

CXRAgent: Director-Orchestrated Multi-Stage Reasoning for Chest X-Ray Interpretation

Jinhui Lou, Yan Yang, *Member, IEEE*, Zhou Yu, *Member, IEEE*, Zhenqi Fu, Weidong Han, Qingming Huang, *Fellow, IEEE*, and Jun Yu, *Senior Member, IEEE*

Abstract—Chest X-ray (CXR) plays a pivotal role in clinical diagnosis, and a variety of task-specific and foundation models have been developed for automatic CXR interpretation. However, these models often struggle to adapt to new diagnostic tasks and complex reasoning scenarios. Recently, LLM-based agent models have emerged as a promising paradigm for CXR analysis, enhancing model’s capability through tool coordination, multi-step reasoning, and team collaboration, etc. However, existing agents often rely on a single diagnostic pipeline and lack mechanisms for assessing tools’ reliability, limiting their adaptability and credibility. To this end, we propose CXRAgent, a director-orchestrated, multi-stage agent for CXR interpretation, where a central director coordinates the following stages: (1) **Tool Invocation**: The agent strategically orchestrates a set of CXR-analysis tools, with outputs normalized and verified by the Evidence-driven Validator (EDV), which grounds diagnostic outputs with visual evidence to support reliable downstream diagnosis; (2) **Diagnostic Planning**: Guided by task requirements and intermediate findings, the agent formulates a targeted diagnostic plan. It then assembles an expert team accordingly, defining member roles and coordinating their interactions to enable adaptive and collaborative reasoning; (3) **Collaborative Decision-making**: The agent integrates insights from the expert team with accumulated contextual memories, synthesizing them into an evidence-backed diagnostic conclusion. Experiments on various CXR interpretation tasks show that CXRAgent delivers strong performance, providing visual evidence and generalizes well to clinical tasks of different complexity. Code and data are valuable at this link.

Index Terms—Chest X-ray interpretation, medical agent, large multimodal model, multi-stage reasoning.

I. INTRODUCTION

CHEST X-ray (CXR) is among the most widely used imaging modalities in clinical practice due to their affordability, rapid acquisition, and diagnostic utility across a wide range of thoracic conditions. However, accurate interpretation demands substantial clinical expertise and years

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

Jinhui Lou, Yan Yang, Zhou Yu are with the School of Computer Science, Hangzhou Dianzi University, Hangzhou, 310018, China (e-mail: 22081034@hdu.edu.cn; yangyan@hdu.edu.cn; yuz@hdu.edu.cn).

Zhenqi Fu is with the Department of Automation, Tsinghua University, Beijing 100084, China (e-mail: fuzhenqi@mail.tsinghua.edu.cn).

Weidong Han is with the Department of Colorectal Medical Oncology, Zhejiang Cancer Hospital, Hangzhou, 310022, China (e-mail: hanwd@zjcc.org.cn).

Qingming Huang is with the School of Computer and Control Engineering, University of Chinese Academy of Sciences, Beijing, 101408, China (e-mail: qmhuang@ucas.ac.cn).

Jun Yu is with the School of Intelligence Science and Engineering, Harbin Institute of Technology (Shenzhen), 518055, China. (e-mail: yu-jun@hit.edu.cn).

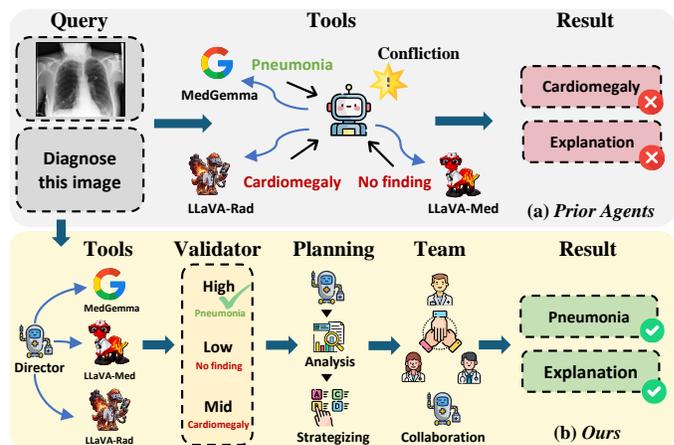


Fig. 1: Comparison between prior agent models and our CXRAgent. (a) Prior agents often struggle to deal with tool conflicts in complex cases and rely on a single diagnosis pipeline. (b) CXRAgent adopts a multi-stage pipeline, guided by a core director, to flexibly coordinate tool invocation, tool-output validation, diagnostic planning, and team-integrated collaborative decision-making. Each tool’s confidence is validated with visual evidence to ensure reliability.

of specialized training. In addition, the heavy workload of radiologists makes the diagnostic process susceptible to delays and errors, including missed or incorrect findings. These challenges have motivated a growing interest in AI-powered systems that aim to assist clinicians and improve the efficiency and precision of CXR interpretation.

Previous efforts have primarily focused on developing task-specific models for automated CXR interpretation. Representative efforts in CXR report generation include RGRG [1], R2GenGPT [2], MAIRA-2 [3], and LLaVA-Rad [4]. RGRG [1] proposed to first detect anatomical structures and then integrate the region-level descriptions into a unified radiology report. R2GenGPT [2] introduced a lightweight fine-tuning strategy to achieve efficient vision-language alignment for radiology report generation. MAIRA-2 [3] advanced grounded radiology report generation by incorporating other reporting contexts as additional inputs. LLaVA-Rad [4] trains a domain-adapted chest X-ray encoder and connects it to pre-trained language models through a lightweight adapter for vision-language alignment. While these task-specific models have achieved notable success in individual tasks, they struggle to generalize across different tasks (for example, from report generation to question answering).

Recent progress in deep learning has further fostered the development of medical foundation models for CXR analysis. Examples include RadFM [5], CheXagent [6], Ark⁺ [7], and MedGemma [8]. RadFM [5] introduced a generalist foundation model capable of modality recognition, disease diagnosis, visual question answering, etc. CheXagent [6] curated a large-scale CheXinstruct dataset and developed a vision-language foundation model for CXR interpretation. Ark⁺ [7] employed a cyclical pretraining strategy to accumulate and reuse knowledge embedded in expert labels for CXR analysis. MedGemma-4B [8] is a large multimodal model designed for medical applications and aims to serve as a tool for agent-based diagnostic systems. However, although these foundation models can perform well on predefined diagnostic tasks, they still struggle to generalize to new diagnostic objectives and complex reasoning scenarios.

Most recently, agent-based medical diagnostic models have gained increasing attention for their ability to unify the strengths of diverse models (e.g., task-specific models and foundation models) and mimic clinical reasoning workflows. These agents aim to advance diagnostic intelligence by enabling multi-step, context-aware reasoning and integrating diverse diagnostic tools. In particular, agent models decompose complex diagnostic tasks into a series of subtasks. By orchestrating tools step by step and solving individual subtasks, agents enhance the transparency of the diagnostic pipeline. Representative agent models in medical image analysis such as MedRAX [9], MDAgents [10] and MMedAgent [11] have demonstrated the potential of agent-based systems by intelligently coordinating multiple analysis tools. However, they typically treat tools as equivalent components, lacking effective mechanisms to assess the reliability and expertise of each tool’s output. As shown at the top of Fig. 1, they struggle to reconcile the conflicting outputs across tools, thereby failing to prioritize more trustworthy information. This risks diluting critical diagnostic signals and undermines the diagnosis credibility. Moreover, prior agents often rely on a rigid single-stage workflow, lacking flexibility. In this paper, as shown at the bottom of Fig. 1, we introduce CXRAgent, a multi-stage, director-orchestrated agent tailored for evidence-backed and adaptive CXR interpretation. Firstly, the agent orchestrates a set of CXR analysis tools, with their outputs reformatted and validated by an Evidence-driven Validator (EDV). The EDV evaluates the confidence of each tool’s output by examining supporting or refuting visual evidence, ensuring reliable inputs for downstream diagnosis. Then, the agent performs flexible planning by assembling the expert team and enabling adaptive collaborative reasoning tailored to task requirements, thereby improving decision-making in complex clinical scenarios. The main contributions of our paper are as follows:

- We propose CXRAgent, a multi-stage, director-orchestrated agent for chest X-ray diagnosis that employs a powerful multi-modal LLM as its core director to flexibly coordinate tool invocation, diagnostic planning, and team-based collaborative decision-making.
- We propose the Evidence-driven Validator (EDV) that grounds tool results with visual supportive and refuting

evidence, provides confidence assessments to evaluate diagnostic reliability and unifies the output formats of various tools to support reliable downstream diagnosis.

