

# How Well Can Modern LLMs Act as Agent Cores in Radiology Environments?

Anonymous ACL submission

## Abstract

Large language models (LLMs) hold promise for building accurate and interpretable agentic systems in complex domains like radiology. To evaluate whether modern LLMs can serve as agent cores in radiology settings, we introduce **RadA-BenchPlat**, a comprehensive platform built on 2,200 patient records spanning 6 anatomical regions, 5 imaging modalities, and 2,200 diseases. The dataset includes 24,200 QA pairs and 10 tool categories for radiology task-solving. Our benchmarking of 7 leading LLMs reveals significant gaps: while models such as Claude-3.7-Sonnet achieve 67.1% task completion in routine scenarios, they struggle with complex reasoning and tool coordination. We then apply prompt engineering strategies, yielding an overall 48.2% performance gain ( $p < 0.001$ ) on complex tasks-with **prompt backpropagation** and **multi-agent collaboration** contributing 16.8% ( $p < 0.01$ ) and 30.7% ( $p < 0.001$ ) improvements, respectively. We further enhance robustness via automated tool building, reaching 65.4% success. Our work provides critical benchmarks and actionable strategies for developing reliable radiology AI agents, moving closer to fully automated clinical applications. Code and data are prepared and will be available upon publication.

## 1 Introduction

Recent advancements in Large Language Models (LLMs) have revolutionized domains ranging from natural language processing to computer vision (Wu et al., 2023; Cui et al., 2024a; Liu et al., 2024). Particularly noteworthy is the emergence of LLM-powered agent systems (Tang et al.; Jin et al., 2024; Qin et al., 2023; Cui et al., 2024b), which orchestrate complex task planning and dynamically invoke specialized tools to deliver robust analyses (Mehandru et al., 2024). By positioning the LLM as a central reasoning core capable of integrating external resources, these systems find a compelling yet challenging testbed in

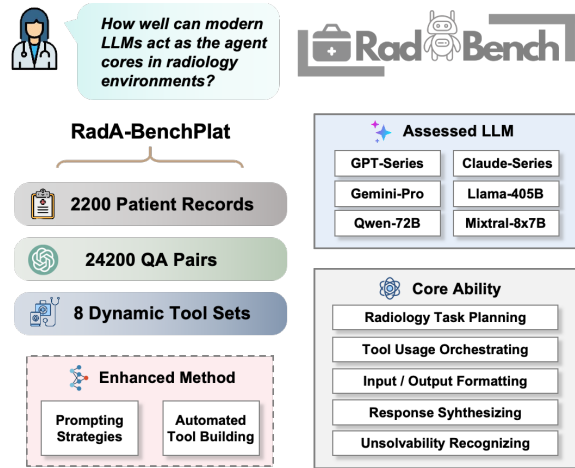


Figure 1: **Benchmark Overview.** The motivation and dataset statistics of the benchmark & Architectural components are integrated into the evaluation platform.

radiology, a field demanding accurate interpretation of complex medical images and detailed textual reporting (Rao et al., 2025; Zhao et al., 2023; Zhou et al., 2025; Lei et al., 2024). While recent efforts have introduced “radiology generalist” models designed to handle diverse imaging tasks within a single framework (Wu et al., 2023; Li et al., 2023; Liu et al., 2023; Zhang et al., 2024, 2023), these monolithic models often struggle with the full spectrum of radiological complexity, which spans imaging modalities, anatomical regions, and evolving diagnostic protocols.

Motivated by the necessity for radiological agents to collaborate with external environments, this paper investigates the feasibility of developing an agent-based system capable of handling varied imaging tasks. Specifically, we explore how effectively current LLMs, serving as an agentic core, can understand professional medical contexts, translate diverse clinical queries into executable plans, and sequentially invoke appropriate tools to complete complex radiological workflows.

To systematically explore this question, we in-

067 introduce **RadA-BenchPlat**, a comprehensive re-  
068 source for evaluating LLM-based agents in radiol-  
069 ogy (Figure 1). The platform is built on three key  
070 components: **First**, we construct a dataset of 2,200  
071 synthetic patient records, spanning 6 anatomical  
072 regions and 5 imaging modalities, all verified as  
073 clinically reasonable by a senior radiologist. **Sec-**  
074 **ond**, for each record, we generate 11 diverse ra-  
075 diology tasks using GPT-4, resulting in a total of  
076 24,200 QA pairs. **Third**, we design a dynamic  
077 toolset simulation strategy to mimic real-world  
078 cross-center variations. This involves dynamically  
079 combining or excluding 10 core tool categories to  
080 create 8 distinct evaluation environments, such as  
081 redundant or insufficient tool sets.

082 Based on the platform, we evaluate 8 lead-  
083 ing LLMs as agentic cores, including the GPT-  
084 series (OpenAI, Accessed by October 12, 2024),  
085 Claude-series (Anthropic, Accessed by October  
086 12, 2024), Gemini Pro (Team et al., 2023),  
087 LLaMA (Touvron et al., 2023), Mixtral (Jiang  
088 et al., 2024), and Qwen (Bai et al., 2023). We fo-  
089 cus on 5 key competencies: *planning tasks, orches-*  
090 *trating tool usage, formatting inputs/outputs, syn-*  
091 *thesizing responses, and recognizing unsolvable*  
092 *cases*. Our results reveal a significant gap between  
093 current LLM capabilities and real-world demands.  
094 While models perform reasonably on basic tasks  
095 like diagnosis, they struggle with integrative tasks;  
096 for instance, Claude 3.7 achieves a 68.3% comple-  
097 tion rate in diagnosis but only 30.1% in compre-  
098 hensive report generation.

099 Building on these findings, we briefly in-  
100 vestigate four prompting strategies—*back-*  
101 *propagation, self-reflection, few-shot learning,*  
102 *and multi-agent collaboration*—which collec-  
103 tively drive a 19.8% overall improvement in  
104 task completion and a substantial 48.2% gain  
105 ( $p < 0.001$ ) in complex scenarios. Beyond these  
106 prompt-level optimizations, we further explore  
107 the potential for self-evolution in radiology agents,  
108 inspired by recent advances in automated tool  
109 building (AutoTB) (Feng et al., 2025; Wölflein  
110 et al., 2025; Cai et al., 2023). We design an  
111 enhancement pipeline that integrates the agentic  
112 system with external AutoTB modules, enabling  
113 the dynamic creation of tools when critical  
114 components are missing. This combined system  
115 achieves an additional 65.4% success rate in  
116 previously unsolvable tasks, underscoring the  
117 vast potential of fully autonomous, self-evolving  
118 agent-based healthcare systems.

## 2 Methods 119

### 2.1 Building RadA-BenchPlat 120

**Patient Records.** A “patient record” here refers 121  
122 to a radiology-centric medical record that includes  
123 patient demographics, medical history, annotated  
124 imaging data (with anatomical and pathological  
125 details), and extensive clinical findings. We cat-  
126 egorize patient records by **anatomy, modality,**  
127 **and disease**. This taxonomy covers 22 common  
128 anatomy-modality combinations derived from 6  
129 anatomical regions and 5 imaging modalities:

**Head & Neck:** {*X-ray, CT, MRI, US*}, 130

**Chest:** {*X-ray, CT, MRI, US*}, 131

**Limb:** {*X-ray, CT, MRI, US*}, 132

**Abd. & Pelv.:** {*X-ray, CT, MRI, US*}, 133

**Spine:** {*X-ray, CT, MRI*}, 134

**Breast:** {*Mammo., MRI, US*}. 135

136 We select 100 common diseases for each of the  
137 22 anatomymodality combinations, resulting in a  
138 total of 2,200 patient records ( $22 \times 100$ ). De-  
139 tailed procedures for record synthesis, expert ver-  
140 ification, and a sample patient record are provided  
141 in the Appendix. A radiologist with over 10  
142 years of experience manually reviewed the records,  
143 confirming that 97.32% met information validity  
144 standards and 96.73% were clinically consistent.  
145 Figure 2a present the demographic distributions  
146 across sex. While the dataset is generally sex-  
147 balanced, breast-related cases are predominantly  
148 female, with 99% of such records corresponding  
149 to women. Analysis of these cases also displays an  
150 age concentration between 45 and 65 years, align-  
151 ing with the broader middle-aged range of 35 to 70  
152 years. Normal distributions in height and weight  
153 further support the datasets authenticity and diver-  
154 sity, establishing a robust foundation for evaluat-  
155 ing LLM agents in complex radiology workflows.

**Task-related QA pairs.** To comprehensively eval- 156  
157 uate LLM-based agents in radiology, we propose  
158 a taxonomy of 11 distinct tasks: (a) organ seg-  
159 mentation, (b) anomaly detection, (c) standard  
160 end-to-end diagnosis, (d) organ and anomaly joint  
161 grounding, (e) diagnosis with grounding clues, (f)  
162 organ-wise biomarker calculation, (g) anomaly-  
163 wise biomarker calculation, (h) standard report  
164 generation, (i) report generation focused on spe-  
165 cific biomarkers, (j) report generation focused on  
166 both biomarkers and indicators<sup>1</sup>, and (k) detailed  
167 treatment planning.

<sup>1</sup>refers to a specific quantification of health status or dis-  
order grading, e.g. CURB-65, tumor grading

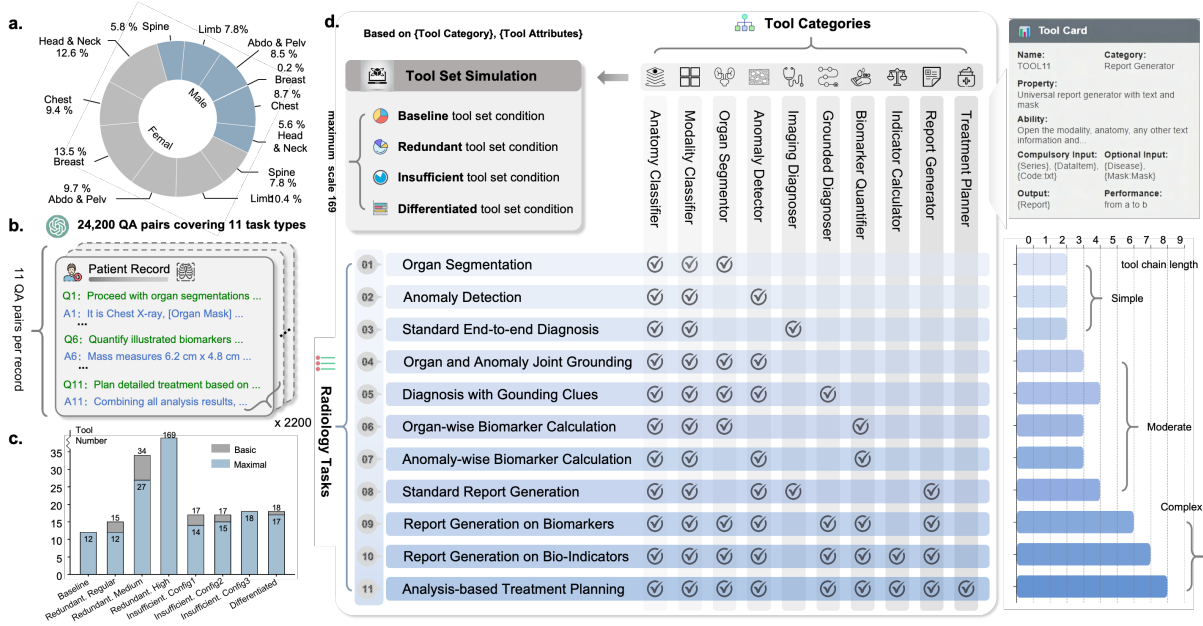


Figure 2: **Statistics in RadA-BenchPlat.** **a.** Distribution of patient sex across different anatomical regions. **b.** Ten categories of tools and eleven types of radiology tasks, with checkmarks indicating tools required for task completion. Tool sets are simulated through various tool attribute combinations. **c.** 24,200 generated QA-pairs based on radiology task types using GPT-4. **d.** Number of tools across eight simulated tool set conditions. **f.** Maximum token length of the open-source and closed-source LLMs used in this study.

As illustrated in Figure 2d, we categorize these tasks into three complexity levels ‘Simple’, ‘Moderate’, and ‘Complex’ based on the number of reasoning steps: fewer than 4, between 4 and 6, and 6 or more, respectively. For instance, organ segmentation is relatively simple, involving anatomy and modality classification followed by segmentation. In contrast, detailed treatment planning is complex, requiring multi-step reasoning including diagnosis, quantification, report synthesis, *etc.* We further simulate each of the 11 tasks across 2,200 patient records (Figure 2b), generating 24,200 QA pairs ( $2,200 \times 11$ ) using a prompt template. For example, a task may prompt “What disease can be inferred?” or “Please write a radiologic report for the image”. The QA pairs are evenly distributed across tasks to ensure balanced and diverse representation, and all undergo manual verification.

