

CP-Env: Evaluating Large Language Models on Clinical Pathways in a Controllable Hospital Environment

Anonymous ACL submission

Abstract

Medical care follows complex clinical pathways that extend beyond isolated physician-patient encounters, emphasizing decision-making and transitions between different stages. Current benchmarks focusing on static exams or isolated dialogues inadequately evaluate large language models (LLMs) in dynamic clinical scenarios. We introduce CP-Env, a controllable agentic hospital environment designed to evaluate LLMs across end-to-end clinical pathways. CP-Env simulates a hospital ecosystem with patient and physician agents, constructing scenarios ranging from triage and specialist consultation to diagnostic testing and multidisciplinary team meetings for agent interaction. Following real hospital adaptive flow of healthcare, it enables branching, long-horizon task execution. We propose a three-tiered evaluation framework encompassing Clinical Efficacy, Process Competency, and Professional Ethics. Results reveal that most models struggle with pathway complexity, exhibiting hallucinations and losing critical diagnostic details. Interestingly, excessive reasoning steps can sometimes prove counterproductive, while top models tend to exhibit reduced tool dependency through internalized knowledge. CP-Env advances medical AI agents development through comprehensive end-to-end clinical evaluation. We provide the benchmark and evaluation tools for further research and development.

1 Introduction

Delivering effective and compassionate medical care extends far beyond isolated physician-patient encounters. Instead, it constitutes a complex clinical pathway involving repeated interactions among health-care providers and patients, ultimately forming a coherent service continuum. This pathway may include triage guidance (“Which department should I visit?”), specialty consultations (“What is causing my symptoms?”) and diagnostic workups (“What tests do I need?”), multidisciplinary team



Figure 1: CP-Env introduces a comprehensive agentic hospital environment designed to address the full spectrum of patients’ healthcare needs through an integrated evaluation framework. In contrast to existing benchmarks that rely on static examinations or isolated dialogue scenarios, CP-Env provides dynamic, interactive environments with sophisticated tool-use capabilities. This approach enables both the delivery of comprehensive patient care and the rigorous evaluation of LLM-based agents across their performance in the complex, pathway-based clinical workflows characteristic of real-world hospital settings.

(MDT) discussions, treatment planning, and prognosis counseling (“How should I recover?”). Crucially, the process emphasizes decision-making and adaptive transitions between steps rather than executing a predetermined linear sequence.

Recently, AI-based agents have begun to demonstrate their potential in complex real-world scenarios, where large language models (LLMs) execute long-horizon tasks in dynamic environments through sustained interactive engagement (Yue

044
045
046
047
048
049
050
051
052
053

054	et al., 2024; Jia et al., 2025; Lu et al., 2025; Li	Competency: Does the agent demonstrate sound	106
055	et al., 2025). This paradigm is already reshaping	and logically coherent problem-solving processes?	107
056	biomedical domains such as genetic experiment	(3) Professional Ethics: Does the agent maintain	108
057	design and clinical diagnosis (Qu et al., 2025; Jin	ethical compliance and deliver humanistic care in	109
058	et al., 2025; Qiu et al., 2025a), offering a new lens	patient interactions? This framework comprehen-	110
059	through which to rethink the role of AI in medicine.	sively assesses LLMs' capabilities in this complex	111
060	However, a fundamental question remains: how	healthcare environment.	112
061	can end users and AI developers determine which		113
062	systems perform best in health-care contexts?	We conducted a comprehensive evaluation of	114
063	Although a growing number of medical bench-	outstanding models using CP-Env. Our experimen-	115
064	marks have been introduced to evaluate LLM per-	tal results reveal a clear performance stratification:	116
065	formance, they are limited in scope. Existing	(1) Proprietary models demonstrate significant ad-	117
066	benchmarks either focus on medical knowledge	vantages in Clinical Efficacy, showing the capabil-	118
067	and reasoning in examination-style formats (Jin	ity to reliably complete complex, branching clinical	119
068	et al., 2019, 2021) or assess conversational abili-	pathways. (2) We identified that a primary	120
069	ties in patient-oriented dialogues (Schmidgall et al.,	failure point for other models is the emergence of	121
070	2024; Fan et al., 2024). These approaches are insuf-	hallucinations in extended pathways, where they	122
071	ficient for evaluating LLM-based agents because	occasionally become myopically focused on imme-	123
072	they: (1) lack dynamic environments and tool-use	diate situational analysis rather than maintaining	124
073	capabilities necessary for realistic, controlled com-	broader diagnostic workflow awareness. (3) Our	125
074	parisons; and (2) fail to capture real-world tasks	multidimensional framework also reveals nuanced	126
075	that reflect the intricate clinical pathways of actual	insights: GPT-5 exhibits exceptional comprehen-	127
076	clinical practice.	sive clinical efficacy, while certain models demon-	128
077	To address these limitations, we introduce CP-	strate specific strengths in professional ethics (such	129
078	Env, an open-ended environment in which agents	as Seed-OSS's performance in empathy).	130
079	play diverse healthcare roles and collaboratively	The main contributions of this paper are summa-	131
080	engage in clinically realistic pathways to deliver	rized as follows:	132
081	patient care. Specifically, CP-Env encompasses:		133
082	(1) Patient role simulation, where patients are	• We introduce CP-Env, the first controllable	134
083	provided with clinical presentations and assume	agentic environment for evaluating LLMs in	
084	seeking medical consultation. They interact with	dynamic, end-to-end clinical pathways.	
085	each attending physician and accurately report		135
086	their known conditions during inquiries. (2) Clin-	• We propose a multidimensional evaluation	136
087	ical pathway navigation: Under expert guidance,	framework comprising Clinical Efficacy, Process	137
088	we design four clinical scenarios following real-	Capability, and Professional Ethics to	138
089	world healthcare pathways, including registration,	comprehensively assess medical agents be-	139
090	specialist consultation, diagnostic testing, and ad-	yond mere diagnostic accuracy.	
091	vanced diagnosis and treatment. Physicians in each		140
092	scenario are assigned tasks that mirror real clinical	• We provide a comprehensive report on current	141
093	practice. Throughout this healthcare pathway, pa-	LLMs, uncovering their characteristics in re-	142
094	tient behavior is dynamic and responsive to physi-	alistic medical scenarios, such as failures in	143
095	cian requests—for instance, when a specialist or-	pathway navigation and non-linear tool depen-	144
096	ders laboratory tests, the patient proceeds to the	ency patterns.	
097	laboratory department. Information between path-		145
098	way nodes is managed through electronic medical	2 Related Works	
099	records, and physicians can integrate clinical deci-	Agentic Environment. Agentic Environments re-	146
100	sion support tools to enhance their diagnostic and	fer to LLM-driven dynamic simulations that repli-	147
101	treatment decisions.	cate real-world scenarios. Generative Agents (Park	148
102	Subsequently, we establish a progressive three-	et al., 2023) pioneers this field by developing cog-	149
103	tier evaluation framework tailored to this environ-	gnitive architectures with memory, reflection, and	150
104	ment: (1) Clinical Efficacy: Can the agent suc-	planning capabilities, enabling multiple agents to	151
105	cessfully resolve medical problems? (2) Process	exhibit believable social behaviors within a vir-	152
		tual town. Subsequent research evolves along two	153