- Inspired by clinical Multidisciplinary Team (MDT) practices, we propose a flexible team-based collaborative paradigm that intelligently assembles and coordinates specialized agent teams, ensuring robust adaptability to diverse tasks and diagnostic complexities.
- Experimental results demonstrate that CXRAgent achieves substantial performance gains across various CXR interpretation tasks including visual question answering and report generation, validating its diagnostic accuracy and adaptability.

II. RELATED WORKS

In this section, we review related works on chest X-ray analysis models and medical agent systems.

A. Chest X-ray Analysis Models

Recent studies in chest X-ray analysis have introduced numerous task-specific models aimed at distinct clinical objectives, such as classification [12], detection [13], localization [14], segmentation [15], visual question answering [16], [17], and report generation [2]–[4], [18]. PIXEL [14] introduced a prompt-driven constrained generative framework that synthesizes anatomically aligned healthy–diseased image pairs to enable supervised learning of pathology localization. Building on visual prompting, FAVP [17] proposed an Adaptive Visual Prompt Creator that dynamically generates eight types of region-level visual prompts to improve CXR visual question answering. LLaVA-Rad [4] adopted a three-stage training paradigm with domain-specific encoder pre-training for efficient and clinically deployable CXR report generation.

The advent of multimodal foundation models has further advanced CXR understanding and interpretation. LLaVA-Med [19] constructed a curated biomedical image–text instruction dataset and optimized a vision–language model tailored for biomedical tasks. MedFILIP [20] achieved fine-grained cross-modal alignment through entity extraction and knowledge injection, while Med-Flamingo [21] leveraged high-quality interleaved data and few-shot learning to enhance adaptability. CheXagent [6] curated the large-scale CheXinstruct dataset and developed a vision–language foundation model specialized for CXR interpretation. Lingshu [22] presented a generalist foundation model for unified multimodal medical understanding and reasoning. DeepMedix-R1 [23] introduced a holistic medical foundation model for CXR interpretation, following a sequential training strategy: it is first fine-tuned on curated instruction data to acquire fundamental interpretive skills, then trained on synthetic reasoning samples for cold-start reasoning, and finally refined via online reinforcement learning to enhance reasoning precision and generative quality. MedGemma-4B [8] proposed a multimodal model designed for medical image understanding and reasoning, integrating visual and textual modalities to support diverse clinical tasks.

Although these models have achieved impressive results across benchmarks, they encounter challenges in handling

complex, ambiguous, or multi-condition cases. Their architectures often lack the flexibility required to dynamically adapt to diverse diagnostic tasks. However, these advances lay the groundwork for agent-based CXR analysis, which aspires to achieve adaptive, interpretable, and complex reasoning.

B. Medical Agent Systems

Medical agent systems [24]–[26] have emerged as a promising paradigm to overcome the limitations of monolithic models by orchestrating multiple specialized components. MedAgents [27] introduced a multi-disciplinary collaboration framework, where LLM-based agents engage in role-playing and multi-round discussions to enhance reasoning in medical tasks. MDAgents [10] proposed a multi-agent framework that automatically assigns task-specific solo or group collaboration structures to LLMs, emulating real-world medical decision-making processes. ClinicalAgent [28] applied this agentic paradigm to clinical trial prediction, while DoctorAgent-RL [29] modeled clinical consultations as dynamic decision-making under uncertainty using multi-agent reinforcement learning. Deep-DxSearch [30] proposed end-to-end reinforcement learning to jointly optimize retrieval and reasoning, enabling retrieval-aware diagnostic strategies beyond inference-only agentic systems.

Recently, multi-modal agents for medical image analysis have also been developed. MMedAgent [11] proposed to manage diverse medical tasks by integrating open-source medical models. It consists of (1) an instruction-tuned multi-modal LLM serving as an action planner and results aggregator, and (2) a set of task-specific medical tools. MMA [31] conducted medical diagnosis question answering by deploying specialized LLM-based agents that processed multi-source data and collaborated via structured workflows. MedRAX [9] demonstrated the potential of agent-based systems by orchestrating multiple diagnostic tools for CXR interpretation. MAM [32] proposed a modular multi-agent architecture for multi-modal medical diagnosis, leveraging role decomposition and collaborative decision-making. CT-Agent [33] incorporated anatomy-aware reasoning with hierarchical token compression to efficiently handle 3D CT analysis. PathoAgenticRAG [34] presented a multimodal retrieval-augmented generation framework for pathology vision-language models, mitigating hallucinations via a database of page-level embeddings from pathology textbooks. Agentic paradigms have also been extended to other medical applications. For example, Xie et al. [35] proposed a decentralized multi-agent reinforcement learning strategy to interact dynamically with a foundation model, mitigating its undue influence and progressively refining priors for low-count PET reconstruction.

However, current medical imaging diagnosis Agent systems face persistent challenges: (1) insufficient mechanisms to reconcile conflicting outputs across tools, (2) reliance on a single, rigid workflow that ignore case complexity, and (3) weak or absent connections between diagnostic conclusions and supporting visual evidence.

III. METHOD

As illustrated in Fig. 2, CXRAgent is designed as a multi-stage, director-orchestrated agent for chest X-ray interpretation. Upon receiving a user query, it leverages a multi-modal LLM as the central “director” to coordinate the following key stages: (1) **Tool Invocation** — The agent orchestrates multiple CXR-analysis tools through iterative reasoning cycles. Multimodal inputs are analyzed step by step, and the outputs of various tools are normalized and validated by the Evidence-driven Validator, which aligns tool responses with corresponding visual evidence to provide reliable insights for downstream reasoning. (2) **Diagnostic Planning** — Guided by task requirements and intermediate results, the agent formulates subsequent diagnostic plans that adapt to evolving clinical needs and varying case complexities. (3) **Collaborative Decision-Making** — Guided by the diagnostic plan and task objectives, the agent assembles a team of specialized experts to execute the required analyses, integrating insights from team decisions and historical reasoning trajectories to deliver a comprehensive, evidence-grounded diagnostic conclusion.

A. Tool Invocation

We first employ the ReAct [36] to iteratively plan and execute tool invocations, followed by the Evidence-driven Validator (EDV) to refine and evaluate the tool outputs with visual evidence. Concretely, ReAct contains the following steps:

Thought Step – State Analysis. At each iteration, the agent decides whether additional tools are needed. The decision is jointly conditioned on the input images I , the textual query q , and the accumulated reasoning memory log, denoted as Log_{i-1} :

$$\text{Thought}_i = \text{Plan}(\text{Log}_{i-1}, I, q). \quad (1)$$

If the agent determines that invoking additional tools is necessary to gather sufficient information under the current iteration i for diagnosis, it proceeds to the subsequent stages.