**Radiology Tool Sets.** To simulate the real-world clinical settings, we define 10 high-level radiology tool categories (detailed in the supplementary): Anatomy Classifier (AC), Modality Classifier (MC), Organ Segmentor (OS), Anomaly Detector (AD), Imaging Diagnoser (ID), Grounded Diagnoser (GD), Biomarker Quantifier (BQ), Indicator Calculator (IC), Report Generator (RG), and

Treatment Planner (TP). A specific tool is denoted by a “tool card” (Figure 2d), detailing its category, properties, capabilities, required/optional inputs, and performance, with values sampled from clinically plausible ranges. For example, a tool may be described as an “anatomy classifier” with a performance score. This design allows for infinite tool variations, enabling diverse evaluation scenarios.

To reflect varied clinical environments, we define four evaluation tool set conditions: “**Baseline**”, “**Redundant**”, “**Insufficient**”, and “**Differentiated**” (Figure 2d). For each patient record and task query, the “**Baseline**” condition provides a minimal, solvable tool set of 12 tools across 10 categories, ensuring both biomarker quantifier and indicator evaluator categories always include two tools (covering organ evaluation and anomaly assessment). The “**Redundant**” condition introduces extra, potentially irrelevant tools at three redundancy levels: regular (1215 tools, with at most one extra per category), medium (2734 tools, adding 23 extras per category), and high (a comprehensive fixed set of 169 tools), testing the agents ability to filter relevant tools in noisy settings. For “**Insufficient**” conditions, the tool set is unsolvable due to missing essentials in three con-

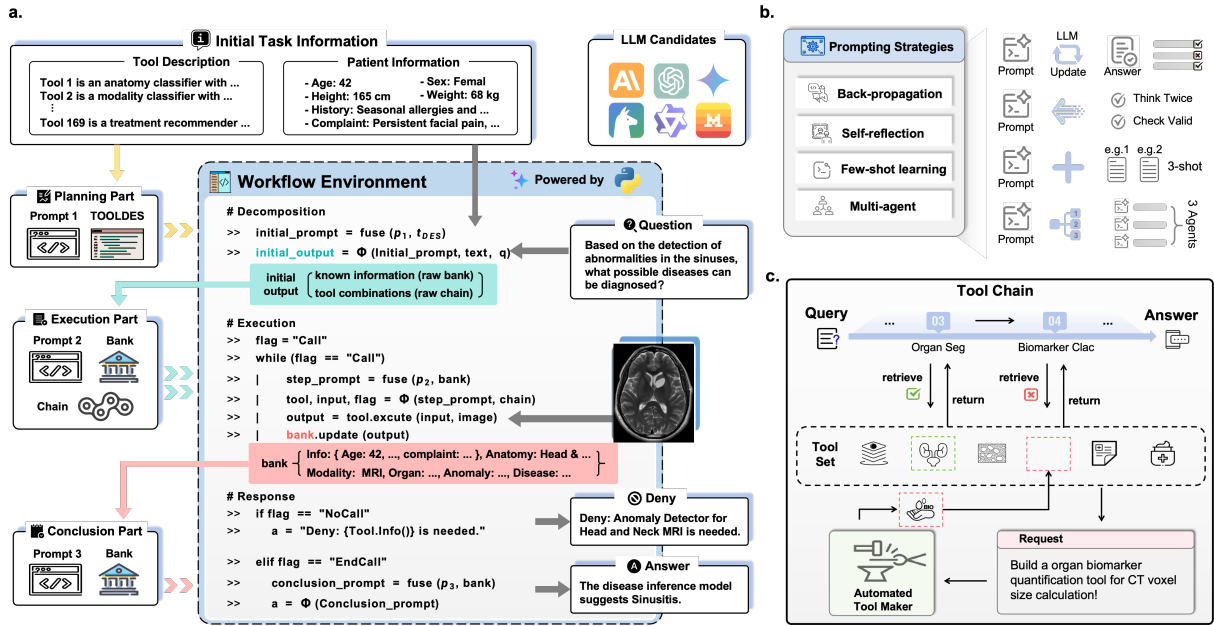


Figure 3: **Detailed agentic workflow and extended methods.** **a.** Three-phase architecture consisting of: initialization phase, circulation phase and conclusion phase. **b.** Four extended prompt engineering methods for performance improvement. **c.** Automated tool building process implemented when no suitable tool is available.

220 figures: config1 omits an entire tool category  
 221 (1417 tools), config2 has modality or anatomy mis-  
 222 matches (1517 tools), and config3 lacks neces-  
 223 sary capabilities (fixed at 18 tools), with increas-  
 224 ing recognition difficulty for mismatched tools.  
 225 Finally, “Differentiated” provides overlapping  
 226 tools (1718 tools) with distinct performance levels,  
 227 requiring agent to choose the most accurate tool.

## 2.2 Benchmarking LLMs as Agent Core

228 **Main Workflow Formulation.** The LLM-based  
 229 agent system consists of three main compo-  
 230 nents (Figure 3a): an LLM agent core ( $\Phi$ ), a set  
 231 of specialized tools ( $\mathcal{T} = \{t_1, t_2, \dots, t_m\}$ ), and a  
 232 memory bank ( $\mathcal{B}$ ). Given a clinical query  $Q$  (e.g.,  
 233 “Please make a diagnosis based on the patient’s in-  
 234 formation”) and patient history, the workflow pro-  
 235 ceeds as follows: The agent core ( $\Phi$ ) first parses  
 236  $Q$  and decomposes it into sub-tasks based on the  
 237 available tools in  $\mathcal{T}$ , initializing  $\mathcal{B}$  with relevant  
 238 patient data. For each sub-task, the core selects  
 239 and executes the appropriate tool from  $\mathcal{T}$ , storing  
 240 outputs in  $\mathcal{B}$  to iteratively refine the process. Fi-  
 241 nally, the agent core synthesizes all information in  
 242  $\mathcal{B}$  to generate a coherent response that addresses  
 243 the clinical query.  
 244

## 2.3 Leveraging Enhanced Methods

245 To address the limited performance of LLMs on  
 246 complex tasks (completion rate: 27.4%), we intro-  
 247

248 duce four prompting enhancement strategies (Fig-  
 249 ure 3b) to boost automation and reliability. In this  
 250 work, we employ **Prompt Back-propagation** by  
 251 iteratively refining prompts through three rounds  
 252 of automated feedback from Claude-3.7-Sonnet  
 253 combined with manual adjustments based on 220  
 254 QA pairs under “Redundant” and “Insufficient”  
 255 tool set conditions, after which the refined prompt  
 256 is applied across the full benchmark. In parallel,  
 257 **Self-reflection** prompts the agent to verify input  
 258 correctness before each tool invocation, a process  
 259 supported by an automated verifier to reduce exe-  
 260 cution errors. Additionally, **Few-shot Learning** is  
 261 used to provide exemplar workflows and tool tem-  
 262 plates that guide planning and standardize execu-  
 263 tion, thereby minimizing misunderstandings and  
 264 formatting issues. Moreover, **Multi-agent Collab-  
 265 oration** restructures the single-agent setup into a  
 266 three-agent system comprising a planner for task  
 267 decomposition, an executor for tool usage and I/O  
 268 management, and a summarizer for final response  
 269 synthesis all leveraging Claude-3.7-Sonnet.

270 In “Insufficient” tool set conditions, tasks be-  
 271 come unsolvable due to missing essential tools.  
 272 Leveraging advances in automated machine learn-  
 273 ing, our system based on the M<sup>3</sup>Builder (Feng  
 274 et al., 2025) framework and enhanced with  
 275 prompting strategies can now detect missing tools  
 276 and generate instructions for dynamic tool con-

Table 1: Target Hit Rates across models under different conditions. Note: Llama-3.1 refers to the 405B version, Mixtral to 8x7B, and Qwen-2.5 to 72B. (**Claude-3.7-sonnet** achieves peak performance in Regular condition).

Condition	GPT-4-Turbo	GPT-4o	Gemini-1.5	Claude-3.5	Claude-3.7	Llama-3.1	Mixtral	Qwen-2.5
Regular	0.286	0.481	0.403	0.519	<b>0.671</b>	0.428	0.056	0.198
Medium	0.217	0.429	0.390	<b>0.450</b>	0.428	0.387	0.000	0.162
High	0.124	0.176	0.085	0.231	<b>0.292</b>	0.093	0.000	0.020

struction (Figure 3c). For instance, if a biomarker quantification tool is absent during lung nodule analysis on chest X-rays, the agent can request a model for biomarker size calculation, triggering automated tool building.

### 3 Results

We perform a quantitative analysis that includes performance on task completion and agentic abilities for RadAgent in radiology tasks.

#### 3.1 Main Observations

##### Observation 1: Decreasing Task Completion Performance from Simple to Complex Tasks

While evaluating LLMs for decomposing radiological tasks and effectively utilize medical tools (Table. 1), task completion rate is used as primary metric, which is defined as successfully solving a query by appropriately utilizing available tools without incurring any intermediate errors during response synthesis. Claude-3.7-Sonnet outperforms GPT-3.5 and GPT-4o, achieving a 67.1% success rate on routine tasks. However, performance drops to 42.8% on advanced tasks and plunges to 29.2% on complex ones, highlighting the increasing difficulty of tasks requiring multi-step reasoning and coordination. A detailed evaluation of the LLMs’ specific capabilities as agent cores is provided in the subsequent sections.

##### Observation 2: Acceptable Planning Ability while Frequent Offset During Execution

To evaluate task planning ability and consistency, we assess tool chain generation under “Redundant” conditions using three metrics: Levenshtein Distance (LD, edit distance between generated and ground-truth tool chains), False Discovery Rate (FDR, proportion of incorrectly included tools), and Tool Matching Accuracy (TMA, step-wise alignment with ground truth). As shown in Table. 2, our findings are: **(1)** closed-source models generally outperform open-source models (except Llama-3.1-405B), achieving LD < 1.4 despite an average chain length of 5; **(2)** differences be-

Table 2: Performance metrics under Redundant Regular, Medium, and High conditions. (↓: Lower is better, ↑: Higher is better).

Model	Exec. ↓ Dist.	Cons. ↓ Dist.	FD Rate ↓	TM Acc. ↑	EC Rate ↑	PFS % ↑
<b>Redundant. Regular Condition</b>						
GPT-4-Turbo	2.02	1.32	0.212	0.640	0.731	0.442
GPT-4o	1.13	0.60	0.106	0.733	0.712	0.667
Gemini-1.5-pro	1.36	0.78	0.154	0.695	0.695	0.519
Claude-3.7-Sonnet	1.31	0.69	0.102	0.763	0.795	0.780
Llama-3.1-405B	1.60	1.43	0.226	0.498	0.682	0.356
Mixtral-8x7B	4.31	3.76	0.117	0.063	0.112	0.093
Qwen-2.5-72B	2.50	2.21	0.313	0.633	0.669	0.669
<b>Redundant. Medium Condition</b>						
GPT-4-Turbo	2.15	1.68	0.209	0.608	0.613	0.413
GPT-4o	1.24	0.65	0.113	0.735	0.679	0.506
Gemini-1.5-pro	1.50	0.86	0.157	0.677	0.713	0.487
Claude-3.7-Sonnet	1.34	0.73	0.159	0.685	0.737	0.512
Llama-3.1-405B	2.28	1.64	0.190	0.475	0.631	0.361
Mixtral-8x7B	0.00	0.00	0.000	0.000	0.000	0.000
Qwen-2.5-72B	2.76	2.41	0.371	0.456	0.560	0.542
<b>Redundant. High Condition</b>						
GPT-4-Turbo	2.23	2.18	0.285	0.519	0.323	0.388
GPT-4o	1.34	0.72	0.121	0.696	0.327	0.266
Gemini-1.5-pro	1.58	3.05	0.417	0.458	0.144	0.231
Claude-3.7-Sonnet	1.35	0.93	0.184	0.599	0.425	0.332
Llama-3.1-405B	2.41	2.72	0.395	0.324	0.158	0.235
Mixtral-8x7B	0.00	0.00	0.000	0.000	0.000	0.000
Qwen-2.5-72B	3.13	3.49	0.331	0.421	0.082	0.284

tween planned and executed chains emerge during multi-iteration execution, with Claude-3.5 and Llama-3.1 showing improved convergence, while others show stable or increased deviations; **(3)** GPT-4o achieves the lowest FDR across all “Redundant” tool set conditions, indicating lower tool selection redundancy; **(4)** GPT-4o and Claude-3.5 yield higher TMA, reflecting better step-wise alignment with ground truth; **(5)** as tool complexity increases, all LLMs exhibit significant performance declines, especially open-source models.

