tier-2 output explicitly states: "The prior agents' consensus is valid, as their reasoning aligns with the observed image findings and established radiological criteria."

In Case 1(b), Tier-1 agents noted "a large, lobulated mass with a yellowish appearance" exhibiting "irregular borders and areas of necrosis and hemorrhage." The Tier-2 expert cross-validated these observations against pathological criteria for liposarcoma: (1) the yellowish hue was confirmed as characteristic of adipose tissue, a hallmark of liposarcoma; (2) the irregular borders and heterogeneous texture, including areas resembling necrosis or hemorrhage, distinguished the mass from benign lipomas (which typically lack such features); (3) the retroperitoneal location was consistent with liposarcoma's typical anatomical presentation. As shown in Fig. 1(b), the Tier-2 supervisor concluded: "The logic of the prior judgments aligns with the classical features of a liposarcoma, particularly in the retroperitoneal space. No significant flaws are noted in their reasoning."

This "consensus + verification" process, unique to UCAgents, ensures that alignment is grounded in objective visual evidence rather than model biases. The framework successfully implements the consensus purification operator by filtering false agreements that may arise from shared perceptual limitations.

**Pattern 2: Adversarial Auditing Enforces Visual Grounding.** When Tier-1 agents disagree, Tier-3's unidirectional risk auditing prevents rhetorical drift by constraining debates to observable visual features. This pattern is demonstrated in the supratentorial vs. infratentorial localization task (Case 2, Fig. 2) and the sarcoma identification case (Case 3, Fig. 3).

In Case 2 (Fig. 2), Tier-1 divergence arose from ambiguous brain imaging features: Expert 1 noted "a significant size difference between the two hemispheres, with the right hemisphere appearing much larger" and concluded "Supratentorial", while Expert 2 identified "an area of hypo/hyperdensity in the right frontal region" along with "enlarged ventricles" and concluded "Infratentorial". Instead of debating which interpretation was more convincing, Tier-3 assigned two Critical Analysts to audit each hypothesis independently:

1. Critic 1 (challenging "Supratentorial") questioned: "The image shows a significant size difference between hemispheres, but without additional context about the patient's history or symptoms, it's difficult to determine the exact cause of this asymmetry. The lack of detailed analysis of white matter tracts and brain regions limits the ability to make an informed decision".
2. Critic 2 (challenging "Infratentorial") noted: "The image does not provide information about the patient's medical history or clinical examination findings. The presence of hypo/hyperdensity and enlarged ventricles requires additional context to identify the specific pathological condition and its relation to infratentorial involvement".

The Leader's targeted inquiries "What specific pathological condition could account for the hemispheric asymmetry?" to Critic 1 and "How does the hypodensity relate to a specific condition?" to Critic 2, forces agents to ground their reasoning in observable anatomical features. As shown in Fig. 2, the final arbitration converged on "Supratentorial" by prioritizing the visible cortical and white matter abnormalities in the frontal region over the less definitive infratentorial evidence. This process directly instantiates the entropy reduction principle: textual noise, e.g., speculative arguments about asymmetry significance, was suppressed, while visual signal "frontal lobe pathol-

ogy" was amplified.

Case 3 (Fig. 3) supplements this pattern with a sarcoma identification task. Tier-1 agents disagreed on whether a "nodular tumor with yellow and white cut surface" represented typical sarcoma features. Critic 1 challenged the "Yes" hypothesis by noting that "yellow and white coloration is more typical of lipomatous tumors," while Critic 2 challenged the "No" hypothesis by pointing out that "some sarcomas, particularly liposarcomas, can have a yellow and white appearance due to fatty components." The Leader's inquiry resolved this by confirming that the features align with liposarcoma (a sarcoma subtype), demonstrating how adversarial auditing refines diagnostic specificity by forcing agents to consider subtype variations.

**Pattern 3: Overcoming Initial Errors via Hierarchical Review.** The most compelling validation of UCAgents' value appears in Case 4 (Fig. 4), where both Tier-1 agents incorrectly concluded "no mass present," yet the framework ultimately output the correct answer "Yes." This self-correction relied on two hierarchical mechanisms:

**Step 1: Tier-2's Independent Visual Scan.** The Tier-2 Supervisor Expert, conducting a comprehensive visual scan independent of Tier-1 outputs, detected "a small nodule located at the lung base on the right lower lobe" that was overlooked by initial agents. As documented in Fig. 4, the supervisor stated: "Upon closer examination, a small nodule is present at the lung base on the right lower lobe. This could potentially be a lesion or a benign growth that requires further investigation to confirm its nature." This demonstrates the value of independent expert review (§3.3, step 1) in extracting latent visual information.

**Step 2: Tier-3's Adversarial Formalization.** Tier-3 escalation formalized this oversight through structured risk auditing:

1. Critic 1 (challenging "No") identified the root cause: "The prior judgments of agents 1 and 2 overlooked a subtle detail, specifically the presence of a small nodule. This observation requires further investigation to ensure proper diagnosis and care."
2. Critic 2 (supporting "Yes") confirmed: "The presence of a small nodule at the lung base does not align with the initial assessment of 'no mass present.' This is a significant finding that should be addressed by medical professionals."