Benchmark	Task Type	Interaction Paradigm	Clinical Scope	Agent Roles	Tool Usage	Evaluation Metrics
MedQA (Jin et al., 2021)	Medical QA	Static	✗	✗	✗	Accuracy
Medbullets (Chen et al., 2025)	Medical QA	Static	✗	✗	✗	Accuracy
CMB-Clin (Wang et al., 2024)	Diagnosis	Static	Consultation	✗	✗	Clinical Efficacy
AI Hospital (Fan et al., 2024)	Diagnosis	Dialogue	Consultation & Exams	Doctor-Patient	✗	Clinical Efficacy
AgentClinic (Schmidgall et al., 2024)	Diagnosis	Dialogue	Consultation & Exams	Doctor-Patient	✓	Accuracy & Compliance
MedChain (Liu et al., 2024)	CDM	Dialogue	Sequential Healthcare	Doctor-Patient	✓	Accuracy & IoU
MedAgentSim (Almansoori et al., 2025)	Diagnosis	Dialogue & Self-Evolving	Consultation & Exams	Doctor-Patient	✓	Accuracy
CP-Env	Full-Process Healthcare	Dynamic Agentic Environment	Full Clinical Pathway	Multi-Agent	✓	Holistic (CE, PC, PE)

Table 1: **Comparisons with existing medical benchmarks.** We categorize existing benchmarks into two types based on the evolution of interaction modalities: **Static Exam-based QA** and **Sequential Interactive Dialogue**. CDM denotes clinical decision making. Unlike previous benchmarks that rely on static exam questions or single-scenario dialogues, *CP-Env* presents a pathway-based, dynamic environment. This is achieved through dynamic physician interactions, full-process healthcare delivery spanning interconnected pathway stages, and multi-dimensional evaluation metrics encompassing Clinical Efficacy, Process Competency, and Professional Ethics.

primary trajectories. The first focuses on large-scale social dynamics, scaling agent populations and grounding behaviors in real-world data to study macroscopic phenomena (Piao et al., 2025; Park et al., 2024; Mou et al., 2024). The second trajectory emphasizes sophisticated professional environments, transforming simulations from behavioral observation platforms into agent evaluation and optimization systems (Almansoori et al., 2025; Zhang et al., 2025). Agent Hospital (Li et al., 2024) enables physician agents to self-evolve through doctor-patient interactions and knowledge repository integration, while AgentsCourt (He et al., 2024) simulates judicial processes to evaluate and improve verdict prediction accuracy. However, current medical LLMs provide only chat functionality, failing to address patients’ comprehensive healthcare needs. Our objective is to implement, evaluate, and optimize a hospital agent environment that guides patients through complete, end-to-end clinical pathways.

Medical Benchmark. Early medical benchmarks (Jin et al., 2019, 2021), derived from academic papers and licensing examinations, primarily employ multiple-choice questions to assess models’ medical knowledge (Pal et al., 2022; Liu et al., 2023). The advent of LLMs drives significant evolution in benchmark design. First, researchers move away from traditional examination questions toward authentic clinical cases that better align with real-world scenarios (Chen et al., 2024; Wang et al., 2023; zhao zy15, 2024). Second, as reasoning models evolve (OpenAI, 2024; DeepSeek-AI, 2025), enhanced reasoning requirements emerge. Researchers begin exploring scenarios and questions that demand stronger analytical

capabilities (Zuo et al., 2025; Qiu et al., 2025b; Wu et al., 2025; Zhu et al., 2025a). Third, some studies move beyond static question-answering toward authentic doctor-patient conversations (Zhu et al., 2025b; Schmidgall et al., 2024; Fan et al., 2024). However, existing benchmarks either focus on medical knowledge and reasoning through examination-style formats or assess conversational abilities in patient-oriented dialogues. Yet, delivering effective and compassionate medical care extends far beyond isolated physician-patient encounters. *CP-Env* simulates the multi-party clinical pathway and comprehensively evaluates multiple dimensions of care delivery throughout the patient’s journey.