Action Step – Tool Invocation. The model selects the most appropriate tool(s) based on the current reasoning context, formulated as:

$$\begin{aligned} \text{Action}_i &= \{a_{i1}, a_{i2}, \dots, a_{in}\} \\ &= \bigcup_{j=1}^n \text{Tool}_{ij}(I, q, \text{Thought}_i), \quad \text{Tool}_{ij} \in \mathcal{T} \end{aligned} \quad (2)$$

where \mathcal{T} is the tool set, Tool_{ij} is the j -th tool of i -th iteration. $a_{i1}, a_{i2}, \dots, a_{in}$ is the selected tools (i.e., arrangement of tool invocation), n is the total number of tools invoked. All invocations are logged for further reasoning.

Observation Step – Result Aggregation. The selected tools $\{a_{i1}, a_{i2}, \dots, a_{in}\}$ are applied to the relevant image-query pairs, producing results $\{o_{i1}, o_{i2}, \dots, o_{in}\}$, which are then aggregated as observations to update the reasoning log for the next iteration. These observations will provide new information for the next iteration of the reasoning process:

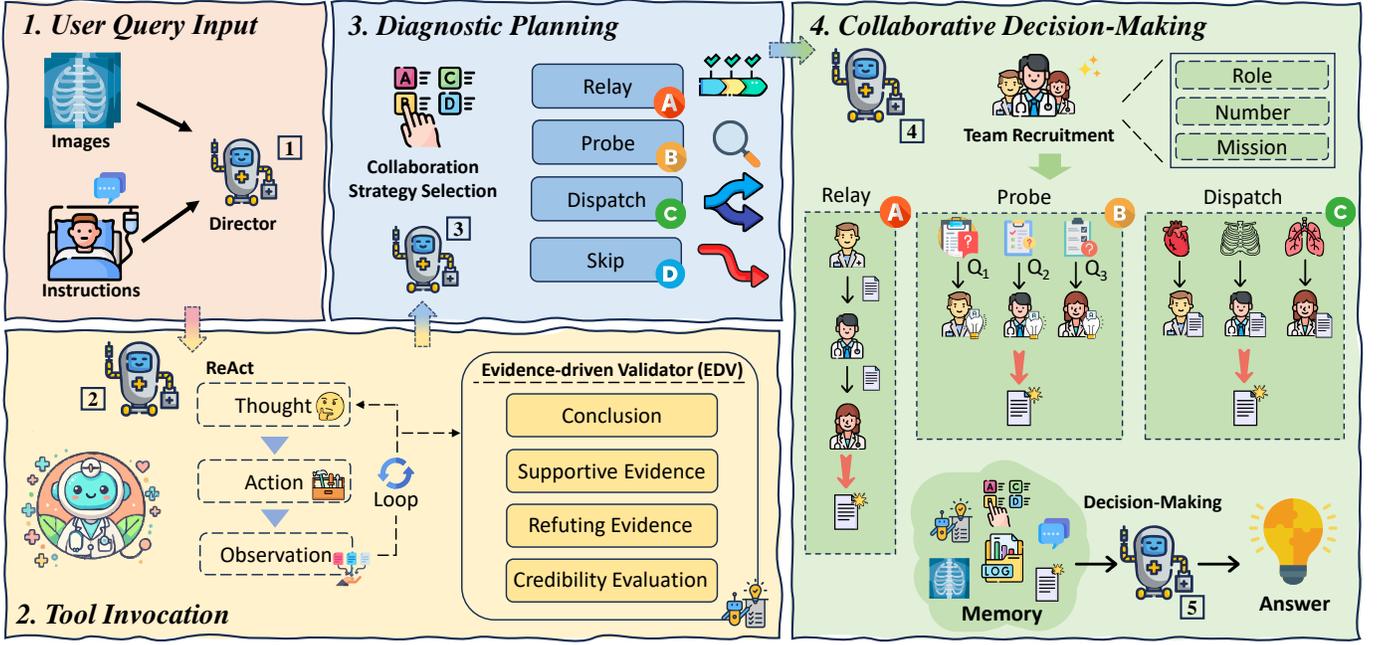


Fig. 2: **Overview of the CXRAgent.** The central director coordinates the following key stages: (1) ingesting user queries, (2) orchestrating CXR-analysis tools with outputs normalized and validated by the Evidence-driven Validator (EDV) to ensure visually grounded reliability, (3) making diagnostic plans to adaptively assemble specialized expert teams for collaborative reasoning, and (4) synthesizing evidence-backed conclusions by integrating expert collaborative insights and the contextual memories.

$$\text{Obs}_i = \{o_{i1}, o_{i2}, \dots, o_{in}\}, \quad (4)$$

$$\text{Log}_i \leftarrow \text{Log}_{i-1} \cup \text{Obs}_i. \quad (5)$$

Thought Step – Reflection and State Analysis. After gathering the results, the model reflects on the new observations and updates its reasoning state. This process determines whether further tool invocations are necessary:

$$\text{Thought}_{i+1} = \text{Plan}(\text{Log}_i, I, q). \quad (6)$$

The model repeatedly performs the thought, action, and observation steps until sufficient information is accumulated.

Evidence-driven Validator (EDV). To mitigate diagnostic errors arising from unreliable tool outputs, we introduce Evidence-driven Validator, a tool-agnostic verification module. Unlike methods that merely summarize tool results, EDV actively grounds results with visual evidence from the input image. For each diagnostic statement produced by a tool, EDV identifies supporting and contradicting evidence from the image and estimates the overall confidence of the statement. It also standardizes the statements of tools. After validating and normalizing outputs from individual tools, EDV enhances overall reliability of the diagnostic process. Concretely, EDV generates a set of structured components that explain, validate, and assess the credibility of each diagnostic result:

- **Conclusion:** A concise reformatted restatement of the diagnostic statement produced by the tool.
- **Supportive Evidence:** Key image findings that correspond to and reinforce the diagnostic statement.

- **Refuting Evidence:** Observations that contradict the statement, including absent key evidence or contradictory cues.
- **Confidence Assessment:** A qualitative judgment of the statement’s reliability, based on both the supportive and refuting evidence.

The structured and reformatted outputs produced by EDV establish a unified framework for validating and explaining predictions across different diagnostic tools. By explicitly grounding each conclusion in the visual evidence from the input data, EDV not only improves transparency but also enables consistent, interpretable cross-tool evaluation for the downstream interpretation.

B. Diagnostic Planning

In real-world chest X-ray diagnostics, user queries vary widely in complexity and requirements, necessitating adaptable reasoning strategies and collaboration among multidisciplinary teams. To this end, we design Diagnostic Planning, a stage guided by the “director” that dynamically determines collaboration modes (i.e., how to adaptively assemble specialized expert teams and formulate reasoning strategies) based on task-specific characteristics. We define four key collaboration modes:

- 1) **[Skip]** — Directly generates diagnostic results without team collaboration for those straightforward cases with unambiguous visual evidence, maximizing the diagnosis efficiency.
- 2) **[Relay]** — Employs a sequential refinement process where each specialist incrementally improves the diag-

nosis. Each expert builds on the outputs of the previous one, progressively refining conclusions. This mode is ideal for routine cases that require step-by-step refinement.

- 3) **[Dispatch]** — Decomposes complex tasks into specialized subtasks, assigning each to domain-specific experts for parallel analysis, enabling focused feature analysis. For example, assigning a heart expert for cardiac diagnosis and a bone expert for skeletal diagnosis.
- 4) **[Probe]** — Breaks down the case into targeted probing questions, prompting each team member to respond individually, facilitates the gathering of diverse and focused insights, supporting comprehensive understanding and diagnostic reasoning.

For cases with clear diagnostic features, the system employs the Skip strategy to deliver immediate results without activating the team, maximizing efficiency for straightforward findings. When presented with well-defined diagnostic queries supported by unambiguous visual evidence, the system applies the Relay strategy, sequentially refining the diagnosis through stepwise expert analysis to ensure efficient and incremental improvement. For tasks requiring multi-faceted feature analysis, the Dispatch strategy engages specialized experts working in parallel, each focusing on distinct subtasks for targeted and efficient processing. Finally, for complex and ambiguous cases demanding high-level reasoning, the Probe strategy actively generates targeted probing diagnostic questions to be addressed, enabling them to uncover subtle clues and collaboratively build a comprehensive and well-supported diagnostic conclusion.