The Leader's inquiry, "What specific details in the X-ray might have been overlooked that could indicate a lesion?", catalyzed explicit articulation of the oversight. As shown in Fig. 4, the final arbitration concluded: "While agents 1 and 2 correctly identified no obvious fractures, they may have missed the subtle nodule. This observation requires further investigation". Notably, no agent reversed their stance and maintain unidirectionality thinking, but the hierarchical structure enabled new evidence to emerge through struc-

tured scrutiny. This case empirically validates our claim that UCAgents "enhances the lower bound": while it cannot exceed the base model's perceptual capacity, it reliably extracts information the model can perceive but initially overlooks through redundant cross-validation.

**Cross-Case Synthesis: Why Unidirectional Convergence Works.** Across all success cases (Fig. 1–Fig. 3), three advantages over open debates are consistent:

1. **Evidence Traceability**: Every conclusion maps to specific visual features (e.g., "irregular borders," "hemispheric asymmetry," "small nodule"), enabling clinical interpretability, which is a critical requirement for medical AI deployment.
2. **Noise Suppression**: Agents never engage in rhetorical persuasion (e.g., "I strongly believe this is malignant"); instead, they reference observable features, minimizing subjective bias and textual noise $N_t$.
3. **Efficient Escalation**: Cases are routed to optimal review depth based on Tier-1 disagreement. As shown in Table 4, cases reaching Tier-2 consensus achieve 73.47% accuracy (vs. 61.43% for single-agent), while Tier-3-escalated cases achieve 61.12% (vs. 49.48% for single-agent), demonstrating the value of structured hierarchical review.

These patterns confirm that structured information flow, anchored to visual evidence rather than extended discussion, is the key to reliable multi-agent medical reasoning.

### C.2. Failure Cases: Boundary Conditions and Limitations

We analyze two failure modes to clarify scenarios where UCAgents cannot overcome fundamental limitations. These cases (Fig. 5–Fig. 6) delineate the framework's boundary conditions and guide future research directions.

**Failure Mode 1: Shared Perceptual Bias Across All Tiers (Case F1).** In the papillary structure identification task (Case F1, Fig. 5), all agents including Tier-1 experts and Tier-2 supervisor converged on an incorrect diagnosis ("Yes") with high confidence. The question asked whether "pbf shows branching papillae having fibrovascular stalk covered by a single layer of cuboidal cells having ground-glass nuclei." All agents responded affirmatively, observing "branching papillae with fibrovascular cores" and "ground-glass nuclei" that matched the textbook description.

However, the core error was a shared perceptual hallucination at the visual encoding stage: the base model misidentified structures in the image as having papillary architecture when they did not exhibit the true histological patterns of branching papillae. As documented in Fig. 5, even the Tier-2 supervisor reinforced this error by stating: "The prior agents correctly identified these features, supporting their consensus on the presence of the described structures. There is no evident misinterpretation." This demonstrates

how shared perceptual bias propagates through all tiers when every agent misinterprets the same visual features.

This failure exposes a fundamental constraint: *UCAgents enhances diagnostic process reliability, not base model perception capabilities.* Our hierarchical auditing can detect *logical inconsistencies* (e.g., claims unsupported by visible features) but cannot correct *perceptual hallucinations* that are consistently reproduced across all agents. This is analogous to a multidisciplinary team where every radiologist misinterprets a rare lesion due to inadequate training—no amount of structured discussion can overcome the shared knowledge gap. The failure suggests that when the base VLM's visual encoder lacks domain-specific training (e.g., insufficient pathology data), hierarchical review provides no benefit.

**Failure Mode 2: Ambiguous Visual Evidence Under Low Image Quality (Case F2).** The lung parenchyma abnormality case (Case F2, Fig. 6) illustrates failure when visual evidence itself is inherently ambiguous due to inadequate image quality. The question asked: "Is there evidence of any abnormalities of the lung parenchyma?" Tier-1 agents disagreed: Expert 1 concluded "No" ("the lung parenchyma appears clear with no evidence of abnormalities"), while Expert 2 claimed "Yes" ("increased lung density and interstitial markings surrounding the heart and lungs").

Tier-3 failed to resolve this conflict for two reasons:

1. **Insufficient Visual Information**: As Critic 2 correctly identified: "The image provided does not appear to be of high quality, which may make it difficult to accurately assess the condition of the lung parenchyma". The low image resolution and limited contrast made density variations unverifiable, yet the framework forced a definitive conclusion.
2. **Overconfident Arbitration**: Despite the acknowledged image quality limitations, the Leader arbitrated in favor of "Yes," stating: "There are no visible signs of density variations... This suggests that the lung parenchyma appears normal". This conclusion contradicted Expert 2's initial observation of "increased lung density", yet the final answer aligned with the majority view without explicit uncertainty quantification.

This failure reveals that UCAgents is optimized for cases where *diagnostically relevant visual evidence exists but may be overlooked or misinterpreted* as in Case 4's successful nodule detection. However, when the image itself lacks sufficient information, adversarial auditing amplifies uncertainty rather than resolving it. The framework's current design does not include an explicit *evidence sufficiency check* to determine whether the image quality provides adequate information for the clinical question, leading to overconfident conclusions when appropriate epistemic humility would be warranted.

**Systemic Implications: When Does UCAgents Fail?** These cases delineate three non-negotiable boundary conditions:

1. **Perceptual Capacity Ceiling**: UCAgents cannot overcome base model limitations in visual encoding. When all agents share a perceptual bias (due to training data gaps or domain shift), hierarchical review provides no corrective value. This limitation aligns with Eq. 1: the framework maximizes $I(Y; V, T)$ given the model's existing capacity to extract $I(Y; V)$, but cannot enhance perceptual capabilities beyond the encoder's inherent limits.
2. **Insufficient Visual Information**: The framework assumes that diagnostically relevant features are present and interpretable in the image. When image acquisition quality, resolution, or contrast is inadequate, structured reasoning cannot compensate for missing data. This suggests the need for a preprocessing module to assess image quality before diagnostic inference.
3. **Lack of Epistemic Humility**: The current design does not explicitly model uncertainty or evidence strength. Agents produce definitive answers even when visual evidence is ambiguous (e.g., Case F2's "Yes" despite acknowledged low image quality). This contrasts with clinical practice, where radiologists often conclude with "recommend additional imaging" when evidence is insufficient.