##### Observation 3: Consistent Optimal Tool Selection amidst Tool Diversification

To assess how well agent cores select the best tools from multiple candidates, we use the Optimal Tool Score (OTS), which shows how closely the chosen tools rank in performance relative to available alternatives. Evaluation is conducted under the “Differentiated” tool set condition. As shown in Figure 4a, most LLMs demonstrate strong performance in selecting high-quality tools

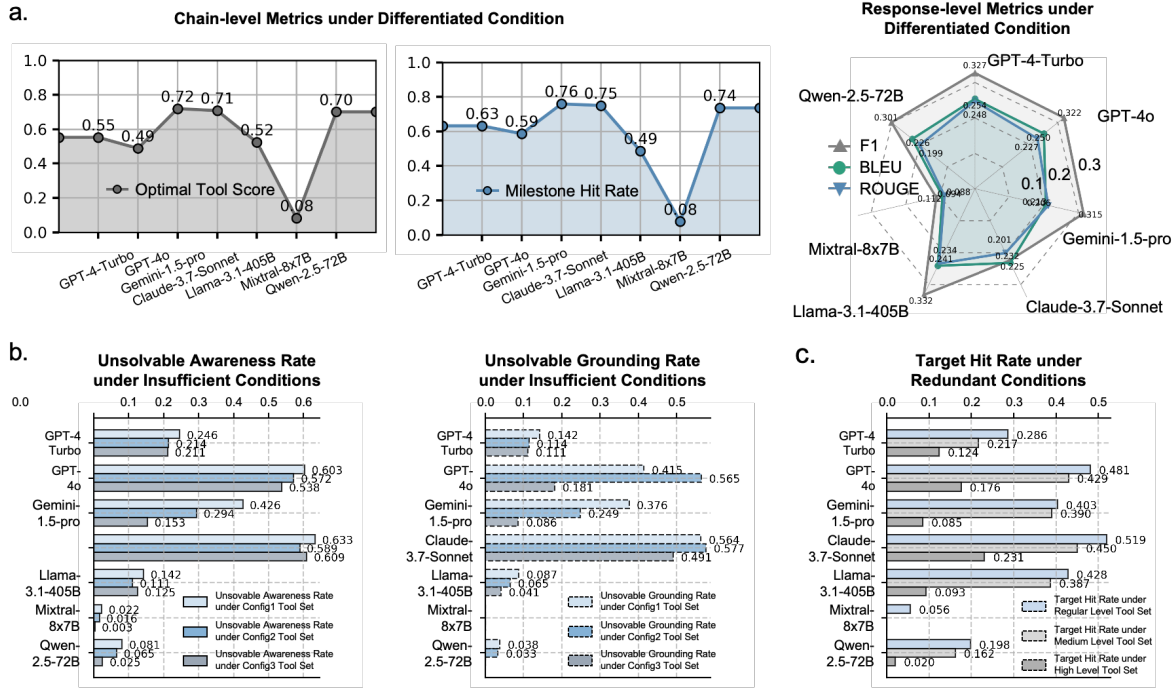


Figure 4: Evaluation results under “Differentiated” and “Redundant” tool set conditions. **a.** Comparative results across leading LLMs under “Differentiated” condition, showing both tool chain-level performance and free text final response-level results. **b.** Performance results under “Insufficient” conditions across LLMs. **c.** Performance results under “Redundant” conditions for leading LLMs.

and completing key task steps. Gemini-1.5-Pro achieves the highest score at 0.720, closely followed by the open-source Qwen-2.5 at 0.701, highlighting their robust comparative reasoning and tool selection capabilities in diverse settings.

#### Observation 4: Limited I/O Management at Large Tool Set Scale

For evaluating data flow management between sequential tools, we use two metrics under “Redundant” tool set conditions: Execution Completion Rate (ECR, the percentage of tool chains completed without I/O errors) and Pre-Failure Success Percentage (PFSP, the proportion of completed chains before failure in unsuccessful cases). Our findings are: (1) as shown in Table. 1, Claude-3.7-sonnet achieves the highest ECR across all “Redundant” conditions (0.795 and 0.737 for regular and medium levels), demonstrating superior execution reliability over other LLMs; (2) in medium and high-level cases, about half of the LLMs reach PFSP > 0.5, indicating that even if the full task isn’t completed, they can execute a substantial part of the tool chain before failure; (3) as tool set complexity increases, all LLMs experience sharp drops in both ECR and PFSP, especially in high-level settings (Table. 1), highlighting major chal-

lenges in managing data flow under complex conditions; (4) Mixtral consistently shows poor performance on both metrics, reflecting its limited ability to accurately organize and manage I/O relationships between tools.

#### Observation 5: Inadequate Response Synthesis due to Target Missing

To assess the generation of correct and coherent responses, we evaluate Target Hit Rate (THR, whether the final tool aligns with the required answer format), Milestone Hit Rate (MHR, success in producing key intermediate outputs), and text similarity metrics (BLEU, ROUGE, F1, measuring alignment with reference responses) under both “Redundant” and “Differentiated” tool set conditions. Our analysis shows: (1) Claude-3.5-Sonnet achieves the highest THR in Redundant” conditions (0.519, 0.450, and 0.231 for regular, medium, and high-level cases, respectively; Figure 4a); (2) even so, the best model only succeeds in 51.9% of tasks under regular “Redundant” settings, mainly due to deviations from the correct tool path and I/O errors; (3) for MHR (Figure 4a), Claude-3.5, GPT-4o, and Qwen-2.5 perform similarly well ( $\approx 0.8$ ) in the “Differentiated” condition, demonstrating strong ability to accomplish

key task steps; (4) all models exhibit low BLEU, ROUGE, and F1 scores (rarely exceeding 0.35), indicating persistent challenges in synthesizing intermediate outputs into coherent final responses, largely due to planning and execution errors.

### Observation 6: Impressive Unsolvability Awareness but Limited at Grounding

To evaluate agent core performance on unsolvable tasks, we use Unsolvability Awareness Rate (UAR, the ability to recognize when a task is impossible) and Unsolvability Grounding Rate (UGR, the ability to identify what is missing) under increasingly challenging “Insufficient” conditions: config1 (missing tool category), config2 (missing modality-specific tools), and config3 (tools present but insufficient). As shown in Figure 4c, Claude-3.5-Sonnet and GPT-4o achieve the highest UARs (around 0.6), indicating strong capability to detect unsolvable cases. However, accurate grounding of failure causes is more challenging; while these two models maintain relatively high UGR across all levels, most others such as Gemini-1.5-Pro show notable drops, especially under subtle “Deficient” settings. Success requires not only recognizing infeasibility but also specifying what is missing (e.g., absent tool categories in config1, missing modality-region combinations in config2, or minimum tool capabilities in config3). Many models, including GPT-4-Turbo, LLaMA-3.1-405B, Mixtral, and Qwen-2.5, frequently attempt to solve tasks with inappropriate tools, reflecting limited awareness of tool constraints and difficulty in accurately judging task solvability.

## 3.2 Performance Enhancement

### Integrated Prompting Strategies Yield Significant Performance Gains

Integrating four prompting strategies (detailed in supp.) significantly boosts task completion rates (Figure. 5). In simple and moderate radiology tasks where Claude-3.7 and 3.5 already perform well, the strategies yield extra gains of 0.095 and 0.058, leading to total improvements of 14.4% and 13.9% over baseline. For complex tasks, the impact is more pronounced with improvements up to 48.2%. Each strategy contributes independently, but their combined effect is not strictly additive due to diminishing returns. The results highlight that prompt refinement, self-reflection, and few-shot learning are essential for effective reasoning and guidance, while multi-agent collaboration is

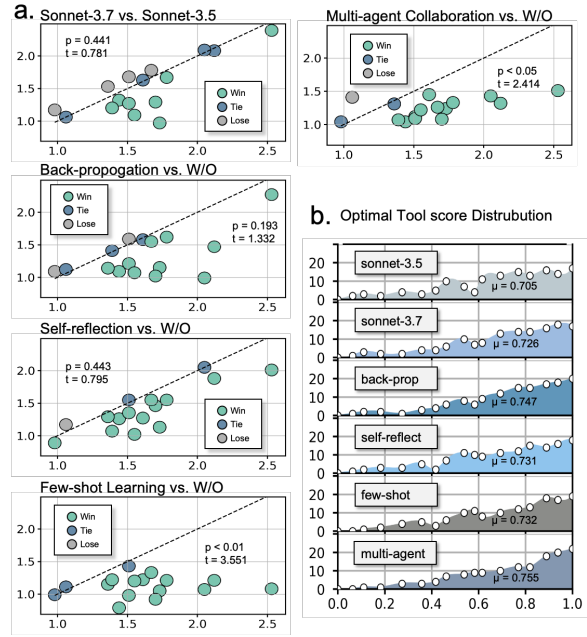


Figure 5: Performance improvement comparison among prompt engineering strategies. a. Significance comparison between strategy groups, where “W/O” represents pure Sonnet-3.7 without engineering, and “Win” indicates the performance method before “vs.” outperforming the latter. b. Performance distribution across six prompting strategies under “Differentiated” tool set condition.

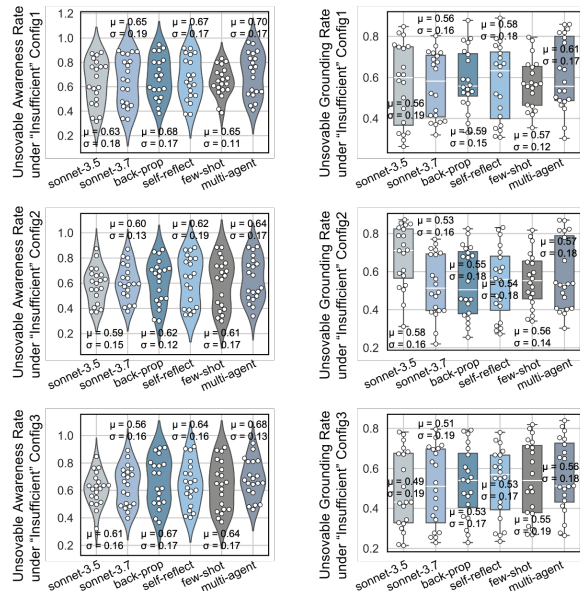


Figure 6: Performance distribution across six prompting strategies under several “Insufficient” tool set conditions.

critical for managing long-context tasks. However, the added prompt complexity can occasionally lead to performance inconsistencies.

## Multi-agent and Back-propagation Strategies Drive the Most Impact

In our main evaluation (Figure 5a), we compared Claude-3.5, Claude-3.7, and Claude-3.7 with different prompting strategies on 220 QA pairs from 20 patient records over five trials, using average Levenshtein distance per patient and defining a win as outperforming the baseline. Self-reflection and back-propagation showed some improvement, but without statistical significance ( $p > 0.05$ ); in contrast, few-shot learning and multi-agent collaboration achieved statistically significant gains ( $p < 0.05$ ), with multi-agent collaboration leading in most comparisons. Further tests under “Insufficient” and “Differentiated” conditions (Figure 6) show that multi-agent collaboration improved unsolvability awareness and grounding rates by about 5%, and back-propagation by about 3%, while occasional declines were observed with self-reflection and back-propagation alone. In the “Differentiated” setting, among 139 QA pairs and six methods, prompting strategies consistently raised optimal tool scores, with back-propagation and multi-agent contributing the largest improvements (score increases of 0.021 and 0.029, Figure 5b).

### Relaxing Tool Constraints through Automated Tool Building

Evaluated on 220 QA pairs, the system achieves a 65.4% completion rate for routine tasks under “Insufficient” conditions, dropping to 44.4% for advanced and 26.2% for complex tasks. The main limitation is the difficulty in accurately identifying missing tool capabilities, especially under “Insufficient config3”, leading to erroneous tool building requests. Nonetheless, the system shows strong potential by autonomously building tools and executing end-to-end tasks without human intervention, representing a significant step toward fully automated radiology workflows. Further details are provided in the supplementary material.

## 4 Related Works

**LLM-based Agents in Medicine.** The evolution of medical AI has transitioned from static classifiers to dynamic LLM-based agents capable of autonomous reasoning. To enhance clinical validity, recent frameworks have adopted collaborative multi-agent strategies (Tang et al.) or integrated external tools for evidence-based planning (Jin et al., 2024; Wang et al., 2025), thereby

addressing the hallucination and precision limitations of generalist models. In the imaging domain, architectures such as ChatCAD (Wang et al., 2024) utilize LLMs as controllers to orchestrate specialized vision models. However, existing benchmarks (Jiang et al., 2025) predominantly focus on textual EHR tasks or isolated classification, lacking a unified evaluation of the multi-modal, tool-integrated workflows essential for radiological agency.