C. Collaborative Decision-Making

Based on the collaboration strategy determined by the Diagnostic Plan and the characteristics of the clinical case, the agent recruits a diagnostic team, intelligently configuring the number of team members, their roles, and missions. This configuration enables context-aware collaboration across diverse clinical scenarios. By aligning team composition with case complexity and diagnostic requirements, the system enables flexible workflows and delivers more reliable diagnostic result. Team recruitment can be formulated as:

$$\text{Team} = \{(\text{agent}_i, \text{role}_i, \text{mission}_i)\}_{i=1}^k = \mathcal{R}(I, q, s), \quad (7)$$

where \mathcal{R} denotes the team recruitment strategy, $s \in \{\text{Skip}, \text{Relay}, \text{Dispatch}, \text{Probe}\}$ is the selected collaboration strategy, and each tuple $\text{Agent}_i = (\text{agent}_i, \text{role}_i, \text{mission}_i)$ specifies the identity, functional role, and assigned mission of the i -th team member. Here, k indicates the team size. This formulation allows dynamic configuration of diagnostic teams tailored to the complexity and specific demands of each clinical case. The variable mission_i represents the sub-mission allocated to each agent, which varies by strategy: refinement tasks in Relay, dedicated subtasks in Dispatch, and probing questions in Probe mode.

Algorithm 1 demonstrates the complete workflow of team formation and collaborative diagnosis. After collecting EDV-validated tool outputs, team collaboration results, and the

Algorithm 1 Collaborative Decision-Making

Input: Query q , CXR images I , Strategy s , Memory M

Output: Diagnosis Result

```

1: Team  $\leftarrow \mathcal{R}(q, I, s)$ 
2: if  $s = \text{Skip}$  then
3:   TeamOutput  $\leftarrow \text{None}$ 
4: else if  $s = \text{Relay}$  then
5:    $r_0 \leftarrow \emptyset$ 
6:   for  $i = 1$  to  $k$  do
7:      $r_i \leftarrow \text{Agent}_i(I, r_{i-1}, \text{mission}_i := \text{task}_i)$ 
8:   end for
9:   TeamOutput  $\leftarrow r_k$ 
10: else if  $s = \text{Dispatch}$  then
11:   for  $i = 1$  to  $k$  do
12:      $r_i \leftarrow \text{Agent}_i(I, \text{mission}_i := \text{sub-task}_i)$ 
13:   end for
14:   TeamOutput  $\leftarrow \text{Concat}(\{r_1, r_2, \dots, r_k\})$ 
15: else if  $s = \text{Probe}$  then
16:   for  $i = 1$  to  $k$  do
17:      $r_i \leftarrow \text{Agent}_i(I, \text{mission}_i := \text{question}_i)$ 
18:   end for
19:   TeamOutput  $\leftarrow \text{Concat}(\{r_1, r_2, \dots, r_k\})$ 
20: end if
21: Diagnosis Result  $\leftarrow \text{Diagnose}(q, I, s, M, \text{TeamOutput})$ 
22: Return Diagnosis Result

```

contextual memories, CXRAgent makes the final diagnostic decision through evidence-based reasoning.

IV. EXPERIMENTS

A. Experiment Setup

Datasets and Tasks. We evaluate CXRAgent across multiple CXR interpretation tasks on the following three datasets: (1) CheXbench [6] serves as a comprehensive benchmark crafted to assess three aspects of CXR analysis including image perception, image-text reasoning, and text generation. We focus on image perception (i.e., view classification and disease identification), image-text reasoning (fine-grained reasoning, visual-question answering (VQA)). (2) Medical-CXR-VQA [37] is a benchmark for CXR visual question answering. Following its official tasks, we randomly sample 594 Presence, 256 Abnormality, and 150 View questions to assess the model. (3) MIMIC-CXR [38] provides 377,110 CXRs paired with 227,835 reports. We randomly select 400 cases from the testing split for the evaluation of CXR report generation. The data will be available for reproducibility.

Implementations. We use GPT-4o or Qwen-VL-Max as the core director, deployed on two NVIDIA RTX A6000 GPUs. To enhance diagnostic capabilities of our CXRAgent, we have integrated a set of advanced tools:

- **MedGemma-4B [8]:** An open-source medical multimodal model finetuned on Gemma3, featuring a SigLIP vision encoder for clinical imaging analysis and visual-language reasoning.
- **LLaVA-Rad [4]:** A multimodal model, optimized for automated radiology report generation from CXR images.

TABLE I: Performance comparison on the CheXbench across tasks (accuracy in %). Tasks: SDI (Single Disease Identification), MDI (Multiple Disease Identification), FGR (Fine-Grained Reasoning), BDC (Binary Disease Classification), VQA (Visual Question Answering), VC (View Classification). Overall is the average over all tasks. Best results are in bold.

Categories	Models							
	LLaVA-Med	CheXagent	GPT-4o	Qwen-VL-Max	MedRAX (GPT)	MedRAX (Qwen)	Ours (GPT)	Ours (Qwen)
SDI <small>MIMIC-CXR</small>	26.7	30.3	29.7	31.7	28.8	32.8	39.3	42.9
SDI <small>CheXpert</small>	26.0	29.6	49.7	51.4	34.9	45.6	52.1	55.6
MDI <small>MIMIC-CXR</small>	28.7	55.3	53.6	53.3	51.6	44.0	64.6	61.3
MDI <small>CheXpert</small>	15.0	52.1	62.1	68.5	49.3	51.4	73.2	68.9
VQA <small>Rad-Restruct</small>	34.9	57.3	60.8	58.2	68.7	57.4	62.6	59.1
VQA <small>SLAKE</small>	55.5	78.1	85.1	78.8	82.9	79.6	86.2	87.0
FGR <small>OpenI</small>	45.6	59.0	41.8	51.3	52.6	48.1	53.4	59.4
BDC <small>CheXpert</small>	47.6	76.0	64.8	61.8	55.8	58.8	64.8	60.5
BDC <small>SHM</small>	49.0	64.0	57.0	51.0	67.0	65.0	77.0	72.0
BDC <small>RSNA</small>	44.0	81.0	76.0	71.0	78.0	76.0	83.0	83.0
VC <small>MIMIC-CXR</small>	23.8	97.5	79.0	86.6	58.7	65.3	74.0	87.7
Overall	36.0	61.8	59.9	60.3	57.1	56.7	66.3	67.0

TABLE II: Performance comparison of multimodal models on Medical-CXR-VQA across clinical question types (accuracy in %). All questions are binary classification tasks. Overall is the average over all tasks. Best results are in bold.

Categories	Models							
	LLaVA-Med	CheXagent	GPT-4o	Qwen-VL-Max	MedRAX (GPT)	MedRAX (Qwen)	Ours (GPT)	Ours (Qwen)
Presence	50.7	63.3	67.2	63.8	67.5	66.4	70.2	68.0
Abnormality	37.1	75.0	70.4	67.9	71.0	66.6	73.4	75.7
View	44.4	48.0	70.1	80.6	66.0	78.7	71.4	83.3
Overall	44.1	62.1	69.2	70.8	68.2	70.5	71.7	75.6

TABLE III: Performance comparison on CXR report generation. RaTEScore measures entity recognition, while LLM-Score (1–5) assesses overall quality. Best results are in bold.