3 Interactive Hospital Environment

To reflect real-world clinical pathways, we introduce *CP-Env*, an interactive environment based on real-world cases that integrates comprehensive information for evaluating the capabilities of LLMs in agentic hospital settings. This section elaborates the essential characteristics, including patient role simulation (Section 3.1), clinical pathway navigation (Section 3.2), and healthcare delivery mechanism (Section 3.3).

3.1 Patient Role Simulation

The effectiveness of the agentic hospital relies fundamentally on realistic patient simulation. We anchor our simulations in authentic clinical cases, with each patient role derived from comprehensive medical records. This rich, reliable medical data ensures accurate patient representation and authentic doctor-patient interactions, thereby maintaining clinical validity. We source data from top-tier medical journals containing detailed clinical encounter

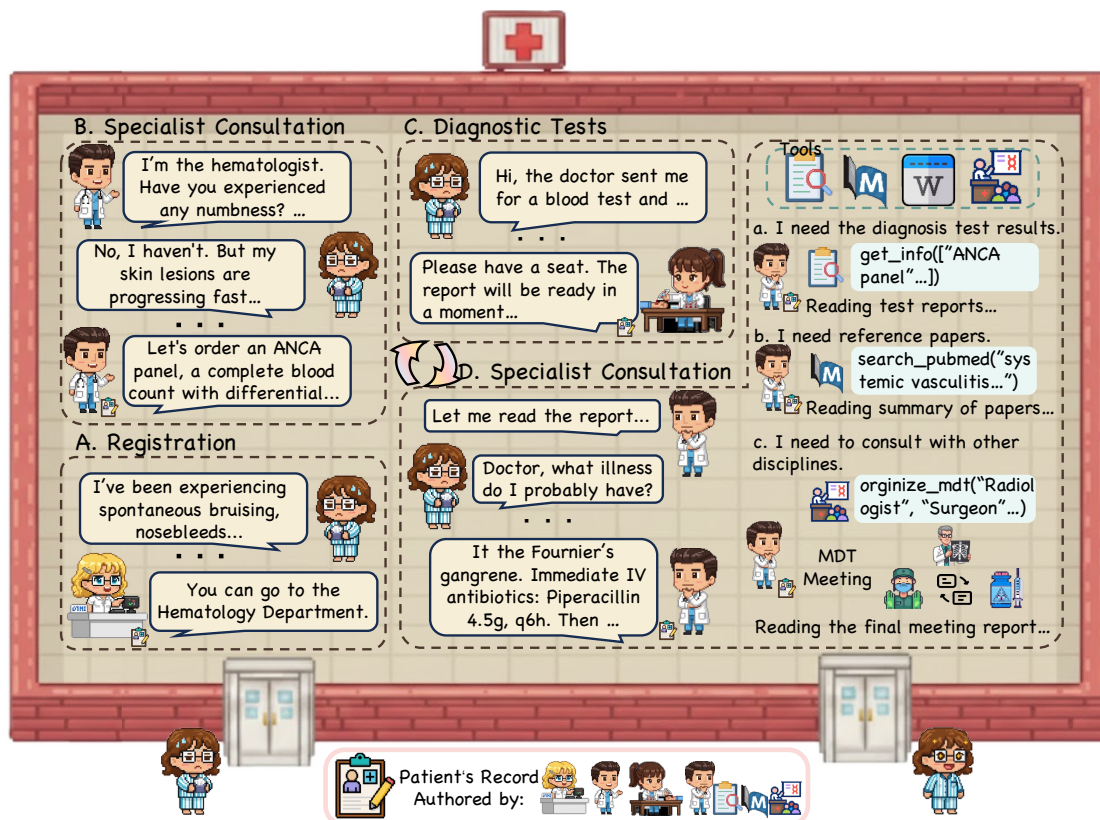


Figure 2: **Illustration of the Interactive Hospital Environment Clinical Pathway.** CP-Env integrates multiple physician roles, each executing specialized tasks through patient interactions. Upon hospital admission, patients are guided through the adaptive care pathway by different healthcare physicians, with medical records being progressively documented at each decision node. Physicians also utilize tools to collect multi-source information for clinical decision-making, ultimately facilitating patient recovery.

information (Zhu et al., 2025a).

Patients are configured to present to the hospital outpatients with specific physical complaints. During the clinical pathway, they engage in dialogue with physicians and are instructed to faithfully replicate the information they know. Each patient possesses knowledge of their general physical condition, including primary symptoms, medical history, and observable physical characteristics—all extracted from case records. To reflect real-world constraints, patients have limited access to comprehensive data, particularly laboratory and diagnostic test results, which is available only after physicians order the examinations.

3.2 Clinical Pathway Navigation

To comprehensively evaluate LLM capabilities in the simulated hospital environment, we design a multi-stage simulation scenario that mirrors real-world clinical pathways. As illustrated in Figure 2, the simulation encompasses the patient's journey through several decision nodes within the clinical pathway (A-D). In each stage, the evaluated LLMs

assume specific physician roles and perform defined clinical tasks through interactions with both the patient agent and the hospital environment. Importantly, the pathway is adaptive and branching, allowing for dynamic interactions and iterative reasoning that reflect the complexity of actual clinical decision-making.

Stage A: Registration and Triage. This initial stage involves the interaction between the patient and the triage nurse. Upon arrival at the outpatient department, the patient presents their symptoms (e.g., spontaneous bruising). The triage nurse needs to conduct a preliminary assessment through empathetic dialogue, evaluate the patient's general condition and symptom severity, and document all pertinent information in the medical record. Subsequently, a primary task is to recommend the appropriate specialist department for the patient's next stage of care.