**Synthetic Data Generation and Validation.** Parallel to agent development, addressing the data scarcity-privacy paradox has driven a shift from traditional de-identification to ab-initio synthesis. While early methods relied on masking real clinical notes (Ren et al., 2025), contemporary approaches utilize persona-based prompting to generate high-fidelity patient cohorts and dialogues from scratch (Haider et al., 2025; Yilmaz et al., 2025). To ensure reliability, validation frameworks have evolved beyond simple n-gram overlap, adopting rigorous metrics such as the Quality-Privacy Score (QPS) to quantify the trade-off between semantic realism and re-identification risk (Sella et al., 2025; Adams et al., 2025). Despite these advancements in general clinical text, frameworks specifically designed to synthesize the structured, multi-modal reporting required for complex radiological error detection remain under-explored (Sun et al., 2025).

## 5 Conclusion

We introduced RadA-BenchPlat to evaluate the feasibility of Modern LLMs acting as agent cores within radiology workflows. Our extensive benchmarking reveals that while current SOTA models—particularly the Claude-3 series—demonstrate potential in simplified scenarios, they currently lack the robustness required for autonomous execution of complex clinical tasks due to planning deviations and tool orchestration failures. However, we demonstrate that integrating adaptive prompting strategies, such as multi-agent collaboration and prompt back-propagation, alongside automated tool building, significantly mitigates these limitations. These findings highlight a promising trajectory for future research: moving beyond standalone LLMs toward engineered, self-correcting agentic systems capable of meeting the rigorous demands of real-world radiology.

## 6 Limitations

Our study is intentionally scoped to establish foundational benchmarks for LLM-based agents in radiology. Consequently, we focus on a controlled synthetic dataset to isolate and evaluate reasoning capabilities without the confounding noise of real-world data; future work will involve validating our framework in clinical settings. Furthermore, we benchmark general-purpose LLMs due to their superior tool-use and planning abilities, which are prerequisites for our agentic framework. Extending this evaluation to medical-specific models and integrating multi-modal (imaging) inputs are critical next steps as these technologies mature.

## References

- Tim Adams, Colin Birkenbihl, Karen Otte, Hwei Geok Ng, Jonas Adrian Rieling, Anatol-Fiete Näher, Ulrich Sax, Fabian Prasser, and Holger Fröhlich. 2025. On the fidelity versus privacy and utility trade-off of synthetic patient data. *iScience*, 28(5).
- Anthropic. Accessed by October 12, 2024. [\[link\]](#).
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. 2023. Large language models as tool makers. *arXiv preprint arXiv:2305.17126*.
- Hejie Cui, Lingjun Mao, Xin Liang, Jieyu Zhang, Hui Ren, Quanzheng Li, Xiang Li, and Carl Yang. 2024a. Biomedical visual instruction tuning with clinician preference alignment. *arXiv preprint arXiv:2406.13173*.
- Hejie Cui, Zhuocheng Shen, Jieyu Zhang, Hui Shao, Lianhui Qin, Joyce C Ho, and Carl Yang. 2024b. LLMs-based few-shot disease predictions using ehr: A novel approach combining predictive agent reasoning and critical agent instruction. *arXiv preprint arXiv:2403.15464*.
- Jinghao Feng, Qiaoyu Zheng, Chaoyi Wu, Ziheng Zhao, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2025. M<sup>3</sup>builder: A multi-agent system for automated machine learning in medical imaging. *arXiv preprint arXiv:2502.20301*.
- Syed Ali Haider, Srinivasagam Prabha, Cesar Abraham Gomez-Cabello, Sahar Borna, Ariana Genovese, Maissa Trabelsy, Bernardo G Collaco, Nadia G Wood, Sanjay Bagaria, Cui Tao, and 1 others. 2025. Synthetic patient-physician conversations simulated by large language models: A multi-dimensional evaluation. *Sensors*, 25(14):4305.

- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, and 1 others. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Yixing Jiang, Kameron C Black, Gloria Geng, Danny Park, James Zou, Andrew Y Ng, and Jonathan H Chen. 2025. Medagentbench: A realistic virtual ehr environment to benchmark medical llm agents. *arXiv preprint arXiv:2501.14654*.
- Qiao Jin, Zhizheng Wang, Yifan Yang, Qingqing Zhu, Donald Wright, Thomas Huang, W John Wilbur, Zhe He, Andrew Taylor, Qingyu Chen, and 1 others. 2024. Agentmd: Empowering language agents for risk prediction with large-scale clinical tool learning. *arXiv preprint arXiv:2402.13225*.
- Jiayu Lei, Xiaoman Zhang, Chaoyi Wu, Lisong Dai, Ya Zhang, Yanyong Zhang, Yanfeng Wang, Weidi Xie, and Yuehua Li. 2024. Autorg-brain: Grounded report generation for brain mri. *arXiv preprint arXiv:2407.16684*.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Zhengliang Liu, Yiwei Li, Peng Shu, Aoxiao Zhong, Longtao Yang, Chao Ju, Zihao Wu, Chong Ma, Jie Luo, Cheng Chen, and 1 others. 2023. Radiology-llama2: Best-in-class large language model for radiology. *arXiv e-prints*, pages arXiv-2309.
- Nikita Mehandru, Brenda Y Miao, Eduardo Rodriguez Almaraz, Madhumita Sushil, Atul J Butte, and Ahmed Alaa. 2024. Evaluating large language models as agents in the clinic. *NPJ digital medicine*, 7(1):84.
- OpenAI. Accessed by October 12, 2024. [\[link\]](#).
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, and 1 others. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Vishwanatha M Rao, Michael Hla, Michael Moor, Subathra Adithan, Stephen Kwak, Eric J Topol, and Pranav Rajpurkar. 2025. Multimodal generative ai for medical image interpretation. *Nature*, 639(8056):888-896.
- Libo Ren, Samuel Belkadi, Lifeng Han, Warren Del-Pinto, and Goran Nenadic. 2025. Synthetic4health: generating annotated synthetic clinical letters. *Frontiers in Digital Health*, 7:1497130.

650	Nadir Sella, Florent Guinot, Nikita Lagrange, Laurent-Philippe Albou, Jonathan Desponds, and Hervé Isambert. 2025. Preserving information while respecting privacy through an information theoretic framework for synthetic health data generation. <i>npj Digital Medicine</i> , 8(1):49.	
656	Cong Sun, Kurt Teichman, Yiliang Zhou, Brian Critelli, David Nauheim, Graham Keir, Xindi Wang, Judy Zhong, Adam E Flanders, George Shih, and 1 others. 2025. Generative large language models trained for detecting errors in radiology reports. <i>Radiology</i> , 315(2):e242575.	
662	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. In <i>ICLR 2024 Workshop on Large Language Model (LLM) Agents</i> .	
668	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	
674	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
680	Sheng Wang, Zihao Zhao, Xi Ouyang, Tianming Liu, Qian Wang, and Dinggang Shen. 2024. Interactive computer-aided diagnosis on medical image using large language models. <i>Communications Engineering</i> , 3(1):133.	
685	Ziyue Wang, Junde Wu, Linghan Cai, Chang Han Low, Xihong Yang, Qiaxuan Li, and Yueming Jin. 2025. Medagent-pro: Towards evidence-based multi-modal medical diagnosis via reasoning agentic workflow. <i>arXiv preprint arXiv:2503.18968</i> .	
690	Georg Wölflein, Dyke Ferber, Daniel Truhn, Ognjen Arandjelović, and Jakob Nikolas Kather. 2025. Llm agents making agent tools. <i>arXiv preprint arXiv:2502.11705</i> .	
694	Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Towards generalist foundation model for radiology. <i>arXiv preprint arXiv:2308.02463</i> .	
698	Abdurrahim Yilmaz, Furkan Yuceyalcin, Rahmetullah Varol, Ece Gokyayla, Ozan Erdem, Donghee Choi, Ali Anil Demircali, Gulsum Gencoglan, Joram M Posma, and Burak Temelkuran. 2025. A synthetic data generation framework for scalable and resource-efficient medical ai assistants. <i>medRxiv</i> , pages 2025–05.	
	Kai Zhang, Rong Zhou, Eashan Adhikarla, Zhiling Yan, Yixin Liu, Jun Yu, Zhengliang Liu, Xun Chen, Brian D Davison, Hui Ren, and 1 others. 2024. A generalist vision–language foundation model for diverse biomedical tasks. <i>Nature Medicine</i> , pages 1–13.	705 706 707 708 709 710
	Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-vqa: Visual instruction tuning for medical visual question answering. <i>arXiv preprint arXiv:2305.10415</i> .	711 712 713 714 715
	Ziheng Zhao, Yao Zhang, Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. One model to rule them all: Towards universal segmentation for medical images with text prompts. <i>arXiv preprint arXiv:2312.17183</i> .	716 717 718 719 720
	Xiao Zhou, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2025. Knowledge-enhanced visual-language pretraining for computational pathology. In <i>Computer Vision – ECCV 2024</i> , pages 345–362, Cham. Springer Nature Switzerland.	721 722 723 724 725 726
	<b>A Appendix</b>	727
	<b>A.1 Terminologies Explanation</b>	728
	Key terminologies used throughout our work are defined in Extended Figure. 7.	729 730
	<b>A.2 Diseases/abnormalities Names for Patient Records</b>	731 732
	We consider 100 common conditions for each anatomy-modality combination in our synthetic patient records. The considered diseases/abnormalities names are listed in the CSV file which will be made available upon publication.	733 734 735 736 737
	<b>A.3 Detailed Explanations for Tool Categories</b>	738 739
	We explain the 10 tool categories here:	740
	<ul style="list-style-type: none"> <li>• <b>Anatomy Classifier (AC)</b> is a classifier to predict what anatomy region the input image is shot on.</li> <li>• <b>Modality Classifier (MC)</b> is a classifier to predict what imaging modality the input image is based on.</li> <li>• <b>Organ Segmentor (OS)</b> is a tool to predict the dense segmentation masks for a certain organ on the input image.</li> <li>• <b>Anomaly Detector (AD)</b> is a detection tool to predict the anomaly region masks or box coordinates for a certain anomaly type based on the input image.</li> </ul>	741 742 743 744 745 746 747 748 749 750 751 752 753