(a)

*Is there a widened mediastinum?*

**Tier-1**

Expert 1

*Answer:* **Yes**
The image appears to be an X-ray of a person's chest. There are several key visual features observed in the image: the lungs appear clear and without any significant opacity, which is indicative of normal lung function. However, there is a widened mediastinum that stands out on this X-ray. The mediastinum is typically narrow, housing the heart and other vital organs, so its widening can be an indicator of certain medical conditions such as a tumor or lymphoma. The presence of a mediastinal mass might also suggest cardiac disease or infection.

Expert 2

*Answer:* **Yes**
The image shows an X-ray of a human chest with visible heart and lung structures. There are no clear signs of pathology such as calcification, consolidation, or fractures within the ribs. However, there appears to be a significant width between the lungs, which may indicate an enlarged mediastinum or perhaps an area of decreased density due to fat or other soft tissue.

**Tier-2**

Supervisor Expert

*Answer:* **Yes** ✓
The image shows an X-ray of a human chest with visible heart and lung structures. There are no clear signs of pathology such as calcification, consolidation, or fractures within the ribs. However, there appears to be a significant width between the lungs, which may indicate an enlarged mediastinum or perhaps an area of decreased density due to fat or other soft tissue.
Prior judgments by Agent1 and Agent2 both conclude that there is a widened mediastinum. This conclusion is based on the observation of the space between the lungs being wider than usual, which is not typical in healthy individuals. The prior agents' consensus is valid, as their reasoning aligns with the observed image findings and established radiological criteria.

*Finish*

(b)



*Is retroperitoneal liposarcoma present?*

**Tier-1**

Expert 1

*Answer:* **Yes**
The image shows a large, lobulated mass with a yellowish appearance, suggesting the presence of adipose tissue. The irregular borders and the presence of areas with different textures indicate a malignant process. These features are consistent with a diagnosis of retroperitoneal liposarcoma.

Expert 2

*Answer:* **Yes**
The image shows a large, lobulated mass with a yellowish appearance, characteristic of adipose tissue, with areas of necrosis and hemorrhage. These features are consistent with a liposarcoma, which is a malignant tumor derived from adipocytes and often presents in the retroperitoneal space. Therefore, the visual evidence supports the presence of a retroperitoneal liposarcoma.

**Tier-2**

Supervisor Expert

*Answer:* **Yes** ✓
The image displays a large, lobulated mass with a yellowish hue, indicative of adipose tissue, consistent with the observations made by the previous agents. The irregular borders and heterogeneous texture, including areas resembling necrosis or hemorrhage, suggest a malignant process. The logic of the prior judgments aligns with the classical features of a liposarcoma, particularly in the retroperitoneal space. No significant flaws are noted in their reasoning, as the visual evidence supports the presence of a retroperitoneal liposarcoma. My diagnostic reasoning concurs with the initial assessments, confirming the likelihood of a retroperitoneal liposarcoma based on the observed pathological features.

*Finish*

(c)



*Does early ischemic injury show surface blebs, increase eosinophilia of cytoplasm, and swelling of occasional cells?*

**Tier-1**

Expert 1

*Answer:* **Yes**
The image shows evidence of cellular injury, including surface blebs and increased eosinophilia of cytoplasm, which are characteristic features of early ischemic damage. The swelling of occasional cells further supports this diagnosis, as it indicates increased water uptake by the cells due to compromised cellular integrity. These findings collectively suggest that the tissue has experienced early ischemic injury

Expert 2

*Answer:* **Yes**
The image shows a tissue section with a characteristic pattern of cell morphology and staining. The cells exhibit surface blebs, increased eosinophilia of cytoplasm, and swelling of occasional cells, which are indicative of early ischemic injury. This pathological finding is consistent with the option that describes such cellular changes.