Stage B: Specialist Consultation. Following triage, the patient enters the specific department and interacts with the designated specialist. In this stage, the specialist needs to conduct a com-

269 prehensive anamnesis through multi-turn dialogue, 320
270 exploring the patient’s medical history and specific 321
271 signs (e.g., numbness or lesion progression) while 322
272 documenting findings in the medical record. Un- 323
273 like initial triage, this encounter requires deeper 324
274 domain expertise to differentiate nuanced presen- 325
275 tations. The primary task here is hypothesis gen- 326
276 eration and examination ordering. Based on the 327
277 dialogue, the specialist must formulate initial di- 328
278 agnostic hypotheses and order appropriate inves- 329
279 tigations (e.g., an ANCA panel) to verify these 330
280 suspicions, facilitating progression to subsequent 331
281 diagnostic procedures. 332

282 **Stage C: Diagnostic Testing.** Following the spe- 333
283 cialist’s test orders, the hospital environment gen- 334
284 erates corresponding laboratory and imaging results. 335
285 To access these results from the medical records, 336
286 the specialist needs to employ information retrieval 337
287 tools. At this stage, a critical task is result inter- 338
288 pretation and synthesis. Specifically, the LLM must 339
289 parse raw medical results, identify abnormal indi- 340
290 cators, and integrate these objective findings with 341
291 the subjective information collected during Stage 342
292 B. This scenario evaluates the LLM’s capacity to 343
293 ground its reasoning in multimodal clinical data, 344
294 rather than depending solely on conversation. 345

295 **Stage D: Advanced Diagnosis and Treatment.** 346
296 This final stage yields the definitive clinical out- 347
297 come after the comprehensive multi-stage assess- 348
298 ment. With test results, the specialist needs to vali- 349
299 date previous hypotheses and refine the differential 350
300 diagnosis. When evidence supports high diagnostic 351
301 certainty, they can determine the definitive diagno- 352
302 sis. If evidence remains inconclusive, additional 353
303 investigations may be warranted, returning the pa- 354
304 tient to Stage B. For complex cases, the LLM can 355
305 facilitate a MDT meeting, enabling collaboration 356
306 with experts from complementary disciplines (e.g., 357
307 radiology, surgery) to deliberate on diagnostic and 358
308 treatment strategies. The specialist can also re- 359
309 trieve pertinent literature from external databases 360
310 like PubMed to strengthen the evidence base for 361
311 decision-making. This stage culminates in a com- 362
312 prehensive Final Clinical Report containing the 363
313 confirmed diagnosis, evidence-based medication 364
314 regimen, and structured follow-up protocol. 365

315 3.3 Healthcare Delivery Mechanism 366

316 Hospital patient care encompasses complex, clin- 367
317 ical pathways with adaptive branching spanning 368
318 multiple stages. Patients frequently undergo mul- 369
319 tiple examinations, require iterative consultations,

and attend unscheduled follow-up visits, creating 320
dynamic care trajectories. To authentically simu- 321
late real-world healthcare processes, CP-Env imple- 322
ments medical record management to ensure seam- 323
less transitions across scenarios, complemented by 324
comprehensive tool support. 325

Medical Record Management. Dynamic health- 326
care delivery generates complex, nonlinear patient 327
data throughout the care continuum. In real-world 328
clinical settings, comprehensive record manage- 329
ment protocols are rigorously implemented, includ- 330
ing mandatory documentation for every clinical 331
encounter. CP-Env adopts this medical record man- 332
agement paradigm by requiring physician agents 333
to document clinical reports after each patient in- 334
teraction, with all reports stored in the patient’s 335
medical record. During subsequent follow-up vis- 336
its or appointments, incoming physician agents can 337
directly assess the patient’s medical history and 338
current status through previous clinical reports. 339

Multidisciplinary Team Collaboration. Complex 340
clinical diagnoses often require multidisciplinary 341
collaboration. When cases exceed single-specialty 342
capabilities, multidisciplinary team (MDT) meet- 343
ings convene physicians from various disciplines to 344
provide diverse perspectives. We replicate this ap- 345
proach by enabling attending physicians to assem- 346
ble MDT teams with specialized expertise during 347
the diagnostic process. Through iterative discus- 348
sions, these teams generate analyses stored in the 349
medical record, informing attending physicians’ 350
decision-making across key clinical domains. 351

Clinical Tool Orchestration. Clinical diagnostic 352
decision-making requires physicians to synthesize 353
heterogeneous data from multiple sources. CP-Env 354
incorporates a suite of clinical tools that mirror real- 355
world workflows. The diagnostic process begins 356
with patient-physician dialogues for symptom elic- 357
itation through verbal communication. Physicians 358
then access laboratory results and retrieve perti- 359
nent reports from medical records via information 360
tools. To support evidence-based practice, CP-Env 361
integrates real-time queries to medical knowledge 362
bases, including PubMed and Wikipedia. Further- 363
more, it facilitates MDT consultations, enabling 364
physicians to leverage cross-departmental exper- 365
tise through discussions and MDT reports. 366

367 4 Agent Evaluation Framework 368

Leveraging the Interactive Hospital Environment 368
established in the previous chapter, we system- 369