## Detailed Explanation of Important Terminologies In this Paper

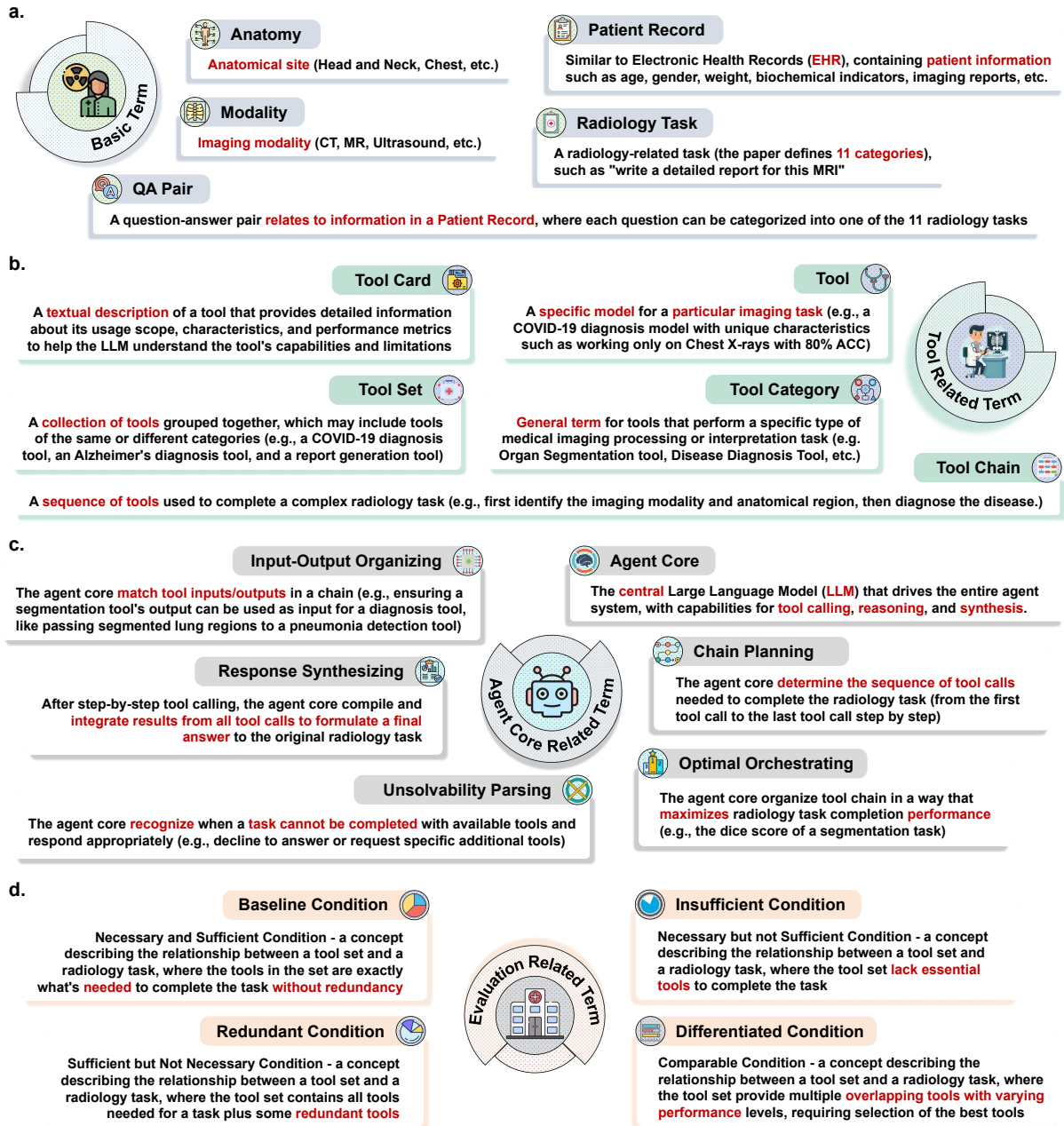


Figure 7: A detailed explanation of 20 important terminologies used in this paper. **a.** Basic terms commonly used throughout the paper. **b.** Sophisticated definitions of tool-related concepts. **c.** Agent-related terms including five key abilities. **d.** Four radiology environments featuring various tool combinations for different evaluation purposes.

754	• <b>Imaging Diagnoser (ID)</b> is a diagnosis tool to predict related diseases solely based on the input image.	804
755		805
756		806
757	• <b>Synthetic Diagnoser (SD)</b> is a diagnosis tool to predict related diseases synthetically based on the input image and text information.	807
758		808
759		809
760	• <b>Biomarker Quantifier (BQ)</b> is a biomarker calculation tool to exactly calculate the radiology biomarker based on the input image and organ-or-anomaly-wise dense masks.	810
761		811
762		812
763		813
764	• <b>Indicator Evaluator (IE)</b> is an indicator (like tumor grading) calculation tool to exactly calculate or predict some medical indicators based on the input image and organ-or-anomaly-wise dense masks.	814
765		815
766		816
767		817
768		818
769	• <b>Report Generator (RG)</b> is a text generation tool to predict radiology reports.	819
770		820
771	• <b>Treatment Recommendation (TR)</b> is a treatment recommendation tool (or system) to provide treatment recommendation based on current clinical findings.	821
772		822
773		823
774		824
775		825
776		826
777		827
778		828
779		829
780		830
781		831
782		832
783		833
784		834
785		835
786		836
787		837
788		838
789		839
790		840
791		841
792		842
793		843
794		844
795		845
796		846
797		847
798		848
799		849
800		850
801		851
802		852
803		853
		854

#### A.4 Detailed Explanations for Task Types

We explain the 11 radiology tasks here with their detailed tool category chains:

- **Organ mask annotation** is a task targeting predicting a dense mask for a certain organ. **Chain:** *anatomy classification modality classification organ segmentation*
- **Anomaly region annotation** is a task targeting predicting a regional mask for a certain anomaly type. **Chain:** *anatomy classification modality classification anomaly detection*
- **Organ-wise biomarker calculation** is a task to calculate or measure some organ-wise biomarkers. **Chain:** *anatomy classification modality classification organ segmentation biomarker quantification*
- **Anomaly-wise biomarker calculation** is a task to calculate or measure some anomaly-wise biomarkers. **Chain:** *anatomy classification modality classification anomaly detection anomaly quantification*
- **Organ-and-anomaly-wise image interpretation** is a task to interpret the image and answer the user query on organs or anomaly regions or both. **Chain:** *anatomy classification modality classification [organ segmentation, anomaly detection]*
- **Disease diagnosis without visual clues** is a task to end-to-end make a diagnosis based

- on a radiology-central patient record. **Chain:** *anatomy classification modality classification disease diagnosis*
- **Disease diagnosis with visual clues** is a task to make a diagnosis based on not only original images but also relevant segmentation annotations (organs or anomalies). **Chain:** *anatomy classification modality classification [organ segmentation, anomaly detection] disease inference*
- **Common report generation** is a basic radiology routine task to interpret a certain radiology image into free-text descriptions. **Chain:** *anatomy classification modality classification anomaly detection disease diagnosis report generation*
- **Report generation focusing on specific biomarkers** is a report generation task variant emphasizing reflecting certain biomarker conditions in the final report. **Chain:** *anatomy classification modality classification [organ segmentation, anomaly detection] [organ biomarker quantification, anomaly biomarker quantification] report generation*
- **Report generation focusing on specific biomarkers and indicators** is a report generation task variant emphasizing reflecting certain biomarkers and indicator conditions in the final report. **Chain:** *anatomy classification modality classification [organ segmentation, anomaly detection] disease diagnosis [organ biomarker quantification, anomaly quantification] indicator evaluation report generation*
- **Treatment planning** is a task to provide a treatment plan for a patient. **Chain:** *anatomy classification modality classification [organ segmentation, anomaly detection] disease diagnosis [organ biomarker quantification, anomaly quantification] indicator evaluation report generation treatment recommendation*

#### A.5 Feature Distributions in Patient Record

Extended Figure 8e visualizes anomalies, diseases, biomarkers, and indicators within patient records. Here, “biomarkers” refer to imaging features (e.g., dimensions, textures), while “indicators” correspond to clinical classifications (e.g., cancer staging) or scoring systems (e.g., CURB-65). We embed these attributes using BioLORD and MedCPT. The t-SNE plot reveals a uniform distribution and confirming that our dataset covers a wide range of clinical scenarios.

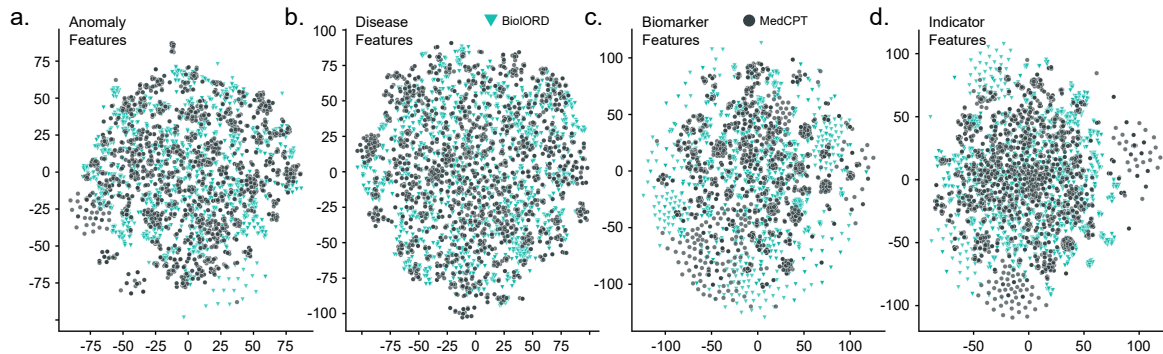


Figure 8: Features distribution on anomaly/disease/biomarker/indicator extracted by BioLORD and MedCPT.

### A.6 Context Token Consumption Analysis

As shown in Extended Figure 9, we present the multi-turn response token lengths generated by two LLMs, GPT-4o and Llama-3.1, representing closed-source and open-source models, respectively. In most conditions, the total token lengths for both models are comparable, ranging from 3,000 to 30,000 tokens. However, in the Insufficient tool set conditions, GPT-4o exhibits significantly longer context lengths than Llama, primarily due to Llama’s frequent failure to execute proper denial responses, resulting in premature terminations.

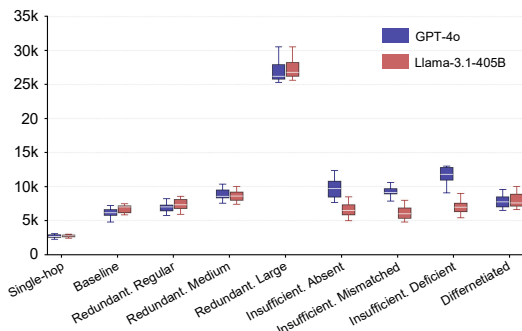


Figure 9: Context token length in 8 simulated conditions using GPT-4o/Llama-3.1-405B as the agent core.

### A.7 Performance on Baseline Tool Set Condition

1) closed-source models generally outperform open-source models (except Llama-3.1-405B) across both “Baseline” tool set conditions, with closed-source models achieving LD between planning and ground truth less than 1.5 (average 1.5 edits needed to match ground truth chains) as

shown in Extended Figure 10 - this represents relatively good performance given the average chain length of 5 step; 2) as shown in Extended Figure 10, differences emerge between planned and executed chains during multi-iteration execution, where models can adjust tool selections - Claude-3.5 and Llama-3.1 show convergence toward ground truth chains, while others maintain or increase deviations.

### A.8 A Patient Record Example

We, here, demonstrate a concrete synthetic patient record in our RadABench-Data:

Listing 1: Patient Record Example

```

"Information": {
  "Age": "42",
  "Sex": "Female",
  "Height": "165",
  "Weight": "68",
  "History": "Patient has a history of
    seasonal allergies and recurrent upper
    respiratory infections",
  "Complaint": "Persistent facial pain, nasal
    congestion, and headache for the
    past 2 weeks"
},
"Anatomy": "Head and Neck",
"Modality": "X-ray",
"Anomaly": {
  "Part": "Maxillary sinuses",
  "Symptom": "Opacification"
},
"Disease": "Sinusitis",
"OrganBiomarker": {
  "OrganObject": "Maxillary sinus",
  "OrganDim": "density",
  "OrganQuant": "+40 Hounsfield Units"
},
"AnomalyBiomarker": {
  "AnomalyObject": "Opacification",
  "AnomalyDim": "intensity",
  "AnomalyQuant": "80% increase compared to
    normal airspace"
}

```

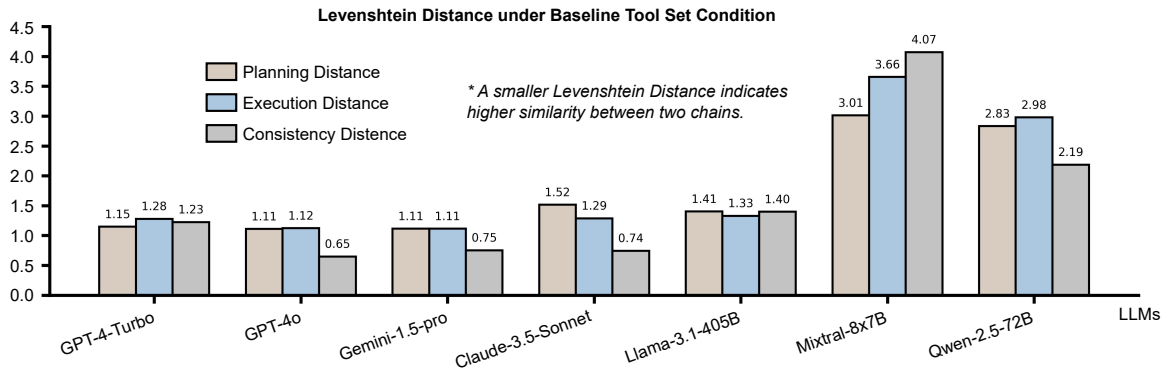


Figure 10: Three Levenshtein Distance between pairs of decision, execution, and ground truth tool chains under NS. Condition among 7 LLMs.

```

919 "Indicator": {
920   "Name": "Lund-Mackay Score",
921   "Value": "8 (Moderate sinusitis)"
922 },
923 "Report": {
924   "Finding": "X-ray of the paranasal sinuses
925     demonstrates bilateral maxillary
926     sinus opacification. The right maxillary
927     sinus shows complete opacification,
928     while the left maxillary sinus demonstrates
929     air-fluid levels. Frontal and ethmoid
930     sinuses appear clear. No evidence of bone
931     erosion or destruction. Nasal septum
932     appears midline. Soft tissues of the face
933     and neck are unremarkable.",
934   "Impression": "Findings consistent with
935     bilateral maxillary sinusitis, more
936     pronounced on the right side. No evidence of
937     complications such as orbital or
938     intracranial involvement."
939 },
940 "Treatment": "Given the patient's symptoms and
941   radiographic findings, a diagnosis
942   of acute bacterial sinusitis is likely. Initial
943   treatment should include a 10-14 day
944   course of broad-spectrum antibiotics such as
945   amoxicillin-clavulanate or, in case of
946   penicillin allergy, a respiratory
947   fluoroquinolone. Adjunctive treatments
948   include nasal
949   saline irrigation, intranasal corticosteroids,
950   and oral decongestants for symptom
951   relief. The patient should be advised to stay
952   well-hydrated and use over-the-counter
953   pain relievers as needed. If symptoms persist
954   or worsen after 72 hours of antibiotic
955   therapy, reassessment is warranted. A follow-up
956   appointment should be scheduled in 2-3
957   weeks to ensure resolution of symptoms. If
958   recurrent episodes occur, further evaluation
959   with CT imaging and potential referral to an
960   ENT specialist for consideration of
961   endoscopic sinus surgery may be necessary."