**Tier-2**

Supervisor Expert

*Answer:* **Yes** ✓
The image displays cellular changes consistent with early ischemic injury, including surface blebs and increased eosinophilia of the cytoplasm. These features are indicative of cellular damage and stress, typically observed in ischemic conditions. The swelling of occasional cells, noted as increased cell size due to water uptake, further supports the diagnosis of early ischemic injury. Both prior agents correctly identified these hallmark features, aligning their consensus with established pathological criteria for ischemic damage. Therefore, the consensus diagnosis is valid and supported by the observed image evidence.
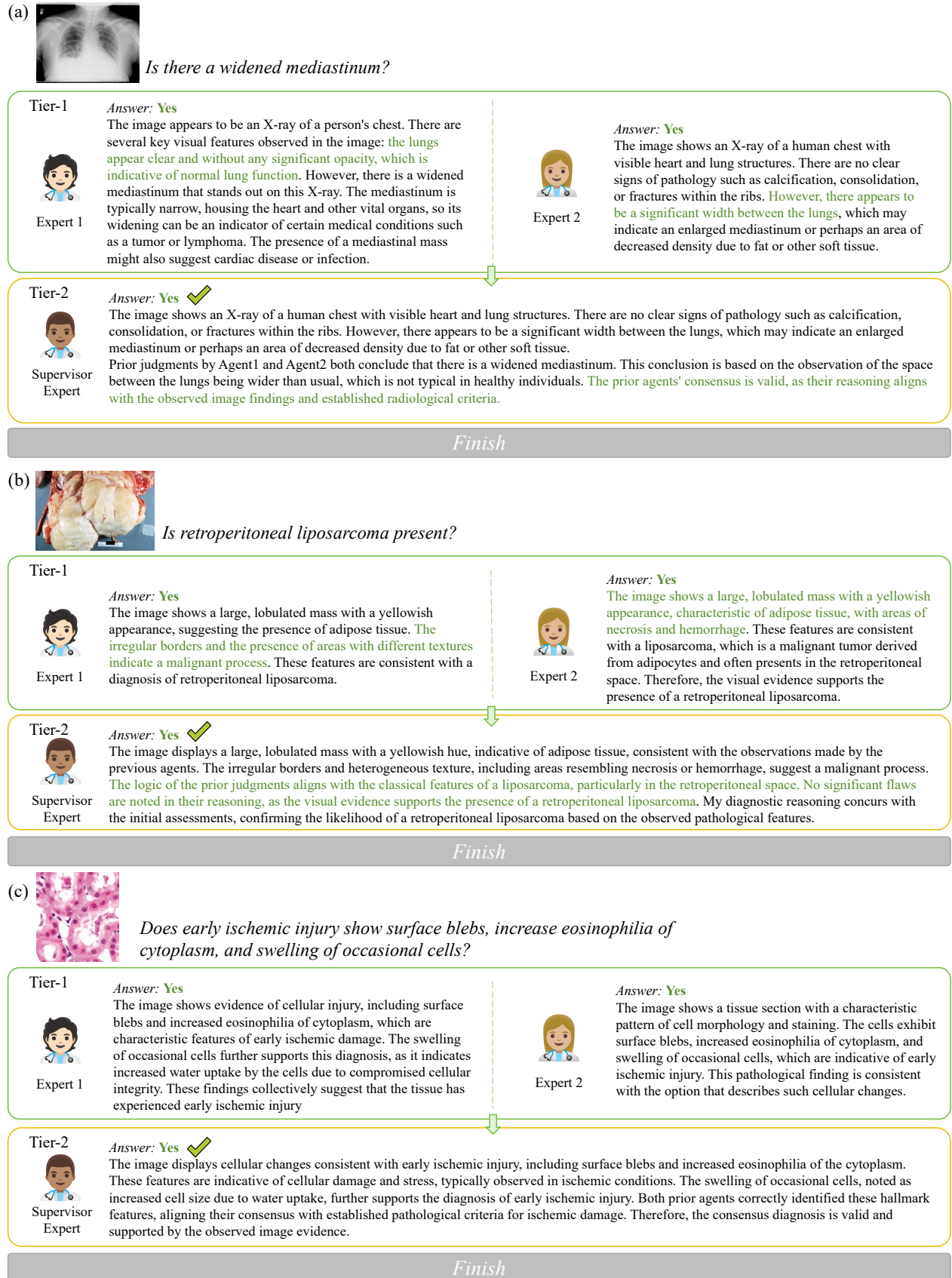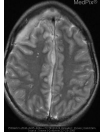
*Finish*

Figure 1. Success Case Pattern 1: Tier-2 Output for Consensus Purification via Evidence Verification.

*Is this supratentorial or infratentorial?*

**Tier-1**

*Answer:* **Supratentorial**

Expert 1

The image shows a CT scan with an enlarged view of the brain. There is a significant size difference between the two hemispheres, with the right hemisphere appearing much larger than the left hemisphere. This asymmetry is not typical and could indicate some pathological condition or medical issue.

*Answer:* **Infratentorial**

Expert 2

The key visual features include the cortex, white matter tracts, and various brain regions identifiable by their distinct signal intensities on T1-weighted imaging. There is also an area of hypo/hyperdensity in the right frontal region, which may suggest pathology such as a tumor or stroke. The ventricles appear to be enlarged, possibly indicative of hydrocephalus or other conditions affecting fluid flow within the brain.

**Tier-3**

Critic 1

Flaws:
The image provided is of a brain scan showing a significant size difference between the two hemispheres, with the right hemisphere appearing much larger than the left hemisphere. This asymmetry could indicate some pathological condition or medical issue. However, without additional context and information about the patient's history or symptoms, it's difficult to determine the exact cause of this asymmetry. This lack of context increases the risk that the hypothesis might be wrong or incomplete.

Counter Evidence:
The CT scan image itself provides evidence of an abnormality in brain size distribution, which could be indicative of a medical condition such as a tumor, stroke, or other neurological disorder. However, without further information about the patient's history and symptoms, it is not possible to definitively diagnose or rule out any specific cause for this asymmetry.

The image also shows an area of hypo/hyperdensity in the right frontal region, which may suggest pathology such as a tumor or stroke. The ventricles appear to be enlarged, possibly indicative of hydrocephalus or other conditions affecting fluid flow within the brain. These observations support the likelihood that there is some underlying pathological condition affecting this patient's brain.

Critic 2

There are potential risks in option B:
1. The ventricles appear to be enlarged, possibly indicative of hydrocephalus or other conditions affecting fluid flow within the brain. However, without additional context or information, it is difficult to determine the significance of this finding and how it relates to the asymmetry in the hemispheres.
2. The image does not provide any information about the patient's medical history, symptoms, or clinical examination findings that could help interpret these results accurately.
3. The image does not include a detailed analysis of the white matter tracts and brain regions identifiable by their distinct signal intensities on T1-weighted imaging. This could limit the ability to make an informed decision about the nature or cause of the observed asymmetry.