370	atically collected comprehensive interaction data	420
371	from LLMs throughout the complete healthcare	421
372	workflow. To conduct a rigorous and multifaceted	422
373	evaluation of LLM capabilities within agentic hos-	423
374	pital settings, we developed an Agent Evaluation	424
375	Framework guided by three progressive research	425
376	questions: (1) Clinical Efficacy: Can the agent suc-	426
377	cessfully resolve medical problems? (2) Process	427
378	Competency: Does the agent demonstrate sound	428
379	and logically coherent problem-solving processes?	429
380	(3) Professional Ethics: Does the agent maintain	430
381	ethical compliance and deliver humanistic care in	431
382	patient interactions?	432
383	4.1 Clinical Efficacy	433
384	Clinical efficacy in real-world settings constitutes	434
385	the fundamental benchmark for healthcare evalu-	435
386	ation. Accordingly, LLM agents must prioritize	436
387	optimizing patient outcomes through accurate di-	437
388	agnosis and therapeutic interventions.	438
389	Work Completion (WC) evaluates whether LLMs	439
390	can comprehensively fulfill the whole hospital	440
391	workflow.	441
392	Diagnosis Recall@k (DR@k) evaluates whether	442
393	the top k diagnosis contain the correct diagnosis.	443
394	Triage Precision (TP) measures the appropriateness	444
395	of recommended medical departments.	445
396	4.2 Process Competency	446
397	A competent physician not only provides accurate	447
398	diagnoses but also demonstrates rigorous clinical	448
399	reasoning and effective utilization of diagnostic	449
400	tools. CP-Env evaluates LLMs across information	450
401	inquiry and gathering, clinical reasoning and di-	451
402	agnostic logic, and medical record documentation,	452
403	which comprehensively examines LLMs' ability to	453
404	synthesize complex medical information and uti-	454
405	lize clinical tools, providing a holistic evaluation	455
406	of their medical competency.	456
407	Inquiry Sufficiency (IS) measures the extent of	457
408	essential diagnostic information obtained through	458
409	clinical inquiry.	459
410	Logic Coherence (LC) quantifies the complete-	460
411	ness and consistency of diagnostic reasoning chains	461
412	throughout the healthcare continuum.	462
413	Record Compliance (RC) evaluates the quality	463
414	and completeness of clinical documentation.	464
415	Investigation Coverage (IC) quantifies the IoU	465
416	between physician-ordered tests and the ground-	466
417	truth case's diagnostic tests.	467
418	Result Utilization (RU) measures the proportion	468
419	of ordered test results utilized by the physician.	
	4.3 Professional Ethics	
	Practicing patient-centered care requires physicians	
	to extend their role beyond accurate diagnosis to ad-	
	dress patients' psychological vulnerability with em-	
	pathy, and appropriate professional boundaries. To	
	evaluate these abilities in LLMs, CP-Env conducts	
	comprehensive assessments of patient encounter	
	dialogues.	
	Privacy Safeguard (PS) assesses LLMs' capacity	
	to safeguard patient privacy during diagnosis.	
	Treatment Individualization (TI) quantifies how	
	treatment plans address patient-specific factors.	
	Empathic Dialogue (ED) evaluates the demonstra-	
	tion of care and compassion toward patients.	
	Follow-up Planning (FP) evaluates the quality and	
	appropriateness of follow-up planning.	
	5 Experiments	
	5.1 Settings	
	Agent Models. To comprehensively evaluate	
	LLMs' capabilities in the hospital environment,	
	we selected multiple models to serve as physi-	
	cian agent backbones. Given that CP-Env re-	
	quires physicians to dynamically leverage ex-	
	ternal tools, function-calling capability consti-	
	tutes the primary selection criterion. Our evalu-	
	ation encompasses both open-source and propri-	
	etary models. The open-source models included	
	Seed-OSS-36B-Instruct (Seed-OSS; Team, 2025a),	
	Qwen3-30B-A3B-Instruct-2507 (Qwen3; Team,	
	2025b), Qwen3-Next-80B-A3B-Instruct (Qwen3-	
	Next), GLM-4.5-Air (Z.ai, 2025), Llama-3.3-70B-	
	Instruct (Llama-3.3; Meta, 2025a), Llama-4-Scout-	
	17B-16E-Instruct (Llama-4; Meta, 2025b), GPT-	
	OSS-120B (OpenAI, 2025a) and GLM-4.7 (Z.ai,	
	2025). Additionally, we incorporated state-of-the-	
	art proprietary models, specifically Gemini-2.5-	
	Pro, Gemini-3-Pro (Deepmind, 2025), GPT-5 (Ope-	
	nAI, 2025b), GPT-5.2. Unfortunately, existing	
	open-source medical models generally lack reli-	
	able function-calling capabilities, and Qwen3 re-	
	asoning models cannot handle long workflow tasks	
	effectively; therefore, they were excluded from our	
	evaluation. Furthermore, to establish a unified and	
	equitable testing platform, we selected GPT-OSS-	
	120B as the patient agent backbone due to its cost-	
	effectiveness, accessibility, and wide recognition.	
	5.2 Main Results	
	Table 2 presents the main evaluation results of	
	LLMs on CP-Env. The results reveal a clear perfor-	