Model	RaTEScore	LLM Score
LLaVA-Med	0.370	1.26
CheXagent	0.510	2.59
LLaVA-Rad	0.503	2.99
MedGemma	0.465	2.45
GPT-4o	0.429	2.06
Qwen-VL-Max	0.415	1.80
MedRAX (with GPT)	0.462	2.25
MedRAX (with Qwen)	0.466	2.24
Ours (with GPT)	0.508	2.77
Ours (with Qwen)	0.513	2.69

- **CheXagent [6]**: A foundation model for interpretation of chest X-rays (e.g., disease classification and localization).
- **LLaVA-Med [19]**: A multimodal assistant enabling interactive image-based medical diagnosis.
- **MAIRA-2 [3]**: A radiology-specific multimodal model for grounded report generation.
- **MedVLM-R1 [39]**: A vision-language model for radiological tasks that generates natural language reasoning alongside the final answer.

Baselines. We compare our framework with three representative categories of baselines: 1) general-purpose multimodal

LLMs (Qwen-VL-Max, GPT-4o) with cross-domain medical adaptation capabilities; 2) radiology-specialized models (LLaVA-Med, CheXagent, LLaVA-Rad, MedGemma) fine-tuned on medical tasks; and 3) MedRAX, a CXR agent that integrates various diagnostic tools for CXR analysis without requiring model fine-tuning.

Evaluation Metrics. For report generation, we use two metrics: 1) RaTEScore [40] with entity-level evaluation of reports by assessing critical elements through its specialized named entity recognition (NER) framework; and 2) LLMscore, implemented with GPT-4o, which rates reports on a 1–5 scale following the ITU five-point standard (1 = Bad, 5 = Excellent). For other tasks, we employ accuracy for assessment. Results are reported as the average over all samples.

B. Performance Comparison

CheXbench. As demonstrated in Table I, our Qwen-based CXRagent achieves state-of-the-art performance on the comprehensive CheXbench benchmark, attaining an overall accuracy of 67.0%. The framework demonstrates remarkable improvements across multiple challenging tasks, particularly excelling in complex diagnostic scenarios. For multi-disease identification, CXRagent achieves 73.2% on CheXpert and 64.6% on MIMIC-CXR, representing significant advancements over previous methods. The model’s robustness is especially evident in fine-grained reasoning tasks, where it reaches 59.4% accuracy on OpenI, outperforming all comparable approaches. This consistent performance across diverse evaluation categories underscores the effectiveness of our multi-stage

adaptive design, particularly the EDV module’s capability to resolve tool inconsistencies through evidence-based validation.

Medical-CXR-VQA. As shown in Table II, CXRAgent delivers exceptional performance on the Medical-CXR-VQA benchmark, achieving a leading overall accuracy of 75.6%. The framework exhibits particularly strong capabilities in abnormality detection (75.7%) and view classification (83.3%), demonstrating its proficiency in handling clinically nuanced questioning. Notably, our approach maintains superior performance across all question types, with presence detection reaching 70.2% accuracy. The adaptive team recruitment mechanism proves crucial in these visual question answering scenarios, enabling dynamic collaboration patterns that effectively address the varied complexity of clinical inquiries. The consistent superiority across both benchmarks demonstrates CXRAgent’s strong adaptability.

MIMIC-CXR. As demonstrated in Table III, our proposed CXRAgent framework achieves superior performance in chest X-ray report generation, attaining the highest RaTEScore of 0.513 among all evaluated models. This metric, which specifically measures the accuracy of clinical entity recognition, highlights our method’s exceptional capability in identifying and describing medically relevant findings with precision. The RaTEScore performance represents a significant advancement over general-purpose multimodal models such as GPT-4o (0.429) and Qwen-VL-Max (0.415), demonstrating the effectiveness of our approach to report generation. In terms of overall report quality assessment using the LLMscore metric, our GPT-based implementation achieves a score of 2.77, ranking second only to LLaVA-Rad (2.99), which is specifically optimized for radiology report generation. This performance is particularly noteworthy given that LLaVA-Rad employs dedicated training strategies tailored specifically for this task. Compared to baseline models, our framework demonstrates substantial improvements, with relative enhancements of 34.5% over GPT-4o (2.06) and 23.1% over MedRAX with GPT backbone (2.25). These improvements underscore the effectiveness of our multi-stage reasoning framework in generating clinically coherent and comprehensive radiology reports. The consistent performance across both evaluation dimensions—entity-level recognition (RaTEScore) and overall quality (LLMScore)—validates the robustness of our approach. The integration of EDV ensures that generated reports maintain high factual accuracy, while the adaptive team collaboration mechanism enables comprehensive coverage of clinical findings. This balanced performance across different assessment criteria demonstrates the framework’s capability to handle the complex requirements of clinical report generation, where both precise entity extraction and overall narrative quality are essential for practical utility.

C. Ablation Study

CheXbench. Shown in Tables IV, V, the ablation results reveal three key findings: First, although tool integration generally improves performance, incorporating unevaluated tool outputs can occasionally degrade the performance. For example, in the view classification task on CheXbench (Qwen),

accuracy drops from 86.6% without tools to 83.6% with tools. Second, the EDV module effectively mitigates such errors by resolving conflicts through visual evidence validation, playing a critical role in ensuring reliability when tool outputs are inconsistent. Third, the full model that combines tools, EDV, and team collaboration delivers optimal performance across tasks, highlighting the synergistic contribution of all components. Under the GPT-4o as director, it yields a relative accuracy improvement of 10.7% over the baseline and shows particular advantages in complex reasoning tasks, increasing multi-disease identification (CheXpert) to 73.2% and fine-grained interpretation (OpenI) to 53.4%. Overall, these results demonstrate that CXRAgent delivers robust and consistent gains across different large language model backbones, confirming that the combination of specialized tools, evidence-driven validation, and collaborative reasoning leads to clinically reliable and generalizable diagnostic performance.

Medical-CXR-VQA. The ablation results presented in Table VI demonstrate consistent performance improvements through the progressive integration of our framework’s components on the Medical-CXR-VQA benchmark. Three key patterns emerge from the analysis. First, the incorporation of specialized CXR-analysis tools establishes a solid foundation for accurate diagnosis, with the Qwen-directed system showing particular strength in abnormality detection (increasing from 67.9% to 73.4%). Second, the EDV module contributes significantly to diagnostic reliability, particularly for presence detection tasks where the GPT-4o-based system improves from 68.0% to 70.0%. This enhancement confirms EDV’s role in validating tool outputs against visual evidence, reducing diagnostic uncertainty. Third, the complete framework achieves optimal performance across all question types, with the Qwen-based configuration reaching 75.6% overall accuracy—representing a 4.8% absolute improvement over the baseline. The framework demonstrates robust performance across different director models, with both Qwen and GPT-4o configurations showing the highest scores in their respective categories. These results validate that the synergistic combination of specialized tools, evidence-based validation, and collaborative reasoning enables more reliable and comprehensive CXR interpretation.

MIMIC-CXR. Table VII presents the ablation results for report generation. While integrating tools alone enhances diagnostic accuracy, the full framework integrating both EDV validation and team collaboration yields the best performance. It reveals three critical findings about our framework’s report generation capability. First, integrating specialized CXR-analysis tools substantially enhances diagnostic performance, confirming their value in capturing domain-specific patterns. Second, the EDV module not only improves diagnostic accuracy but also impacts raw entity extraction, enhancing clinical report quality by enforcing evidence-backed coherence. Third, the complete system achieves optimal performance through collaborative reasoning, where virtual specialist teams demonstrate unique value in synthesizing comprehensive findings beyond the capability of individual tools. Overall, these results validate our core design: combining specialized tools with evidence-based validation and expert team coordination

TABLE IV: Ablation study on CheXbench (accuracy in %, with Qwen as director). Best results are in bold. For abbreviation definitions, please refer to Table 1.