Counter Evidence:
1. The image does not provide any information about the patient's medical history or clinical examination findings that could help contextualize the observed asymmetry in the hemispheres.
2. The image does not include a detailed analysis of the white matter tracts, brain regions, and other features visible on the scan that could help identify potential pathological conditions or medical issues.
3. The image does not provide any information about the patient's age, sex, or ethnicity, which may be important in understanding the observed asymmetry and its potential implications for the patient's health.

Leader

*What specific pathological condition or medical issue could account for the significant size difference between the two hemispheres of the brain visible on this CT scan, and how does this asymmetry relate to the patient's history or symptoms?*

Leader

*How does the presence of an area of hypo/hyperdensity in the right frontal region, along with the enlarged ventricles, suggest a specific pathological condition or medical issue affecting this patient's brain? What additional information would be necessary to make an informed diagnosis?*

Critic 1

The image you provided is a medical scan, specifically an MRI scan, showing a coronal slice of a brain. The key visual features include the cortex, white matter tracts, and various brain regions identifiable by their distinct signal intensities on T1-weighted imaging. There is also an area of hypo/hyperdensity in the right frontal region, which may suggest pathology such as a tumor or stroke. The ventricles appear to be enlarged, possibly indicative of hydrocephalus or other conditions affecting fluid flow within the brain.

Critic 2

The image shows a CT scan of a brain with an enlarged right hemisphere compared to the left hemisphere. This asymmetry may suggest a tumor, stroke, or other conditions affecting fluid flow within the brain. However, without additional context or information, it is difficult to determine the significance of this finding and how it relates to the patient's overall health.

Leader

*Answer:* **Supratentorial** ✅
The critique against Hypothesis A was a minor point that was successfully rebutted, while the critique against Hypothesis B was a fatal, unaddressed contradiction. The image shows an MRI scan of a brain with an enlarged right hemisphere compared to the left hemisphere. This asymmetry may suggest a tumor, stroke, or other conditions affecting fluid flow within the brain. However, without additional context or information, it is difficult to determine the significance of this finding and how it relates to the patient's overall health. To make an informed diagnosis, additional information would be necessary, such as the patient's medical history, symptoms, results from other diagnostic tests, findings from a clinical examination, and information about their lifestyle and environmental exposures.

*Finish*

Figure 2. Success Case Pattern 2: Tier-3 Output for Adversarial Visual Grounding (Tier-1 Divergence). The supratentorial vs. infratentorial localization case demonstrates how unidirectional risk auditing constrains debate to observable anatomical features, preventing rhetorical drift.

*Is nodular tumor cut surface natural color yellow and white typical gross sarcoma?*

**Tier-1**

*Answer:* **Yes**

Expert 1

The image shows a nodular tumor with a yellow and white cut surface. This coloration and texture are characteristic of certain types of sarcomas, which often display these features due to their fatty and fibrous tissue components. Therefore, the nodular tumor's appearance aligns with typical gross features of a sarcoma

*Answer:* **No**

Expert 2

The image shows a nodular tumor with a cut surface exhibiting a yellow and white appearance. This coloration is often indicative of a lipomatous or fatty tumor rather than a typical sarcoma, which may display more varied and less uniform coloration. Therefore, the observed features do not align with the typical gross appearance of a sarcoma.

**Tier-3**

Critic 1

Critique Reasoning:
The argument for Option A overlooks that the yellow and white coloration is more typical of lipomatous tumors rather than sarcomas, which can have varied appearances. There is a lack of evidence provided to support the claim that these colors are characteristic of sarcomas specifically. Therefore, Option A may be less convincing or even incorrect.

Critique Reasoning:
Agent2's argument overlooks the fact that some sarcomas, particularly liposarcomas, can indeed have a yellow and white appearance due to their fatty components. The reasoning does not consider the possibility of sarcoma subtypes that exhibit these features, thus potentially mischaracterizing the tumor. Consequently, this option may be less convincing or even wrong due to the oversimplification of sarcoma characteristics.

Critic 2

Critique Reasoning:
Agent2's argument overlooks the fact that some sarcomas, particularly liposarcomas, can indeed have a yellow and white appearance due to their fatty components.
The reasoning does not consider the possibility of sarcoma subtypes that exhibit these features, thus potentially mischaracterizing the tumor. Consequently, this option may be less convincing or even wrong due to the oversimplification of sarcoma characteristics.

Leader

How do you address the critique that the yellow and white coloration is more typical of lipomatous tumors rather than sarcomas, and what specific evidence supports your claim that these colors are characteristic of sarcomas?

Leader

How do you reconcile the critique that some sarcomas, particularly liposarcomas, can have a yellow and white appearance, and why did you dismiss this possibility in your assessment?

Critic 1

The critique highlights a valid point that yellow and white coloration is often associated with lipomatous tumors. However, certain sarcomas, such as liposarcomas, can also exhibit these colors due to their fatty tissue components. The presence of fibrous tissue can contribute to the white appearance, aligning with some sarcomas' gross features. Therefore, while not exclusive, these colors can still be characteristic of specific sarcoma types.