```

963 which is generated by prompting GPT-4 lever-  
 964 aging the following prompts:

#### Listing 2: Patient Record Generation Prompt

```

965 You are an experienced clinical radiologist.
966 Your task is to generate a detailed
967 medical case based on a hypothetical 256x256
968 medical image. I will provide you with
969 the Anatomy and Modality of the image, as well
970 as an overall Disease. Using this
971 information, you should create a comprehensive
972 case report including patient
973 information, specific anomalies, biomarkers,
974 indicators, a radiology report, and
975 treatment recommendations.
976 Please structure your response using the
977 following template:
978
979 <Case>
980   <Information>
981     <Age> [Number without units] </Age>
982     <Sex> [Male / Female] </Sex>
983     <Height> [Number in cm] </Height>
984     <Weight> [Number in kg] </Weight>
985     <History> [Brief descriptive text]
986     </History>
987     <Complaint> [Brief descriptive text]
988     </Complaint>
989   </Information>
990   <Anatomy> [Head and Neck / Chest / Breast /
991     Abdomen and Pelvis / Limb / Spine]
992   </Anatomy>
993   <Modality> [CT / MRI / X-ray / Ultrasound /
994     Mammography] </Modality>
995   <Anomaly>
996     <Part> [Specific location of anomaly
997       (e.g., right upper lobe of lung)]
998     </Part>
999     <Symptom> [Type of anomaly (e.g.,
1000       nodule)] </Symptom>
1001   </Anomaly>
1002   <Disease> [Corresponding disease name]
1003   </Disease>
1004   <OrganBiomarker>
1005     <OrganObject> [A specific organ serving
1006       as biomarker] </OrganObject>
1007     <OrganDim> [number / length / size /
1008       volume / angle / density / intensity]

```

1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081

```

/
texture] </OrganDim>
<OrganQuant> [Specific quantitative
value] </OrganQuant>
</OrganBiomarker>
<AnomalyBiomarker>
<AnomalyObject> [The same Anomaly
described in the Anomaly Symptom
serving
as biomarker] </AnomalyObject>
<AnomalyDim> [number / length / size /
volume / angle / density / intensity
/
texture] </AnomalyDim>
<AnomalyQuant> [Specific quantitative
value] </AnomalyQuant>
</AnomalyBiomarker>
<Indicator>
<Name> [Name of the indicator (e.g.,
Lung Cancer TNM Staging Score)]
</Name>
<Value> [Specific value or grade (e.g.,
cT2aN0M0 (Stage IB))] </Value>
</Indicator>
<Report>
<Finding> [Findings section in the style
of a MIMIC-CXR report] </Finding>
<Impression> [Impression section in the
style of a MIMIC-CXR report]
</Impression>
</Report>
<Treatment> [A paragraph including
diagnostic procedures, medication
recommendations, and follow-up suggestions]
</Treatment>
</Case>

Guidelines for generating the case report:

Information: Provide realistic patient
demographics, medical history, and chief
complaint.
Anomaly: Describe a specific anomaly consistent
with the given anatomy, modality,
and disease.
Biomarkers: Choose a relevant organ biomarker
and the anomaly biomarker mentioned
before that can be observed in the image and
provide plausible quantitative value.
Indicator: Calculate an appropriate indicator
based on the patient information and
biomarker value. Provide a specific, medically
accurate score or grade.
Report: Generate a concise radiological report
in the style of MIMIC-CXR, with
separate Findings and Impression sections.
Treatment: Offer a comprehensive treatment plan
including diagnostic procedures,
medication recommendations, and follow-up care.
Ensure that all parts of the case are medically
accurate and consistent with each
other. Use your expertise as a radiologist to
provide realistic and detailed
information throughout the case.

Given the following parameters:
- Anatomy: {ANATOMY}
- Imaging Modality: {MODALITY}
- Disease: {DISEASE}
Please generate this patient record.
```

**A.9 A Tool Card Example** 1083

Here is an example of a tool card of a report generation tool focused on Head and Neck X-ray images: 1084  
1085  
1086

Listing 3: Tool Card Example

```

=== Tool Description for TOOL17 ===
Name: TOOL17
Category: Report Generator
Ability: Given the Head and Neck X-ray Image,
any other text information and
organ/anomaly masks and labels, generate a
radiology report.
Property: Report Generator only suitable for
Head and Neck X-ray image with
Text and Mask
Compulsory Input: ['$Image$']
Optional Input: ['$Information$',
'$OrganObject$', '$AnomalyObject$',
'$Disease$',
'$OrganDim$', '$OrganQuant$', '$AnomalyDim$',
'$AnomalyQuant$', '$IndicatorName$',
'$ValueName$', '$OrganMask$', '$AnomalyMask$']
Output: ['$Report$']
Performance: Score from 0.4 to 0.88, increases
with optional inputs
```

1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107

**A.10 A Task Simulation Example** 1109

Here is an example of QA pair generation from a chest X-ray record of a pneumonia patient: 1110  
1111

Listing 4: QA Pair Example

```

<Q1> Identify and segment the lung fields in
this chest X-ray. </Q1>
<A1> The lung fields are segmented [{\em
organ mask}: left lung, right lung].
Both lungs show clear boundaries. </A1>
...
<Q11> Based on the imaging findings and
clinical indicators, what treatment plan
would you recommend? </Q11>
<A11> Given the moderate pneumonia severity
(CURB-65 score: 2) and [anomaly mask],
recommend oral antibiotics and follow-up
chest X-ray in 2 weeks. </A11>
```

1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127

which is generated by the following prompt, 1129  
leveraging GPT-4: 1130

Listing 5: QA Pair Generation Prompt

```

Assume a clinical medical imaging scenario
where you, as the Agent core, play the role
of a doctor. Given a patient's radiological
image, you want to complete different tasks
by calling various tools. There are ten tools
in total (numbered 0 to 9):

TOOLKIT
0. |Modality Classifier|
Property: A classification model
Ability: Determine the modality of the Image.
Input: [Image]
```

1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143

1144	Output: [Modality]	tools.	1214
1145		Input: [Image] & [Modality] & [Anatomy]	1215
1146	1.  Anatomy Classifier	Optional Input: ([Organ Mask] & [Organ Label]), ([Anomaly Mask] & [Anomaly Label]),	1216
1147	Property: A classification model,	[Diseases], ([Object] & [Dim] & [Quant]),	1217
1148	Ability: Determine the anatomy of the Image.	([Indicators] & [Values])	1218
1149	Input: [Image]	Output: [Report]	1219
1150	Output: [Anatomy]		1220
1151			1221
1152	2.  Organ Segmentation Model		1222
1153	Property: A segmentation model	9.  Treatment Recommendation Model	1223
1154	Ability: Given the modality and anatomy,	Property: a language model	1224
1155	segment all organs in the Image	Ability: Recommend personalized treatment	1225
1156	(can not segment any lesion or abnormality).	plans based on all results processed by	1226
1157	Input: [Image] & [Modality] & [Anatomy]	former tools and the patient's information.	1227
1158	Output: [Organ Mask] & [Organ Label]	Input: [Image] & [Priors] & [Modality] & [Anatomy] & [Diseases]	1228
1159		Optional Input: ([Organ Mask] & [Organ Label]), ([Anomaly Mask] & [Anomaly Label]),	1229
1160	3.  Anomaly Detection Model	([Object] & [Dim] & [Quant]), ([Indicators] & [Values])	1230
1161	Property: A detection model	Output: [Treatment]	1231
1162	Ability: Given the modality and anatomy,		1232
1163	determine the location and type of		1233
1164	abnormality.		1234
1165	Input: [Image] & [Modality] & [Anatomy]		1235
1166	Output: [Anomaly Mask] & [Anomaly Label]		1236
1167		There are 11 different tasks and the chain of	1237
1168	4.  Disease Diagnosis Model	tools each task hopes to break down into:	1238
1169	Property: A classification model		1239
1170	Ability: Given the modality and anatomy,	(1) Basic Image Analysis and Organ Segmentation	1240
1171	diagnose diseases directly from the	ToolUse: 012	1241
1172	input	Description: Perform basic image analysis	1242
1173	image.	and segment organs within the medical	1243
1174	Input: [Image] & [Modality] & [Anatomy]	image.	1244
1175	Output: [Diseases]		1245
1176		(2) Basic Image Analysis and Anomaly Detection	1246
1177	5.  Disease Inference Model	ToolUse: 013	1247
1178	Property: A Inference model	Description: Perform basic image analysis	1248
1179	Ability: Infer disease based on organ	and detect anomalies within the medical	1249
1180	segmentation and anomaly detection	image.	1250
1181	results.		1251
1182	Input: [Image] & [Organ Mask] & [Organ Label] & [Anomaly Mask] & [Anomaly Label]	(3) Image-based Direct Disease Diagnosis	1252
1183	Output: [Diseases]	ToolUse: 014	1253
1184		Description: Diagnose disease directly from	1254
1185		the medical image without intermediate	1255
1186		steps.	1256
1187	6.  Biomarker Quantification Model		1257
1188	Property: A quantification model	(4) Organ segmentation and anomaly localization	1258
1189	Ability: Given the organ region or anomaly	ToolUse: 0123	1259
1190	region and the biomarker of interest,	Description: Segment organs and locate	1260
1191	estimate its property. The biomarker can be	anomalies within the medical image.	1261
1192	either organ or anomaly. The dimension		1262
1193	can be one of the number, length, size,	(5) Anomaly-based Disease Diagnosis	1263
1194	volume, angle, density, intensity or	ToolUse: 01235	1264
1195	texture of the organ or anomaly.	Description: Diagnose disease based on	1265
1196	Input: [Image] & [Object] & [Dim] & [Biomarker Mask] & [Biomarker Label]	disease inference by finding	1266
1197	Output: [Quant]	abnormalities	1267
1198		and the organ in which they occur.	1268
1199			1269
1200	7.  Indicator Evaluation Model	(6) Organ Biomarker Quantification	1270
1201	Property: A calculation model	ToolUse: 0126	1271
1202	Ability: Use prior patient information and	Description: Quantify specific biomarkers	1272
1203	several biomarkers values to calculate	related to organs in the medical image.	1273
1204	the indicator, the indicator can be a score		1274
1205	or a grading.	(7) Anomaly Biomarker Quantification	1275
1206	Input: [Priors] & [Indicator] & [Biomarkers] & [Quants]	ToolUse: 0136	1276
1207	Output: [Value]	Description: Quantify specific biomarkers	1277
1208		related to anomalies in the medical	1278
1209		image.	1279
1210	8.  Report Generation Model		1280
1211	Property: A multimodal model	(8) Disease and Anomaly Based Report Generation	1281
1212	Ability: Generate a medical report by	ToolUse: 01348	1282
1213	integrating results processed by former	Description: Generate a medical report based	1283

on detected diseases and anomalies.

(9) Disease and Biomarker Based Report Generation  
 ToolUse: 0123568  
 Description: Generate a comprehensive report incorporating anomaly detection, disease diagnosis and biomarker quantification.

(10) Comprehensive Evaluation Report Generation  
 ToolUse: 01235678  
 Description: Generate a detailed evaluation report including all aspects of the medical image analysis.

(11) Comprehensive Report Generation and Treatment Recommendations  
 ToolUse: 012356789  
 Description: Generate a comprehensive report including all analysis results and treatment recommendations.

Please generate 11 mutually independent question-answer pairs corresponding to tasks above, based on the different task natures and the content of the case. Specifically, strictly avoid including information in the questions that should be determined by the tools (such as imaging modality, specific anatomy, or precise abnormality types). Pay attention: create task-specific question-answer pairs that highlight the unique tool usage patterns for different tasks. The questions should be:

1. Naturally aligned with the task description
2. Representative of real-world scenarios for that task type
3. Questions must not reveal information about modality or anatomy that should be determined by Tools 0 and 1 or other tools. The questions should be phrased in a way that necessitates the use of these basic identification tools.