Models	Clinical Efficacy				Process Competency					Professional Ethics			
	WC	DR@3	DR@5	TP	IS	LC	RC	IC	RU	PS	TI	ED	FP
<i>Open Source LLMs</i>													
Seed-OSS _{36B}	80.11	17.38	17.92	67.76	51.20	64.74	77.60	24.29	46.99	98.14	31.08	50.16	38.31
Qwen3 _{30B}	85.03	39.67	40.55	66.89	44.50	8.81	78.73	27.26	70.89	99.13	17.29	24.75	41.42
GPT-OSS _{120B}	76.28	29.62	31.15	72.79	72.37	11.98	86.01	9.96	18.56	98.80	22.37	49.75	<u>47.65</u>
Qwen3-Next _{80B}	83.83	21.20	22.30	70.05	59.19	<u>54.89</u>	75.75	21.49	66.71	85.90	12.33	45.76	36.12
GLM-4.5-Air	84.70	35.41	36.39	71.58	65.56	42.90	69.42	21.86	<u>93.27</u>	95.41	19.51	<u>57.44</u>	39.34
Llama-3.3 _{70B}	90.71	39.67	40.98	69.84	64.17	19.22	55.14	<u>25.44</u>	93.31	93.88	13.63	49.17	35.85
Llama-4 _{Scout}	74.86	39.67	41.42	68.63	60.85	32.00	56.20	24.24	93.17	92.73	14.37	48.85	17.27
GLM-4.7	85.25	37.92	39.45	73.66	55.36	28.08	77.13	24.89	91.56	91.58	20.08	49.43	35.85
<i>Proprietary LLMs</i>													
Gemini-2.5-Pro	94.86	39.45	40.66	73.22	62.86	46.56	77.88	5.30	12.36	99.34	21.56	64.21	39.84
Gemini-3-Pro	<u>93.99</u>	39.89	40.89	70.16	<u>76.60</u>	37.06	75.48	21.59	83.18	<u>99.73</u>	27.04	54.85	36.45
GPT-5	93.33	<u>44.81</u>	<u>47.43</u>	<u>75.41</u>	85.28	20.69	94.48	3.13	6.96	<u>99.73</u>	<u>30.73</u>	50.71	51.37
GPT-5.2	91.37	58.14	59.67	77.60	67.98	8.75	<u>92.34</u>	1.15	2.18	99.78	29.20	49.96	35.50

Table 2: **Performance of Different Large Language Models on the CP-Env Benchmark.** The evaluation encompasses clinical efficacy, process competency, and professional ethics. Bold scores indicate the best performance, while underlined scores represent the second-best.

mance hierarchy: proprietary models substantially outperform their open-source counterparts. GPT-5.2 and GPT-5 achieve superior performance across key dimensions, including clinical pathway navigation, reasoning consistency, and patient empathy. Gemini-3-Pro and Gemini-2.5-Pro exhibit comparable performance, while open-source models show more limited capabilities.

Performance at Clinical Efficacy. At CE, only proprietary models like GPT-5 perform well on WC, successfully navigating the complete clinical pathways, while other models presented limitations. A primary failure observed is cognitive hallucination during extended workflows—models become entrapped in reasoning loops and fail to advance logically through sequential steps, which is particularly pronounced in Llama-4 and GPT-OSS-120B. Some reasoning models (e.g., Qwen3-30B-A3B-Thinking) exhibit it so severely that they are excluded from evaluation. Additionally, certain models (e.g., GLM-4.5-Air) suffer from intermittent tool-calling format errors. In the DR, GPT-5.2 and GPT-5 achieve superior performance through its extensive knowledge base and well-calibrated reasoning capabilities. Llama and Qwen3-30B ranked second, with Qwen3-30B demonstrating remarkably strong performance despite its smaller size. Notably, Qwen3-Next underperformed relative to expectations. Comparative analysis revealed that while Qwen3-Next provided detailed

responses with comprehensive reasoning chains, occasional hallucinations led to significant deviations from clinical pathways. In contrast, Qwen3-30B’s more concise and direct reasoning approach proved sufficient for achieving favorable outcomes, while GPT-5 demonstrated appropriate reasoning restraint without excessive elaboration. These findings parallel observations in real clinical settings, where excessive analytical complexity does not necessarily correlate with improved patient outcomes. This suggests that our benchmark rewards judicious and targeted reasoning over exhaustive but potentially error-prone analysis.

Performance at Process Competency. Overall, OpenAI models demonstrate superior medical competency with substantial advantages over competitors. In IS, GPT-5 and GPT-OSS-120B extract the most comprehensive patient information through inquiries, highlighting their proficiency in clinical dialogue, while other models clustered within a narrower performance range behind. For RC, OpenAI models maintain leadership by producing standardized medical documentation with comprehensive content and precise terminology, reflecting superior medical literacy compared to other models. However, Seed-OSS and Qwen3-Next excel in logical reasoning tasks (LC). Seed-OSS’s advantage stems from its thinking model method, while Qwen3-Next exhibited robust reasoning capabilities consistent with prior analysis. Notably,

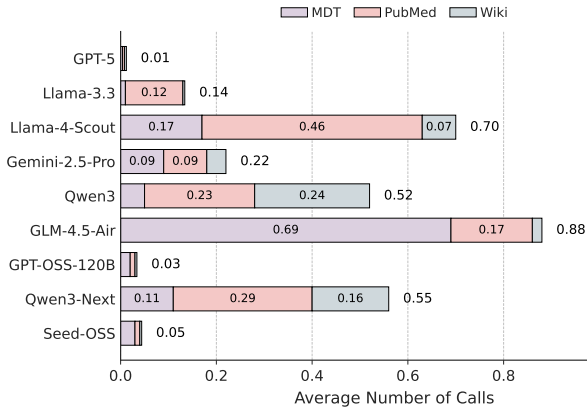


Figure 3: Comparison of average tool usage frequency and types across different models.

hallucinations during extended tasks degrade overall performance, suggesting that extensive thinking processes do not necessarily improve outcomes in long-sequence tasks. It should be noted that GPT-5’s lower LC performance may not reflect its true capabilities, as the model’s reasoning processes are concealed during inference. In IC and RU, GPT-5 and Gemini-2.5-Pro adopted conservative approaches, ordering excessive diagnostic tests compared to real cases, resulting in lower efficiency scores. While other models performed more reasonably, all LLMs demonstrated over-testing tendencies. RU reveals that GLM-4.5-air and LLaMA models showed more judicious diagnostic test utilization, whereas GPT exhibited redundant test-ordering patterns.