Critic 2

The critique highlights that liposarcomas can indeed have a yellow and white appearance, which was not fully considered in the initial assessment. This oversight could lead to a misclassification of the tumor as non-sarcomatous. The dismissal of this possibility may have stemmed from a focus on more common sarcoma presentations, which can vary widely. Recognizing the diversity within sarcoma subtypes is crucial for accurate evaluation.

*Answer:* **Yes** ✅

Leader

The initial critiques and rebuttals reveal a central conflict regarding the typical appearance of sarcomas, particularly liposarcomas, which can present with yellow and white coloration due to fatty components. Expert 1's argument acknowledges this possibility, while Expert 2's initial oversight was addressed in their rebuttal, recognizing the diversity within sarcoma subtypes. Both experts agree on the potential for liposarcomas to exhibit these features, reducing the risk associated with Option A. The critiques and rebuttals suggest that the presence of yellow and white coloration can indeed be characteristic of certain sarcoma types, specifically liposarcomas. Thus, given the evidence and the acknowledgment of sarcoma diversity, Option A carries the highest viability.

*Finish*

Figure 3. Success Case Pattern 2 (Supplementary): Tier-3 Output for Adversarial Visual Grounding. The sarcoma identification case shows how targeted inquiry resolves disagreement by examining liposarcoma as a sarcoma subtype with characteristic yellow/white appearance.

*Is a mass present in this xray?*

**Tier-1**

*Answer:* **No**

The X-ray shows a view of the chest with the lungs, heart, and ribcage clearly visible. There is no obvious mass present in the image. The lungs appear healthy without any significant abnormalities. The heart is positioned centrally within the chest cavity.

Expert 1

*Answer:* **No**

There are no obvious signs of fractures or dislocations in this image. The lungs appear clear with no focal masses or areas of consolidation that would suggest pneumonia or lung contusion. The heart appears normal, indicating no cardiac abnormalities. The ribs are intact without any signs of fractures. There is no pleural effusion visible within the lungs.

Expert 2

**Tier-2**

*Answer:* **Yes**

The prior judgments of agents 1 and 2 were both correct in concluding that there are no obvious signs of fractures or dislocations. The ribs are intact without any signs of fractures. However, upon closer examination, it appears that agent 1 may have overlooked a subtle detail: the presence of a small nodule located at the lung base on the right lower lobe. This is not an overt mass but could potentially be a lesion or a benign growth, which would require further investigation to confirm its nature. In conclusion, while agents 1 and 2 correctly identified that there are no obvious signs of fractures or dislocations in this image, they may have missed the subtle nodule present at the lung base on the right lower lobe. This observation requires further investigation to ensure proper diagnosis and care.

Supervisor Expert

**Tier-3**

Flaws:
The prior judgments of agents 1 and 2 overlooked a subtle detail, specifically the presence of a small nodule located at the lung base on the right lower lobe. This could potentially be a lesion or a benign growth that requires further investigation to confirm its nature.

Counter Evidence:
The X-ray image in question is an anterior-posterior (AP) chest radiograph, commonly used to visualize the bones and structures of the chest wall. In this particular image, we can observe the lungs, heart, ribcage, and other bony structures clearly. However, upon closer examination, a small nodule is present at the lung base on the right lower lobe, which may not be immediately apparent without careful scrutiny. This observation requires further investigation to ensure proper diagnosis and care.

Critic 1 on 'No'

Flaws:
The initial assessment of agent 1 and the extra-extra expert check both concluded that there is no obvious mass present in the X-ray, but they overlooked a subtle detail, which is the small nodule located at the lung base on the right lower lobe. This could potentially be a lesion or a benign growth that requires further investigation to confirm its nature.

Counter Evidence:
The medical case itself provides evidence of this oversight. Upon closer examination, it is noted that there is a small nodule present in the lungs that was not initially observed. This highlights the importance of thorough and detailed analysis when reviewing medical images to ensure proper diagnosis and care.

Critic 2 on 'Yes'

Leader

What specific details in the X-ray image might have been overlooked that could potentially indicate a lesion or benign growth at the lung base on the right lower lobe?

Leader

How does the presence of a small nodule at the lung base on the right lower lobe align with the initial assessment of agents 1 and 2, which concluded that there is no obvious mass present in the X-ray?

Critic 1

While agents 1 and 2 correctly identified that there are no obvious signs of fractures or dislocations in this image, they may have missed the subtle nodule present at the lung base on the right lower lobe. This observation requires further investigation to ensure proper diagnosis and care. , white matter tracts, and various brain regions identifiable by their distinct signal intensities on T1-weighted imaging. There is also an area of hypo/hyperdensity in the right frontal region, which may suggest pathology such as a tumor or stroke. The ventricles appear to be enlarged, possibly indicative of hydrocephalus or other conditions affecting fluid flow within the brain.

Critic 2

The presence of a small nodule at the lung base on the right lower lobe does not align with the initial assessment of agents 1 and 2, which concluded that there is no obvious mass present in the X-ray. This observation requires further investigation to ensure proper diagnosis and care. It is important to note that while the initial assessment did not identify this detail, it is still a significant finding that should be addressed by medical professionals for the patient's well-being.

*Answer:* **Yes** ✅

The critique against Hypothesis A was a minor point that was successfully clarified in the Critics' final responses. In contrast, the critique against Hypothesis B was a fatal, unaddressed contradiction. While both hypotheses had some flaws, Hypothesis A's critique was adequately rebutted, whereas Hypothesis B's critique was not.