Evaluation Task	Component Configuration			
	None	Tools	Tools+EDV	Full Model
SDI <small>MIMIC-CXR</small>	31.7	35.3	40.5	42.9
SDI <small>CheXpert</small>	51.4	52.1	54.4	55.6
MDI <small>MIMIC-CXR</small>	53.3	55.3	58.6	61.3
MDI <small>CheXpert</small>	68.5	66.4	67.5	68.9
VQA <small>Rad-Restruct</small>	58.2	57.3	59.1	59.1
VQA <small>SLAKE</small>	78.8	84.5	86.1	87.0
FGR <small>OpenI</small>	51.3	52.3	53.4	59.4
BDC <small>CheXpert</small>	61.8	62.2	61.8	60.5
BDC <small>SIIM</small>	51.0	68.0	71.0	72.0
BDC <small>RSNA</small>	71.0	76.0	80.0	83.0
VC <small>MIMIC-CXR</small>	86.6	83.6	86.6	87.7
Overall	60.3	63.0	65.4	67.0

TABLE V: Ablation study on CheXbench (accuracy in %, with GPT-4o as director). Best results are in bold.

Evaluation Task	Component Configuration			
	None	Tools	Tools+EDV	Full Model
SDI <small>MIMIC-CXR</small>	29.7	33.3	36.9	39.3
SDI <small>CheXpert</small>	49.7	50.3	51.5	52.1
MDI <small>MIMIC-CXR</small>	53.6	56.7	61.7	64.6
MDI <small>CheXpert</small>	62.1	59.6	65.4	73.2
VQA <small>Rad-Restruct</small>	60.8	55.7	66.1	62.6
VQA <small>SLAKE</small>	85.1	82.9	84.6	86.2
FGR <small>OpenI</small>	41.8	45.3	47.9	53.4
BDC <small>CheXpert</small>	64.8	64.0	65.2	64.8
BDC <small>SIIM</small>	57.0	66.0	74.0	77.0
BDC <small>RSNA</small>	76.0	79.0	81.0	83.0
VC <small>MIMIC-CXR</small>	79.0	72.7	73.3	74.0
Overall	59.9	60.5	64.3	66.3

produces the most clinically reliable reports, highlighting our framework’s strength in generating clinically grounded diagnostic reports rather than merely listing observations.

D. Case Study

Case1. As illustrated in Fig. 3, CXRAgent performs multi-stage reasoning for chest X-ray interpretation, demonstrating strong capability in handling complex diagnostic cases through evidence-based validation and adaptive collaboration. While baseline models such as MedRAX and MedGemma produce conflicting or partially inaccurate findings, CXRAgent ensures diagnostic consistency through systematic reasoning.

The process begins with the invocation of two specialized tools: MedGemma and LLaVA-Rad. MedGemma provides a detailed yet over-sensitive analysis, reporting abnormalities such as “right upper lobe hyperinflation” and “hazy opacity in the right lower lobe,” whereas LLaVA-Rad offers a more conservative assessment, indicating no focal consolidation, pleural effusion, or pneumothorax.

To verify these outputs, the EDV module conducts visual grounding. For MedGemma’s hyperinflation claim, EDV detects no key radiographic signs such as increased lucency

TABLE VI: Ablation study on Medical-CXR-VQA (accuracy in %). Comparison under different directors (Qwen vs. GPT-4o). Best results are in bold.

Evaluation Task	Component Configuration			
	None	Tools	Tools+EDV	Full Model
(a) Using Qwen as Director				
Presence	63.8	66.3	66.6	68.0
Abnormality	67.9	73.4	75.3	75.7
View	80.6	80.0	82.0	83.3
Overall (Qwen)	70.8	73.2	74.6	75.6
(b) Using GPT-4o as Director				
Presence	67.2	68.0	70.0	70.2
Abnormality	70.4	71.8	72.7	73.4
View	70.1	69.3	70.7	71.4
Overall (GPT)	69.2	69.7	71.1	71.7

TABLE VII: Ablation study on report generation (best results in bold).

Evaluation Metric	Component Configuration			
	None	Tools	Tools+EDV	Full Model
RaTEScore (Qwen)	0.415	0.477	0.492	0.513
LLMScore (Qwen)	1.80	2.32	2.49	2.69
RaTEScore (GPT)	0.429	0.489	0.493	0.508
LLMScore (GPT)	2.06	2.64	2.67	2.77

or diaphragmatic flattening and thus assigns a “not well-supported” rating. In contrast, LLaVA-Rad’s negative findings are validated with high confidence, supported by clear lung fields and sharp costophrenic angles. Given the conflicting evidence, CXRAgent initiates diagnostic planning and adopts the *Dispatch* strategy, assembling a virtual multidisciplinary team of a cardiologist and pulmonologist for evaluation. The final synthesis stage integrates validated evidence, expert analyses, and prior reasoning trajectories, effectively resolving initial tool conflicts, discarding unsupported findings, and confirming true abnormalities, including the novel vertebral wedging observation. The resulting report achieves superior clinical accuracy compared with all baselines, highlighting CXRAgent’s ability to convert diverse tool outputs into coherent, evidence-grounded, and reliable diagnostic conclusions.

Case2. As shown in Fig. 4, CXRAgent demonstrates superior capabilities through its multi-stage reasoning framework when analyzing this challenging chest X-ray case. While Qwen-VL-Max, MedRAX, and MedGemma all provided incorrect diagnoses, CXRAgent correctly identified pneumonia through its evidence-based multi-stage reasoning.

In the tool invocation, CXRAgent intelligently coordinated multiple specialized tools while critically evaluating their outputs through the Evidence-driven Validator, performing comprehensive bidirectional validation. For example, it deemed CheXagent’s diagnosis “pneumonia” was plausible based on the right upper lobe opacity, while also recognizing that the statement pertains to a localized assessment, as suggested by the remark “the overall lung fields appear relatively clear”. Such balanced evaluation differs from other models that pro-

<p>Report Generation</p> <p>Which finding is in this chest X-ray?</p> 	<p>Reference Report</p> <p>The lungs are clear of focal consolidation, effusion, or pneumothorax. The heart is enlarged. Right upper extremity vascular stent is present. The aorta is calcified and tortuous. Multiple thoracic compression deformities are seen.</p>	<p>MedRAX</p> <p>The chest X-ray image shows no significant abnormalities except for mild cardiomegaly. The lungs appear clear, and there is no evidence of pneumonia, pleural effusion, or pneumothorax. However, the pathology classification suggests a moderate to high probability of several conditions such as Atelectasis, Infiltration, Fibrosis, and Cardiomegaly.</p>	<p>MedGemma</p> <p>The heart appears enlarged, occupying a significant portion of the chest width, indicating cardiomegaly. The lung markings are more prominent than usual, suggesting possible interstitial lung disease or pulmonary edema. A hazy opacity is seen in the right lower lobe, which could be pneumonia, atelectasis, or another condition. The right upper lobe appears hyperinflated, possibly due to emphysema or another lung disease. The aortic arch shows calcification. There is no clear evidence of pleural effusion or pneumothorax.</p>
<p>Our CXRAgent</p> <p>Tool Invocation:</p> <p>Tool-1 (MedGemma): The right upper lobe appears hyperinflated, which could be due to emphysema or other lung disease. EDV-1: [Conclusion] Possible Right Upper Lobe Hyperinflation, [Supportive Evidence] The right upper lobe appears hyperinflated, which could be due to emphysema or other lung disease, [Refuting Evidence] The right upper lobe does not show signs of hyperinflation such as increased lucency or flattening of the diaphragm. [Confidence Estimation] The conclusion is not well-supported by the image and is likely inaccurate.</p> <p>Tool-2 (LLaVA-Rad): No focal consolidation, pleural effusion, or evidence of pneumothorax is seen. EDV-2: [Conclusion] No focal consolidation, pleural effusion, or evidence of pneumothorax is seen, [Supportive Evidence] The lung fields are relatively clear without areas of opacity that would suggest consolidation. There are no blunted costophrenic angles or meniscus signs indicative of pleural effusion, and no hyperlucent areas suggesting pneumothorax, [Refuting Evidence] None noted, [Confidence Estimation] The conclusion is highly credible based on the absence of these specific findings.</p> <p>Diagnostic Planning:</p> <p>Strategy (Dispatch): The Dispatch strategy was selected because it requires systematic decomposition to obtain targeted evidence for resolving diagnostic uncertainties in imaging findings.</p> <p>Team Collaboration:</p> <p>Agent-1: [Role] Cardiologist, [Task] Assess the cardiac silhouette for signs of enlargement. Answer-1: The chest X-ray shows a mildly enlarged cardiac silhouette, with the heart appearing to occupy a larger proportion of the thoracic cavity than normal.</p> <p>Agent-2: [Role] Pulmonologist, [Task] Evaluate lung markings for signs of congestion, interstitial lung disease, or other pathologies, assess for focal consolidations or opacities. Answer-2: No focal consolidations, pleural effusions, or signs of pneumothorax are evident. Additionally, there is mild anterior wedging of a lower thoracic vertebral body. The lung fields are clear without significant opacities or infiltrates, and the costophrenic angles are sharp.</p> <p>Response from Memory:</p> <p>The chest X-ray image shows a mildly enlarged cardiac silhouette, with the heart appearing to occupy a larger proportion of the thoracic cavity than normal. The aorta is calcified and tortuous, as evidenced by visible calcifications along its course and an irregular, winding path. No focal consolidation, pleural effusion, or evidence of pneumothorax is seen. The lung fields are relatively clear without areas of increased opacity, blunted costophrenic angles, or hyperlucent areas. There is mild anterior wedging of a lower thoracic vertebral body.</p>			