Provide answers in the most concise free-text form possible. If visual results such as masks are involved, please embed them in the text in the form of [Organ Mask] or [Anomaly Mask]. (eg. The organ segmentation result is shown as [Organ Mask].) The 11 generated question-answer pairs should follow this template:

```
<Q1> ... </Q1>
<A1> ... </A1>
<Q2> ... </Q2>
<A2> ... </A2>
...
<Q11> ... </Q11>
<A11> ... </A11>
```

The answer only needs to provide a simple final result based on the information already available in the Case, without showing the thought process. I will now provide you with a Case containing all the information.

The patient record is {Patient Record}, please generate the corresponding QA pairs.

### A.11 Detailed Workflow Prompts in RadABench-EvalPlat

In this part we detail the prompt template used in the main workflow of our RadABench-EvalPlat.

#### Task Decomposition.

Listing 6: Task Decomposition Prompts

```
# Medical Image Analysis Assistant

## Task Overview
You are a radiological agent analyzing medical images.
For each query, you will receive:
  1. A medical imaging examination (Image) of a patient (assume already provided)
  2. Known patient Information including demographics, medical history, and main complaints.

Your task involves three sequential parts:
  1. Problem Decomposition (Part 1)
  - Identify available information
  - Break down the question into sequential steps
  2. Sequential Tool Application (Part 2)
  - Execute one tool at a time
  - Record each tool's output
  - Continue until sufficient information is gathered
  3. Solution Synthesis (Part 3)
  - Integrate all results
  - Generate final answer

## Available Information Categories
The following categories must be used exactly as written:

['$Information$', '$Anatomy$', '$Modality$', '$Disease$', '$OrganObject$', '$OrganDim$', '$OrganQuant$', '$AnomalyObject$', '$AnomalyDim$', '$AnomalyQuant$', '$IndicatorName$', '$IndicatorValue$', '$Report$', '$Treatment$']

Where:
- $Information$: Patient demographics (e.g., "45-year-old male", "BMI: 24", "history of diabetes")
- $Anatomy$: Body part (e.g., "chest", "abdomen", "brain")
- $Modality$: Imaging technique (e.g., "X-ray", "CT", "MRI")
- $Disease$: Medical condition (e.g., "pneumonia", "cancer", "fracture")
- $OrganObject$: Organ to measure (e.g., "liver", "heart")
- $OrganDim$: Organ measurement type (e.g., "number", "length", "size", "volume", "angle", "density", "intensity", "texture")
```

```

1419 - $OrganQuant$: Organ measurement value (e.g.,
1420 "5cm", "120ml")
1421 - $AnomalyObject$: Abnormality to measure
1422 (e.g., "tumor", "fracture")
1423 - $AnomalyDim$: Abnormality measurement type
1424 (e.g., "number", "length", "size",
1425 "volume", "angle", "density", "intensity",
1426 "texture")
1427 - $AnomalyQuant$: Abnormality measurement value
1428 (e.g., "2cm", "5ml")
1429 - $IndicatorName$: Clinical indicator name
1430 - $IndicatorValue$: Clinical indicator value
1431 - $Report$: Medical report content
1432 - $Treatment$: Treatment recommendations
1433
1434 ## Available Tool Categories
1435 Tool categories must be used exactly as written:
1436
1437 [*Anatomy Classification Tool*, *Modality
1438 Classification Tool*,
1439 *Organ Segmentation Tool*, *Anomaly Detection
1440 Tool*,
1441 *Disease Diagnosis Tool*, *Disease Inference
1442 Tool*,
1443 *Organ Biomarker Quantification Tool*,
1444 *Anomaly Biomarker Quantification Tool*,
1445 *Indicator Evaluation Tool*,
1446 *Report Generation Tool*, *Treatment
1447 Recommendation Tool*]
1448
1449 ## Response Format for Part 1
1450 For each query, respond ONLY with:
1451
1452 Known Info: [list any categories explicitly
1453 mentioned in the query]
1454 Tool Chain: [list required tools connected by
1455 ->]
1456
1457 ## Examples
1458
1459 Query 1: "For a straightforward approach to
1460 diagnosing the patient's condition
1461 based on her symptoms and the image, what
1462 diseases can be directly identified?"
1463 Response:
1464 Known Info: []
1465 Tool Chain: [*Anatomy Classification Tool* ->
1466 *Modality Classification Tool*
1467 -> *Disease Diagnosis Tool*]
1468
1469 Query 2: "This 45-year-old male's chest CT
1470 shows a 2cm nodule in the right lung.
1471 Can you give a report?"
1472 Response:
1473 Known Info: ['$Information$', '$Anatomy$',
1474 '$Modality$', '$AnomalyObject$',
1475 '$AnomalyDim$', '$AnomalyQuant$']
1476 Tool Chain: [*Organ Segmentation Tool* ->
1477 *Anomaly Detection Tool*
1478 -> *Disease Inference Tool* ->
1479 *Report Generation Tool*]
1480 (because some information is
1481 provided, so
1482 *Anatomy Classification Tool*,
1483 *Modality Classification Tool*,
1484 *Anomaly Biomarker Quantification
1485 Tool* are optimized.)
1486
1487 ## Important Rules
1488 1. Assume the medical image is already provided

```

```

2. Use exact item category names with $$ as
   listed (e.g., '$Anatomy$')
3. Use exact tool category names with ** as
   listed
   (e.g., '*Anatomy Classification Tool*')
4. Only respond with Part 1 analysis - Parts 2
   & 3 will be addressed
   in subsequent interactions
5. Include only the categories explicitly
   mentioned in the query
6. Connect tools using -> symbol

Please wait for my query.
When provided, analyze it following the format
shown in the examples above.

{Patient Record}
{Query}

```

**Tool Selection & Execution.**

Listing 7: Tool Selection & Execution Prompts

```

# Next Step Planning
## Current Status
Current results dictionary: {value_dict}

## Planning Guidelines
1. Reference your high-level tool chain from
   Part 1 decomposition
2. Consider current results to refine specific
   tool selection
3. Maintain sequential progression according to
   planned workflow
4. Adjust tool selection if needed based on
   intermediate results
5. Check if the tool category is missing when
   this category of tools is required
6. Check if the tool is suitable for the
   detected Anatomy and Modality in reserved
   value dictionary based on the Tool description
   Ability and Property
7. Check if the result in reserved value
   dictionary can be derived from each tool
   used
in former steps based on the limited label list
described in Tool Ability
8. If no suitable tool exists, identify which
   type of denial applies:
   - Missing tool category
   - Missing specific modality-anatomy tool
   - Insufficient tool capability

## Input Requirements
1. Required inputs: Must include all mandatory
   inputs specified in tool description
2. Optional inputs: Include if available and
   beneficial to tool performance
3. Do not include variables that are not
   relevant to the tool's function
4. All variables must exist in the current
   results dictionary
5. Use proper $$ notation for all variables

## Response Format
For ongoing analysis (if not final step):
<Call>
<Purpose>Brief, clear statement of this

```

1557	step's goal in context of overall	descriptions	1627
1558	analysis</Purpose>	3. Mark all variables with \$\$ notation	1628
1559	<Tool>TOOL[number] - must match available	4. Include only existing variables from results	1629
1560	specific tools</Tool>	dictionary	1630
1561	<Input>['\$variable1\$', '\$variable2\$', ...] -	5. Keep purpose statements clear and concise	1631
1562	use only existing variables from	6. Brief response only includes one Call,	1632
1563	results</Input>	EndCall, or NoCall XML part without	1633
1564	</Call>	additional explanations	1634
1565		7. For NoCall responses, use the appropriate	1635
1566	For final step only:	format based on denial type	1636
1567	<EndCall>		1637
1568	<Purpose>Brief, clear statement of this	## Decision Making Process	1638
1569	final step's goal</Purpose>	1. Review planned tool chain from Part 1	1639
1570	<Tool>TOOL[number] - must match available	2. Check current results in value dictionary	1640
1571	specific tools</Tool>	3. Check if tool category is missing when this	1641
1572	<Input>['\$variable1\$', '\$variable2\$', ...] -	category of tools is required	1642
1573	use only existing variables from	4. Check tool Ability and Property in detail to	1643
1574	results</Input>	judge its suitability for detected	1644
1575	</EndCall>	Anatomy and Modality in the value dictionary	1645
1576		5. Check if the result in reserved value	1646
1577	When a tool category is completely missing:	dictionary can be derived from each tool	1647
1578	<NoCall>	used	1648
1579	<Purpose>The purpose requiring a missing	in former steps based on the limited label list	1649
1580	tool category</Purpose>	described in Tool Ability	1650
1581	<Category>The missing category from	6. Evaluate tool availability and capability:	1651
1582	['Anatomy Classifier', 'Modality	- Is the required tool category available?	1652
1583	Classifier',	- Are specific tools available for the	1653
1584	'Organ Segmentor', 'Anomaly Detector',	needed modality-anatomy combination?	1654
1585	'Disease Diagnoser', 'Disease	- Do available tools have sufficient	1655
1586	Inferencer',	capabilities?	1656
1587	'Biomarker Quantifier', 'Indicator	7. If tools are insufficient, use appropriate	1657
1588	Evaluator', 'Report Generator',	NoCall format	1658
1589	'Treatment	8. If tools are available, select and format	1659
1590	Recommender']</Category>	appropriate Call/EndCall	1660
1591	<Anatomy>Universal</Anatomy>	9. Use <EndCall> only if this will be the final	1661
1592	<Modality>Universal</Modality>	step	1662
1593	<Ability>CategoryMissing</Ability>	Please provide your next step based on:	1663
1594	</NoCall>	- Original tool chain plan	1664
1595		- Current results	1665
1596	When specific tools for a modality-anatomy	- Available specific tools	1666
1597	combination are missing:	- Remaining analysis needs	1667
1598	<NoCall>	- Tool availability and capability assessment	1668
1599	<Purpose>The purpose requiring a specific		1669
1600	modality-anatomy tool</Purpose>		
1601	<Category>The required tool		
1602	category</Category>		
1603	<Anatomy>The specific anatomy from		
1604	['Universal', 'Head and Neck', 'Chest',		
1605	'Breast', 'Abdomen and Pelvis', 'Limb',		
1606	'Spine']</Anatomy>		
1607	<Modality>The specific modality from		
1608	['Universal', 'X-ray', 'CT', 'MRI',		
1609	'Ultrasound']</Modality>		
1610	<Ability>SpecificToolMissing</Ability>		
1611	</NoCall>		
1612			
1613	When existing tools lack required capabilities:		
1614	<NoCall>		
1615	<Purpose>The purpose requiring advanced		
1616	capabilities</Purpose>		
1617	<Category>The category of existing but		
1618	insufficient tools</Category>		
1619	<Anatomy>The relevant anatomy</Anatomy>		
1620	<Modality>The relevant modality</Modality>		
1621	<Ability>InsufficientCapability</Ability>		
1622	</NoCall>		
1623			
1624	## Format Requirements		
1625	1. Maintain proper XML structure		
1626	2. Use exact tool numbers as specified in tool		
		Response Generation.	1671
		Listing 8: Response Generation Prompt	
		Based on your Part 1 analysis plan, Part 2 tool	1672
		execution sequence, and the final	1673
		results dictionary {value_dict}, provide:	1674
			1675
		1. A concise answer to the initial question	1676
		2. Key supporting evidence from your results	1677
		3. How your findings align with the planned	1678
		analysis	1679
			1680
		Keep your response brief and focused on	1681
		directly answer the initial question.	1682
			1683
			1684
		A.12 An In-depth Case Study	1686
		Here we select a representative case from the	1687
		agent core evaluation workflow to demonstrate the	1688
		process, where we test the agent cores' perfor-	1689
		mance using a QA pair based on the task type	1690

1691 "Anomaly-wise Biomarker Calculation" and a pa- 1750  
 1692 tient record (a 62-year-old female with hyperten- 1751  
 1693 sion, osteoarthritis of the knees, and chronic neck 1752  
 1694 pain with stiffness radiating to right shoulder) 1753  
 1695 under the NR. Denyl2 tool set condition, using 1754  
 1696 Claude-3.5-Sonnet as the agent core, documenting 1755  
 1697 each decision it makes. We can observe both the 1756  
 1698 strengths and limitations reflected by this example. 1757  
 1699 **Initial Input.** The agent core receives the 1758  
 1700 patient information along with a virtual image. 1759  
 1701 For our evaluation task "Anomaly-wise Biomarker 1760  
 1702 Calculation", we generate a corresponding QA 1761  
 1703 pair based on the patient's complete medical 1762  
 1704 record. The initial input is provided below: 1763