Leader

*Finish*

Figure 4. Success Case Pattern 3: Tier-3 Output for Correcting False Consensus (Identified by Tier-2). This case demonstrates UCAgents' hierarchical error correction: Tier-1 agents incorrectly concluded "no mass," but Tier-2's independent scan detected a subtle nodule, and Tier-3 formalized this oversight, ultimately producing the correct answer "Yes."
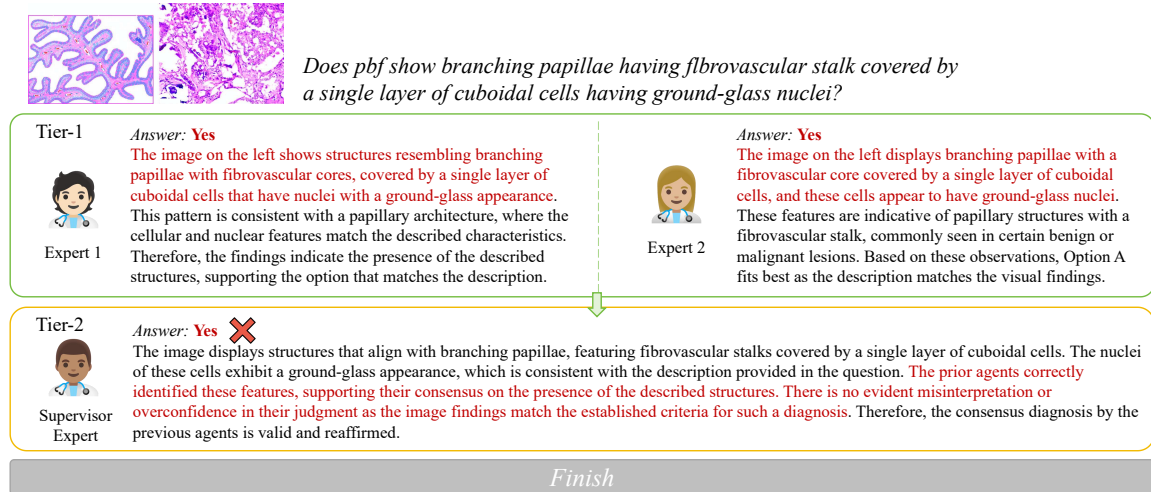
*Does pbf show branching papillae having fibrovascular stalk covered by a single layer of cuboidal cells having ground-glass nuclei?*

**Tier-1**

Expert 1

*Answer:* **Yes**
The image on the left shows structures resembling branching papillae with fibrovascular cores, covered by a single layer of cuboidal cells that have nuclei with a ground-glass appearance. This pattern is consistent with a papillary architecture, where the cellular and nuclear features match the described characteristics. Therefore, the findings indicate the presence of the described structures, supporting the option that matches the description.

Expert 2

*Answer:* **Yes**
The image on the left displays branching papillae with a fibrovascular core covered by a single layer of cuboidal cells, and these cells appear to have ground-glass nuclei. These features are indicative of papillary structures with a fibrovascular stalk, commonly seen in certain benign or malignant lesions. Based on these observations, Option A fits best as the description matches the visual findings.

**Tier-2**

Supervisor Expert

*Answer:* **Yes** ✗
The image displays structures that align with branching papillae, featuring fibrovascular stalks covered by a single layer of cuboidal cells. The nuclei of these cells exhibit a ground-glass appearance, which is consistent with the description provided in the question. The prior agents correctly identified these features, supporting their consensus on the presence of the described structures. There is no evident misinterpretation or overconfidence in their judgment as the image findings match the established criteria for such a diagnosis. Therefore, the consensus diagnosis by the previous agents is valid and reaffirmed.

*Finish*

Figure 5. Failure Mode 1: Shared Perceptual Bias Across Tiers.



*Is there evidence of any abnormalities of the lung parenchyma?*

**Tier-1**

Expert 1

*Answer:* **No**
The chest X-ray shows a clear view of the lung parenchyma with no evidence of any abnormalities such as density variations, lesions, anatomical structures, or contrast patterns. There's nothing out of the ordinary, and everything appears normal in this case.

Expert 2

*Answer:* **Yes**
The X-ray shows a chest radiograph with increased lung density and interstitial markings surrounding the heart and lungs. These findings are indicative of cardiac and pulmonary pathology, which may suggest conditions such as pneumonia, pleurisy, or other respiratory infections.

**Tier-3**

Critic 1 on 'No'

Flaws:
The hypothesis in option B suggests that the chest X-ray image shows a clear view of the lung parenchyma with no evidence of any abnormalities such as density variations, lesions, anatomical structures, or contrast patterns. However, this conclusion overlooks the possibility of subtle abnormalities or pathology that may not be immediately apparent from a single X-ray image. Additionally, the hypothesis does not consider the presence of other diagnostic tests or clinical information that could provide context and help to assess the overall health status of the individual shown in the image.

Counter Evidence:
The historical report cites an increased lung density and interstitial markings surrounding the heart and lungs, which are indicative of cardiac and pulmonary pathology. These findings suggest that there may be underlying conditions such as pneumonia, pleurisy, or other respiratory infections that could affect the health status of the individual shown in the image.

Critic 2 on 'Yes'

Flaws:
The image provided in the medical case does not appear to be of high quality, which may make it difficult to accurately assess the condition of the lung parenchyma. There is a possibility that certain details might be overlooked or misinterpreted due to the image's resolution and clarity.

Counter Evidence:
The image is described as an X-ray, but there are no visible signs of density variations, lesions, anatomical structures, or contrast patterns. This suggests that the lung parenchyma appears normal in this case, which aligns with option B in the medical case.