Fig. 3: Case study of CXRAgent versus state-of-the-art methods on report generation. **Top:** user query, reference report, and baseline predictions. **Bottom:** Prediction of CXRAgent. EDV sections are in purple with the conclusion credibility highlighted. Evidence are in blue, the role and task of team members are in orange. Correct findings are in green and incorrect ones in red.

<p>Disease Identification</p> <p>Which finding is in this chest X-ray?</p> <p>A. Pneumonia B. Fracture C. Pneumothorax D. Edema</p> <p>Answer: Pneumonia</p> 	<p>Qwen-VL-Max</p> <p>The chest X-ray shows bilateral symmetrical increased opacity in the lung fields, which is suggestive of pulmonary edema. There are no clear signs of localized consolidation that would suggest pneumonia, no visible fractures, and no evidence of a pneumothorax.</p> <p>Answer: Edema</p>	<p>MedRAX</p> <p>Based on the classification results showing high probabilities for pneumothorax, consolidation, and lung opacity—along with the visual assessment—pneumothorax appears to be the most likely diagnosis among the given options.</p> <p>Answer: Pneumothorax</p>	<p>MedGemma</p> <p>The left hemithorax shows increased density and blunting of the costophrenic angle, suggestive of a possible pleural effusion—commonly seen in pulmonary edema. The most likely diagnosis is pulmonary edema with associated pleural effusion.</p> <p>Answer: Edema</p>
<p>Our CXRAgent</p> <p>Tool Invocation:</p> <p>Tool-1 (CheXagent): Pneumonia. EDV-1: [Conclusion] Pneumonia, [Supportive Evidence] There is an area of opacity in the right upper lung field, which could represent consolidation consistent with pneumonia, [Refuting Evidence] The overall lung fields appear relatively clear without significant areas of consolidation or infiltrate that are typically seen in more severe cases of pneumonia, [Confidence Estimation] The conclusion is plausible given the opacity in the right upper lung field, but further clinical correlation and possibly additional imaging would be beneficial to confirm the diagnosis.</p> <p>Tool-2 (LLaVA-Rad): No pneumothorax is seen. EDV-2: [Conclusion] No pneumothorax is seen, [Supportive Evidence] The lung fields are fully expanded without any air-crescent sign, suggesting no pneumothorax, [Refuting Evidence] Evaluation of the lung periphery reveals subtle findings that could represent a very small pneumothorax, [Confidence Estimation] The absence of pneumothorax is also correctly noted. The credibility of the conclusion is high.</p> <p>Diagnostic Planning:</p> <p>Strategy (Probe): The Probe strategy was selected because detailed sub-questions are needed to gather focused evidence and resolve diagnostic ambiguity in the image findings.</p> <p>Team Collaboration:</p> <p>Agent-1: [Role] Thoracic Specialist, [Question] Is there any evidence of a pneumothorax, such as air-crescent signs adjacent to the lung edges? Answer-1: There is no evidence of a pneumothorax on this chest X-ray. The lung edges are well-apposed to the ribs without visible air-crescent signs that would suggest the presence of extrapleural air.</p> <p>Agent-2: [Role] Pulmonary Radiologist, [Question] Are there any signs of pulmonary edema or other conditions that could mimic the appearance of pneumonia in the right upper lung field? Answer-2: The chest X-ray demonstrates increased opacity in the right upper lung field, which could be consistent with various conditions including pneumonia. However, specific signs of pulmonary edema such as Kerley B lines, peribronchial cuffing, or a butterfly or bat-wing pattern are not clearly evident.</p> <p>Response from Memory:</p> <p>No signs of pneumothorax are seen, as lung edges are intact without air-crescent signs. In the right upper lung field, the opacity may indicate pneumonia, though typical features of pulmonary edema are absent. The appearance could also reflect atelectasis or a mass, and further imaging would be helpful for clarification. The best match for the finding in this chest X-ray is pneumonia.</p>			

Fig. 4: Case study of CXRAgent versus state-of-the-art methods on visual question answer. **Top:** question, and baseline predictions. **Bottom:** Prediction of CXRAgent. EDV sections are in purple with the conclusion credibility highlighted. Evidence are in blue, the role and task of team members are in orange. Correct results are in green and incorrect ones in red.

duce erroneous outputs without evidence-backed verification.

Also, CXRAgent’s diagnostic planning showed good adaptability. Recognizing the ambiguous findings, it strategically selected the *Probe* collaboration mode. This allowed targeted investigation of critical differential diagnoses through specialized virtual teams. The thoracic specialist definitively ruled out pneumothorax by confirming well-apposed lung edges, while the pulmonary radiologist excluded edema by noting the absence of characteristic signs.

In addition, CXRAgent integrates the memory context for diagnosis, enabling comprehensive synthesis of findings. By

maintaining and referencing prior tool validations and expert analyses, CXRAgent could weigh all evidence to determine the most probable diagnosis while transparently acknowledging diagnostic limitations.

These cases illustrate how CXRAgent’s framework overcomes three fundamental limitations of current systems: (1) the tendency for monolithic models to make absolute but potentially incorrect conclusions, (2) the lack of mechanisms to validate findings against the image evidence, and (3) the inability to adapt reasoning strategies based on case complexity. Through its director-orchestrated, evidence-grounded mul-

tistage reasoning, CXRAgent delivers not only better accuracy but also better clinically realistic, nuanced interpretations that acknowledge diagnostic uncertainty.

V. CONCLUSION

We introduced CXRAgent that conducts evidence-backed and adaptive CXR interpretation through tool invocation, multi-stage reasoning, evidence-driven validation, and team-based collaboration. The core stages are: 1) Tool Invocation: invoking diverse CXR-analysis tools with outputs standardized and verified by an evidence-driven validator; 2) Diagnostic Planning: flexibly planning and assembling specialized expert teams based on task complexity and prior findings; and 3) Collaborative Decision-making: consolidating expert insights and contextual memory into reliable, evidence-backed diagnostic conclusions. Extensive experiments across multiple CXR interpretation tasks demonstrate the superior reliability and adaptability of CXRAgent in handling heterogeneous tool outputs and complex clinical scenarios, highlighting its potential for real-world clinical deployment.