Listing 9: Initial Input

```
1705 $Information$: {
1706   "Age": "62",
1707   "Sex": "Female",
1708   "Height": "165",
1709   "Weight": "72",
1710   "History": "Hypertension, osteoarthritis of
1711             the knees",
1712   "Complaint": "Chronic neck pain and
1713               stiffness, radiating pain to right
1714               shoulder"
1715 },
1716
1717 $Query$: From an anomaly perspective in a
1718         specific medical image, after identifying
1719         the type and area, could you quantify specific
1720         biomarker characteristics?
1721
1722 $Image$: 'PLACEHOLDER_IMAGE'
```

1725 **Task Decomposition.** The agent core processes 1778  
 1726 these inputs by analyzing the patient information 1779  
 1727 and identifying key data for extraction. It de- 1780  
 1728 termines the appropriate high-level tool chain re- 1781  
 1729 quired for task completion. The agent core then 1782  
 1730 stores all relevant information in a memory bank 1783  
 1731 and establishes a tool category chain to guide sub- 1784  
 1732 sequent execution. The process is illustrated be- 1785  
 1733 low: 1786

Listing 10: Stage 1: Task Decomposition

```
1734 Initial Output: Known Info: []
1735
1736 Tool Chain: [*Anatomy Classification Tool* ->
1737             *Modality Classification Tool* ->
1738             *Anomaly Detection Tool* -> *Anomaly Biomarker
1739             Quantification Tool*]
1740
1741 Initial Value Dict: {'$Image$':
1742                     'PLACEHOLDER_IMAGE', '$Information$':
1743                     'PLACEHOLDER_
1744                     INFORMATION'}
1745
1746 Initial Score Dict: {'$Image$': 1.0,
1747                     '$Information$': 1.0}
1748
1749 Initial Fixed Dict:
1750                     frozendict.frozendict({'$Image$': 1.0,
1751                     '$Information$': 1.0})
```

```
High-level Tool chain: Anatomy Classification 1750
Tool -> Modality Classification Tool -> 1751
Anomaly Detection Tool -> Anomaly Biomarker 1752
Quantification Tool 1753
Memory bank: {'$Image$': 'PLACEHOLDER_IMAGE', 1754
              '$Information$': 'PLACEHOLDER_ 1755
              INFORMATION'} 1756
Score bank: {'$Image$': 1.0, '$Information$': 1757
             1.0} 1758
Fixed bank: frozendict.frozendict({'$Image$': 1759
                                   1.0, '$Information$': 1.0}) 1760
```

1764 **Tool Selection & Execution:** The first phase of tool 1765  
 1766 selection and execution commences. At each step, 1767  
 1768 the agent core integrates the planned high-level 1768  
 1769 tool chain with stored information from the mem- 1769  
 1770 ory bank to generate a tool API call, specifying 1770  
 1771 the purpose, tool name, and input parameters. The 1771  
 1772 Platform then executes the API call by activating 1772  
 1773 the selected tool and processing the input. Upon 1773  
 1774 successful execution, the tool's output is stored in 1774  
 1775 the memory bank. We configure the tool set ac- 1775  
 1776 cording to the NR. Denyl2 condition, where there 1776  
 1777 is a mismatch between the anomaly detection tool 1777  
 1778 and the provided *Image* characteristics (Head and 1778  
 1779 Neck X-ray) in terms of anatomy and modality. 1779  
 1780 The established tool set is as follows: 1780

Listing 11: Tool Set Description

```
1781 "TOOL1": {
1782   "Name": "TOOL1",
1783   "Category": "Anatomy Classifier",
1784   "Ability": "Determine the anatomy of the
1785             Image.",
1786   "Property": "Universal Anatomy Classifier",
1787   "Compulsory Input": [
1788     "$Image$"
1789 ],
1790   "Optional Input": [],
1791   "Output": [
1792     "$Anatomy$"
1793 ],
1794   "Performance": "Score from 0.95 to 0.95,
1795                 increases with optional inputs"
1796 },
1797 "TOOL2": {
1798   "Name": "TOOL2",
1799   "Category": "Modality Classifier",
1800   "Ability": "Determine the modality of the
1801             Image.",
1802   "Property": "Universal Modality Classifier",
1803   "Compulsory Input": [
1804     "$Image$"
1805 ],
1806   "Optional Input": [],
1807   "Output": [
1808     "$Modality$"
1809 ],
1810   "Performance": "Score from 0.95 to 0.95,
1811                 increases with optional inputs"
1812 }
```

```

1813 ...
1814
1815 "TOOL5": {
1816   "Name": "TOOL5",
1817   "Category": "Anomaly Detector",
1818   "Ability": "Given the Limb Ultrasound Image,
1819     determine the location and
1820     type of abnormality.",
1821   "Property": "Anomaly Detector only suitable
1822     for Limb Ultrasound image",
1823   "Compulsory Input": [
1824     "$Image$"
1825   ],
1826   "Optional Input": [],
1827   "Output": [
1828     "$AnomalyMask$",
1829     "$AnomalyObject$"
1830   ],
1831   "Performance": "Score from 0.8 to 0.8,
1832     increases with optional inputs"
1833 },
1834 "TOOL6": {
1835   "Name": "TOOL6",
1836   "Category": "Anomaly Detector",
1837   "Ability": "Given the Breast MRI Image,
1838     determine the location and type
1839     of abnormality.",
1840   "Property": "Anomaly Detector only suitable
1841     for Breast MRI image",
1842   "Compulsory Input": [
1843     "$Image$"
1844   ],
1845   "Optional Input": [],
1846   "Output": [
1847     "$AnomalyMask$",
1848     "$AnomalyObject$"
1849   ],
1850   "Performance": "Score from 0.8 to 0.8,
1851     increases with optional inputs"
1852 },
1853 ...
1854
1855 "TOOL8": {
1856   "Name": "TOOL8",
1857   "Category": "Disease Diagnoser",
1858   "Ability": "Given the Head and Neck X-ray
1859     Image, diagnose the disease.",
1860   "Property": "Disease Diagnoser only suitable
1861     for Head and Neck X-ray image",
1862   "Compulsory Input": [
1863     "$Image$"
1864   ],
1865   "Optional Input": [],
1866   "Output": [
1867     "$Disease$"
1868   ],
1869   "Performance": "Score from 0.75 to 0.75,
1870     increases with optional inputs"
1871 },
1872 ...
1873
1874 "TOOL10": {
1875   "Name": "TOOL10",
1876   "Category": "Biomarker Quantifier",
1877   "Ability": "Measure the anomaly biomarker of
1878     the Image.",
1879   "Property": "Universal Anomaly Biomarker
1880     Quantifier",

```

```

"Compulsory Input": [
  "$Image$",
  "$AnomalyObject$",
  "$AnomalyMask$"
],
"Optional Input": [
  "$AnomalyDim$"
],
"Output": [
  "$AnomalyDim$",
  "$AnomalyQuant$"
],
"Performance": "Score from 0.75 to 0.8,
  increases with optional inputs"
},
...

```

In this initial step, the agent core selects the anatomy classification tool. We observe that the anatomical information is successfully stored in the memory bank after the tool execution:

Listing 12: Tool Selection & Execution Step 1

```

Step 1 starts!
Stepoutput: Based on the tool chain from Part 1
  (Anatomy Classification Tool ->
  Modality Classification Tool -> Anomaly
  Detection Tool -> Anomaly Biomarker
  Quantification Tool) and current results
  dictionary, I'll proceed with the first
  step:
<Call>
  <Purpose>Determine the anatomy region of the
    medical image</Purpose>
  <Tool>TOOL1</Tool>
  <Input>['$Image$']</Input>
</Call>
Call Dict: {'Purpose': 'Determine the anatomy
  region of the medical image', 'Tool':
  'TOOL1', 'Input': ['$Image$']}
Memory bank: {'$Image$': 'PLACEHOLDER_IMAGE',
  '$Information$': 'PLACEHOLDER_
  INFORMATION', '$Anatomy$': 'Head and Neck'}
Score bank: {'$Image$': 1.0, '$Information$':
  1.0, '$Anatomy$': 0.95}
Fixed bank: frozendict.frozendict({'$Image$':
  1.0, '$Information$': 1.0})
Step 1 completes!

```

In the second step, the modality classification tool is selected and executed. The successful execution is demonstrated below. As we can observe that the modality information is added to memory bank after tool execution.

Listing 13: Tool Selection & Execution Step 2

```

Step 2 starts!
Stepoutput: Based on the tool chain and current
  results showing Head and Neck anatomy,
  the next step is to determine the image
  modality:

```

1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965

```
<Call>
  <Purpose>Determine the imaging modality of
    the medical image</Purpose>
  <Tool>TOOL2</Tool>
  <Input>['$Image$']</Input>
</Call>
Call Dict: {'Purpose': 'Determine the imaging
  modality of the medical image', 'Tool':
  'TOOL2', 'Input': ['$Image$']}
Memory bank: {'$Image$': 'PLACEHOLDER_IMAGE',
  '$Information$': 'PLACEHOLDER_
  INFORMATION', '$Anatomy$': 'Head and Neck',
  '$Modality$': 'X-ray'}
Score bank: {'$Image$': 1.0, '$Information$':
  1.0, '$Anatomy$': 0.95, '$Modality$':
  0.95}
Fixed bank: frozendict.frozendict({'$Image$':
  1.0, '$Information$': 1.0})
Step 2 completes!
```

1967  
1968  
1969  
1970  
1971  
1972

In the third step, contrary to the expected utilization of the organ segmentation tool, the disease diagnosis tool is invoked instead. While this tool selection is irrelevant to the current task, the execution completes without errors. The tool call is detailed below:

Listing 14: Tool Selection & Execution Step 3

1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999

```
Step 3 starts!
Stepoutput: Given the current results showing
  Head and Neck anatomy and X-ray modality,
  the next step is to detect any anomalies in the
  medical image:
<Call>
  <Purpose>Detect abnormalities in the Head
    and Neck X-ray image</Purpose>
  <Tool>TOOL8</Tool>
  <Input>['$Image$']</Input>
</Call>
Call Dict: {'Purpose': 'Detect abnormalities in
  the Head and Neck X-ray image', 'Tool':
  'TOOL8', 'Input': ['$Image$']}
Memory bank: {'$Image$': 'PLACEHOLDER_IMAGE',
  '$Information$': 'PLACEHOLDER_
  INFORMATION', '$Anatomy$': 'Head and Neck',
  '$Modality$': 'X-ray', '$Disease$':
  'Cervical spine degenerative changes'}
Score bank: {'$Image$': 1.0, '$Information$':
  1.0, '$Anatomy$': 0.95, '$Modality$':
  0.95, '$Disease$': 0.75}
Fixed bank: frozendict.frozendict({'$Image$':
  1.0, '$Information$': 1.0})
Step 3 completes!
```

2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008

During the fourth step, when selecting the anomaly detection tool, the agent core correctly identifies that the available tools in the set lack the capability to detect and localize anomalies in Head and Neck X-ray images. As demonstrated below, the agent core appropriately rejects this invalid operation and successfully recognizes the tools' capability gap.

Listing 15: Tool Selection & Execution Step 4

```
Step 4 starts!
Stepoutput: Before we can quantify biomarkers
  (per the original tool chain), we need
  the anomaly detection results first. Looking at
  the available tools for Head and Neck
  X-ray anomaly detection:
<NoCall>
  <Purpose>Detect specific anomalies in Head
    and Neck X-ray for biomarker
    quantification</Purpose>
  <Category>Anomaly Detector</Category>
  <Anatomy>Head and Neck</Anatomy>
  <Modality>X-ray</Modality>
  <Ability>SpecificToolMissing</Ability>
</NoCall>
Deny Call Dict: {'Purpose': 'Detect specific
  anomalies in Head and Neck X-ray for
  biomarker quantification', 'Category': 'Anomaly
  Detector', 'Anatomy': 'Head and Neck',
  'Modality': 'X-ray', 'Ability':
  'SpecificToolMissing'}Agent denies the call
  because no
  suitable tool is found.
```

2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034

Overall, while the agent core demonstrates correct initialization and tool chain planning, it exhibits a deviation during the tool selection & execution stage. Despite having access to tool chain guidance in each round, the agent core fails to follow these directives, resulting in an unnecessary disease diagnosis tool call. Nevertheless, the execution proceeds without any IO errors, and the agent core successfully identifies the tools' capability limitations, appropriately concluding with a denial of task execution due to the missing functionality.

2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047