Leader

What specific evidence from the chest X-ray image supports the conclusion of increased lung density and interstitial markings surrounding the heart and lungs, and how does this relate to potential underlying conditions such as pneumonia, pleurisy, or other respiratory infections?

Leader

How does the quality of the image provided in the medical case impact the assessment of the lung parenchyma, and what additional diagnostic tests or clinical information would be helpful to provide a more comprehensive evaluation of the individual's health status?

Critic 1

The image provided appears to be an X-ray of a patient's chest. However, due to the low resolution and quality of the image, it is difficult to accurately assess the condition of the lung parenchyma. There are no visible signs of density variations, lesions, anatomical structures, or contrast patterns that would indicate increased lung density or interstitial markings surrounding the heart and lungs.

Critic 2

The quality of the image provided in the medical case can impact the assessment of the lung parenchyma by limiting the visibility of subtle abnormalities or pathology that may not be immediately apparent. A clear and well-lit X-ray image would allow for a more comprehensive evaluation of the individual's health status.

Leader

*Answer:* **Yes** ✗
The critique against Option A was a minor point that was successfully clarified through the Critic's final response, as it was acknowledged that the image provided did not offer sufficient information to support a conclusion of increased lung density or interstitial markings surrounding the heart and lungs. In contrast, the critique against Option B was a fatal contradiction because it failed to address the lack of visible signs of density variations, lesions, anatomical structures, or contrast patterns that would indicate increased lung density or interstitial markings surrounding the heart and lungs.
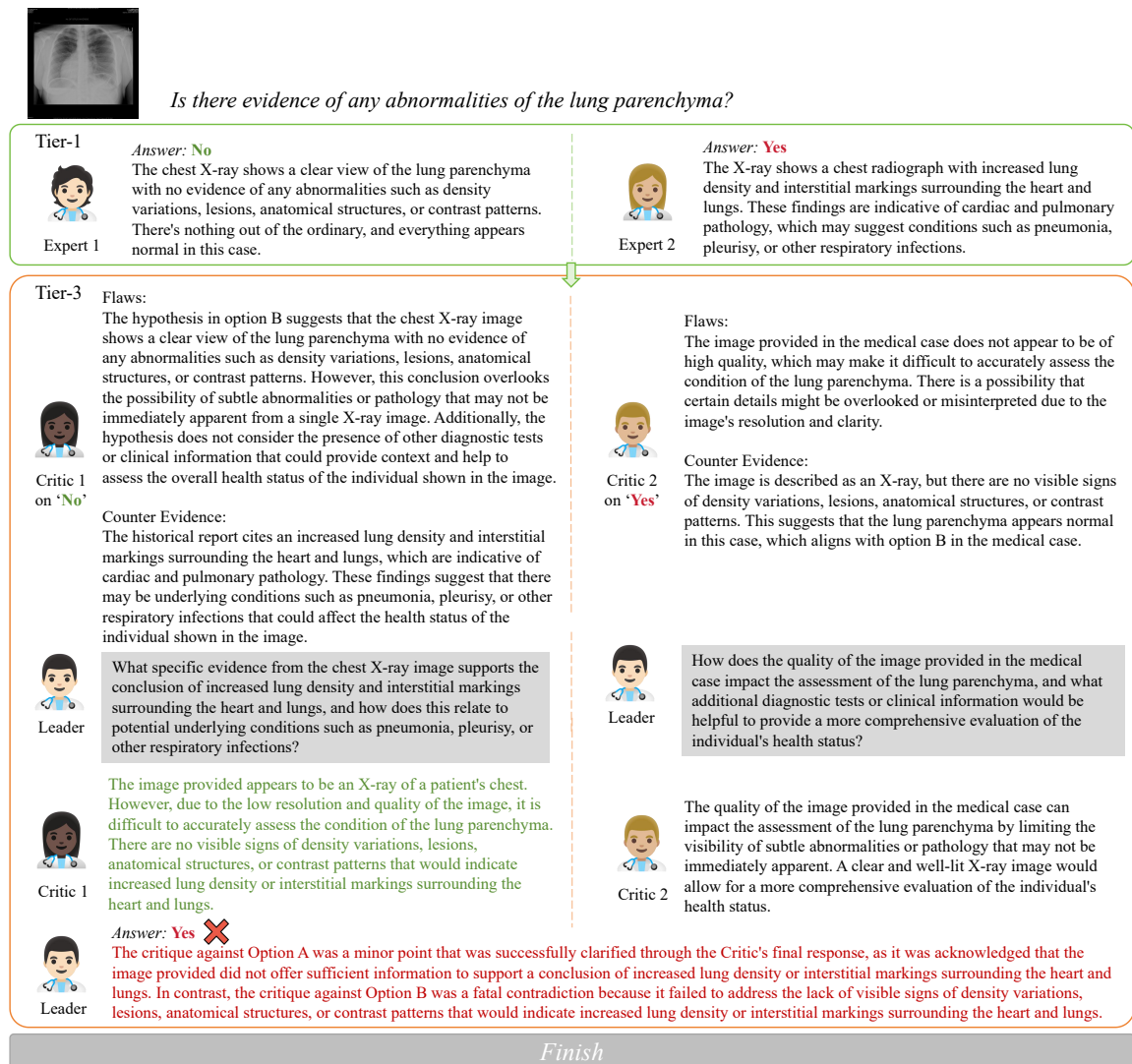
*Finish*

Figure 6. Failure Mode 2: Ambiguous Visual Evidence Under Low Image Quality.