REFERENCES

- [1] T. Tanida, P. Müller, G. Kaissis, and D. Rueckert, "Interactive and explainable region-guided radiology report generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7433–7442.
- [2] Z. Wang, L. Liu, L. Wang, and L. Zhou, "R2GenGPT: Radiology report generation with frozen llms," *Meta-Radiology*, vol. 1, no. 3, p. 100033, 2023.
- [3] S. Bannur *et al.*, "MAIRA-2: grounded radiology report generation," *arXiv preprint arXiv:2406.04449*, 2024.
- [4] Z. Chaves *et al.*, "A clinically accessible small multimodal radiology model and evaluation metric for chest x-ray findings," *Nature Communications*, vol. 16, no. 1, p. 3108, 2025.
- [5] C. Wu, X. Zhang, Y. Zhang, H. Hui, Y. Wang, and W. Xie, "Towards generalist foundation model for radiology by leveraging web-scale 2D&3D medical data," *Nature Communications*, vol. 16, no. 1, p. 7866, 2025.
- [6] Z. Chen *et al.*, "CheXagent: Towards a foundation model for chest x-ray interpretation," in *AAAI 2024 Spring Symposium on Clinical Foundation Models*, 2024.
- [7] D. Ma, J. Pang, M. B. Gotway, and J. Liang, "A fully open AI foundation model applied to chest radiography," *Nature*, pp. 1–11, 2025.
- [8] A. Sellergren *et al.*, "Medgemma technical report," 2025. [Online]. Available: <https://arxiv.org/abs/2507.05201>
- [9] A. Fallahpour, J. Ma, A. Munim, H. Lyu, and B. Wang, "MedRAX: medical reasoning agent for chest x-ray," in *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025.
- [10] Y. Kim *et al.*, "MDAgents: An adaptive collaboration of LLMs for medical decision-making," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [11] B. Li *et al.*, "MMedAgent: Learning to use medical tools with multimodal agent," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, Nov. 2024, pp. 8745–8760.
- [12] J. P. Cohen *et al.*, "Torchxrayvision: A library of chest x-ray datasets and models," in *International Conference on Medical Imaging with Deep Learning*, 2022, pp. 231–249.
- [13] J. Lian *et al.*, "A structure-aware relation network for thoracic diseases detection and segmentation," *IEEE Transactions on Medical Imaging*, vol. 40, no. 8, pp. 2042–2052, 2021.
- [14] K. Dong, Y. Cheng, K. He, and J. Suo, "A generative model uses healthy and diseased image pairs for pixel-level chest x-ray pathology localization," *Nature Biomedical Engineering*, 2025.
- [15] J. Ma, Y. He, F. Li, L. Han, C. You, and B. Wang, "Segment anything in medical images," *Nature Communications*, vol. 15, no. 1, p. 654, 2024.
- [16] T. Yu, B. Ge, S. Wang, Y. Yang, Q. Huang, and J. Yu, "Consistency conditioned memory augmented dynamic diagnosis model for medical visual question answering," *IEEE Journal of Biomedical and Health Informatics*, vol. 29, no. 2, pp. 1357–1370, 2025.
- [17] T. Yu, Z. Tong, J. Yu, and K. Zhang, "Fine-grained adaptive visual prompt for generative medical visual question answering," in *the AAAI Conference on Artificial Intelligence*, 2025, pp. 9662–9670.
- [18] Y. Yang *et al.*, "Token-Mixer: Bind image and text in one embedding space for medical image reporting," *IEEE Transactions on Medical Imaging*, vol. 43, no. 11, pp. 4017–4028, 2024.
- [19] C. Li *et al.*, "LLaVA-Med: Training a large language-and-vision assistant for biomedicine in one day," in *Advances in Neural Information Processing Systems*, vol. 36, 2023, pp. 28 541–28 564.
- [20] X. Liang *et al.*, "MedFILIP: medical fine-grained language-image pre-training," *IEEE Journal of Biomedical and Health Informatics*, 2025.
- [21] M. Moor *et al.*, "Med-flamingo: a multimodal medical few-shot learner," in *Machine Learning for Health (ML4H)*. PMLR, 2023, pp. 353–367.
- [22] L. Team *et al.*, "Lingshu: A generalist foundation model for unified multimodal medical understanding and reasoning," 2025. [Online]. Available: <https://arxiv.org/abs/2506.07044>
- [23] Q. Lin *et al.*, "A foundation model for chest x-ray interpretation with grounded reasoning via online reinforcement learning," 2025. [Online]. Available: <https://arxiv.org/abs/2509.03906>
- [24] Q. Jin *et al.*, "AgentMD: Empowering language agents for risk prediction with large-scale clinical tool learning," *Nature Communications*, vol. 16, no. 1, p. 9377, 2025.
- [25] J. Chen *et al.*, "CoD, towards an interpretable medical agent using chain of diagnosis," in *Findings of the Association for Computational Linguistics: ACL 2025*, Vienna, Austria, Jul. 2025, pp. 14 345–14 368.
- [26] K. Huang *et al.*, "Biomni: A general-purpose biomedical ai agent," *bioRxiv*, pp. 2025–05, 2025.
- [27] X. Tang *et al.*, "MedAgents: Large language models as collaborators for zero-shot medical reasoning," in *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- [28] L. Yue, S. Xing, J. Chen, and T. Fu, "ClinicalAgent: clinical trial multi-agent system with large language model-based reasoning," in *Proceedings of the 15th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*, 2024, pp. 1–10.
- [29] Y. Feng, J. Wang, L. Zhou, and Y. Li, "DoctorAgent-RL: A multi-agent collaborative reinforcement learning system for multi-turn clinical dialogue," *arXiv preprint arXiv:2505.19630*, 2025.
- [30] Q. Zheng *et al.*, "End-to-end agentic RAG system training for traceable diagnostic reasoning," *arXiv preprint arXiv:2508.15746*, 2025.
- [31] Q. Peng *et al.*, "Integration of multi-source medical data for medical diagnosis question answering," *IEEE Transactions on Medical Imaging*, vol. 44, no. 3, pp. 1373–1385, 2025.
- [32] Y. Zhou, L. Song, and J. Shen, "MAM: Modular multi-agent framework for multi-modal medical diagnosis via role-specialized collaboration," in *Findings of the Association for Computational Linguistics: ACL 2025*, Vienna, Austria, Jul. 2025, pp. 25 319–25 333.
- [33] Y. Mao, W. Xu, Y. Qin, and Y. Gao, "CT-Agent: a multimodal-llm agent for 3D CT radiology question answering," *arXiv preprint arXiv:2505.16229*, 2025.
- [34] W. Zhang *et al.*, "Patho-AgenticRAG: Towards multimodal agentic retrieval-augmented generation for pathology vlms via reinforcement learning," *arXiv preprint arXiv:2508.02258*, 2025.
- [35] X. Xie *et al.*, "Prompt-agent-driven integration of foundation model priors for low-count pet reconstruction," *IEEE Transactions on Medical Imaging*, pp. 1–1, 2025.
- [36] S. Yao *et al.*, "ReAct: Synergizing reasoning and acting in language models," in *International Conference on Learning Representations (ICLR)*, 2023.
- [37] X. Hu *et al.*, "Interpretable medical image visual question answering via multi-modal relationship graph learning," *Medical Image Analysis*, vol. 97, p. 103279, 2024.
- [38] A. E. W. Johnson *et al.*, "MIMIC-CXR: A large publicly available database of labeled chest radiographs," *Scientific Data*, vol. 6, p. 317, 2019.
- [39] J. Pan *et al.*, "MedVLM-R1: Incentivizing Medical Reasoning Capability of Vision-Language Models (VLMs) via Reinforcement Learning," in *proceedings of Medical Image Computing and Computer Assisted Intervention – MICCAI 2025*, vol. LNCS 15966, September 2025.
- [40] W. Zhao, C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, "RaTEScore: A metric for radiology report generation," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 15 004–15 019.