Performance at Professional Ethics. All models perform well in privacy protection PS. However, Seed-OSS excels in addressing personalized user needs TI and achieves outstanding ED performance, alongside Gemini-2.5-Pro and GPT-5, which lead in user interactions. GLM-4.5-Air performs commendably in ED as well. Regarding follow-up management, OpenAI’s GPT-5 and GPT-OSS-120B significantly outperform other models. In summary, OpenAI models exhibit robust medical competency, while Seed-OSS excels in user-centered interactions.

5.3 Tool Utilization

To comprehensively analyze tool utilization strategies across different models and their relationship with diagnostic performance, we calculated the average frequency of tool invocations (MDT, PubMed, and Wiki) for each model during medical consultation dialogues, as presented in Figure 3. Models are ordered by DR scores.

Our analysis reveals a U-shaped, rather than

simply positive, relationship between tool usage and diagnostic accuracy. Specifically, (1) High-performing models with minimal tool usage: GPT-5, which achieved the highest diagnostic accuracy, invoked tools negligibly during experiments. Similarly, Gemini-2.5-Pro demonstrated highly selective tool utilization, suggesting that advanced models have internalized sufficient medical knowledge or developed reasoning capabilities that reduce dependence on external resources. (2) Mid-tier models with intensive tool utilization: Models with intermediate accuracy, such as Llama-4-Scout and GLM-4.5-Air, exhibited the highest tool utilization rates. GLM-4.5-Air recorded an average invocation frequency of 0.88, favoring the MDT tool (0.69), while Llama-4-Scout (0.70 total) preferred PubMed retrieval (0.46). This pattern suggests that these models compensate for knowledge uncertainties by actively leveraging tools for validation and supplementation. (3) Low-performing models with limited tool engagement: Models with lower accuracy (e.g., GPT-OSS-120B and Seed-OSS-36B) demonstrated minimal tool invocation rates (0.03 and 0.04, respectively), likely reflecting limitations in intent recognition and instruction-following capabilities. These models often fail to determine when tool assistance is necessary or how to effectively utilize available tools.

These findings have important implications for medical agent design. For medium-scale models, optimizing tool utilization pipelines represents a critical pathway to improved performance. Conversely, for state-of-the-art models, enhancing internal reasoning mechanisms may yield greater benefits than integrating additional external knowledge.

6 Conclusion

This paper introduces CP-Env, a controllable multi-agent hospital environment for evaluating LLMs in end-to-end clinical pathways. Unlike existing benchmarks limited to static tasks, CP-Env assesses models through realistic clinical simulations using the framework: Clinical Efficacy, Process Competency, and Professional Ethics. Our experiments reveal that proprietary models significantly outperform open-source alternatives in navigating complex clinical pathways, with hallucinations during extended processes being the primary failure mode. CP-Env offers essential benchmarks for advancing medical AI beyond isolated evaluations toward comprehensive clinical pathway assessment.

616 Limitations

617 Despite the rigorous design of CP-Env and its care-
618 ful consideration of clinical reproducibility, sev-
619 eral limitations warrant acknowledgment. First,
620 although we systematically simulated real-world
621 hospital clinical pathways to establish a controlled
622 experimental environment, actual clinical scenarios
623 encompass factors beyond purely medical consid-
624 erations, such as physician-patient relationships,
625 physician emotional states, and fatigue. Second,
626 while CP-Env provides a systematic analysis of
627 performance variations across 24 major medical
628 departments, clinical medicine encompasses a vast
629 spectrum of rare conditions. The current frame-
630 work prioritizes simulating clinical pathways for
631 prevalent diseases alongside selected diagnostically
632 challenging cases. Future work will focus on ex-
633 panding coverage to incorporate a broader range
634 of long-tail diseases, enabling more comprehen-
635 sive evaluation of medical agent robustness under
636 atypical clinical scenarios.

637 Ethical Considerations

638 This study develops CP-Env, a controlled hospital
639 environment simulating clinical pathways. All clin-
640 ical case data used in CP-Env were sourced from
641 publicly available medical journals and records,
642 with rigorous de-identification procedures ensuring
643 compliance with data protection standards. We ac-
644 knowledge that LLMs serving as physician agents
645 may reflect practice biases from their training data,
646 and observed hallucinations indicate potential de-
647 viations from sound clinical reasoning. CP-Env
648 is positioned as an agentic evaluation framework
649 for academic research—assessing medical AI per-
650 formance across clinical effectiveness, procedural
651 competence, and professional ethics—rather than a
652 deployable clinical diagnostic system. Any practi-
653 cal application must undergo rigorous clinical vali-
654 dation with human oversight and ultimate decision-
655 making authority retained by qualified healthcare
656 professionals.

657 References

658 Mohammad Almansoori, Komal Kumar, and Hisham
659 Cholakkal. 2025. Self-evolving multi-agent simula-
660 tions for realistic clinical interactions. *arXiv preprint*
661 *arXiv:2503.22678*.

662 Hanjie Chen, Zhouxiang Fang, Yash Singla, and Mark
663 Dredze. 2024. Benchmarking large language models

on answering and explaining challenging medical
questions. *arXiv preprint arXiv:2402.18060*.

Hanjie Chen, Zhouxiang Fang, Yash Singla, and Mark
Dredze. 2025. Benchmarking large language mod-
els on answering and explaining challenging medical
questions. In *Proceedings of the 2025 Conference*
of the Nations of the Americas Chapter of the Asso-
ciation for Computational Linguistics: Human Lan-
guage Technologies (Volume 1: Long Papers), pages
3563–3599.

Google Deepmind. 2025. [Gemini 2.5 pro](#).

DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing rea-](#)
[soning capability in llms via reinforcement learning.](#)
Preprint, arXiv:2501.12948.

Zhihao Fan, Jialong Tang, Wei Chen, Siyuan Wang,
Zhongyu Wei, Jun Xi, Fei Huang, and Jingren Zhou.
2024. Ai hospital: Benchmarking large language
models in a multi-agent medical interaction simulator.
arXiv preprint arXiv:2402.09742.

Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin,
Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang,
Kang Liu, and Jun Zhao. 2024. Agentscourt: Build-
ing judicial decision-making agents with court debate
simulation and legal knowledge augmentation. *arXiv*
preprint arXiv:2403.02959.

Zheng Jia, Shengbin Yue, Wei Chen, Siyuan Wang, Yi-
dong Liu, Yun Song, and Zhongyu Wei. 2025. Ready
jurist one: Benchmarking language agents for legal
intelligence in dynamic environments. *arXiv preprint*
arXiv:2507.04037.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,
Hanyi Fang, and Peter Szolovits. 2021. What disease
does this patient have? a large-scale open domain
question answering dataset from medical exams. *Ap-*
plied Sciences, 11(14):6421.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W
Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset
for biomedical research question answering. *arXiv*
preprint arXiv:1909.06146.

Ruofan Jin, Yucheng Guo, Yuanhao Qu, Ming Yang,
Chun Shang, Qirong Yang, Linlin Chao, Yi Zhou,
Ruilai Xu, Ziyao Xu, and 1 others. 2025. Biolab:
End-to-end autonomous life sciences research with
multi-agents system integrating biological foundation
models. *bioRxiv*, pages 2025–09.

Junkai Li, Yunghwei Lai, Weitao Li, Jingyi Ren, Meng
Zhang, Xinhui Kang, Siyu Wang, Peng Li, Ya-Qin
Zhang, Weizhi Ma, and 1 others. 2024. Agent hospi-
tal: A simulacrum of hospital with evolvable medical
agents. *arXiv preprint arXiv:2405.02957*.

Peiji Li, Jiasheng Ye, Yongkang Chen, Yichuan Ma,
Zijie Yu, Kedi Chen, Ganqu Cui, Haozhan Li, Ji-
acheng Chen, Chengqi Lyu, and 1 others. 2025. In-
ternbootcamp technical report: Boosting llm rea-
soning with verifiable task scaling. *arXiv preprint*
arXiv:2508.08636.

720	Jie Liu, Wenxuan Wang, Zizhan Ma, Guolin Huang, Yihang SU, Kao-Jung Chang, Wenting Chen, Haoliang Li, Linlin Shen, and Michael Lyu. 2024. Medchain: Bridging the gap between llm agents and clinical practice through interactive sequential benchmarking. <i>arXiv preprint arXiv:2412.01605</i> .	773	Pengcheng Qiu, Chaoyi Wu, Junwei Liu, Qiaoyu Zheng, Yusheng Liao, Haowen Wang, Yun Yue, Qianrui Fan, Shuai Zhen, Jian Wang, and 1 others. 2025a. Evolving diagnostic agents in a virtual clinical environment. <i>arXiv preprint arXiv:2510.24654</i> .	774
721		775		776
722		776		777
723		777		
724		778		
725		779		
726	Junling Liu, Peilin Zhou, Yining Hua, Dading Chong, Zhongyu Tian, Andrew Liu, Helin Wang, Chenyu You, Zhenhua Guo, Lei Zhu, and 1 others. 2023. Benchmarking large language models on cmexam-a comprehensive chinese medical exam dataset. <i>Advances in Neural Information Processing Systems</i> , 36:52430–52452.	780	Pengcheng Qiu, Chaoyi Wu, Shuyu Liu, Weike Zhao, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2025b. Quantifying the reasoning abilities of llms on real-world clinical cases. <i>arXiv preprint arXiv:2503.04691</i> .	781
727		781		782
728		782		
729		783		
730		784		
731		785		
732		786		
733	Siyuan Lu, Zechuan Wang, Hongxuan Zhang, Qintong Wu, Leilei Gan, Chenyi Zhuang, Jinjie Gu, and Tao Lin. 2025. Don't just fine-tune the agent, tune the environment. <i>arXiv preprint arXiv:2510.10197</i> .	787	Yuanhao Qu, Kaixuan Huang, Ming Yin, Kanghong Zhan, Dyllan Liu, Di Yin, Henry C Cousins, William A Johnson, Xiaotong Wang, Mihir Shah, and 1 others. 2025. Crispr-gpt for agentic automation of gene-editing experiments. <i>Nature Biomedical Engineering</i> , pages 1–14.	788
734		788		
735		789		
736		790		
737	Meta. 2025a. Llama 3.3. https://www.llama.com/docs/model-cards-and-prompt-formats/llama3_3/ .	791	Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor. 2024. Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments. <i>arXiv preprint arXiv:2405.07960</i> .	792
738		792		793
739		793		
740	Meta. 2025b. Llama 4. https://www.llama.com/docs/model-cards-and-prompt-formats/llama4/ .	794	ByteDance Seed Team. 2025a. Seed-oss open-source models. https://github.com/ByteDance-Seed/seed-oss .	795
741		795		796
742		796		
743	Xinyi Mou, Zhongyu Wei, and Xuanjing Huang. 2024. Unveiling the truth and facilitating change: Towards agent-based large-scale social movement simulation. <i>arXiv preprint arXiv:2402.16333</i> .	797	Qwen Team. 2025b. <i>Qwen3 technical report</i> . <i>Preprint</i> , arXiv:2505.09388.	798
744		798		
745		799		
746		800		
747	OpenAI. 2024. <i>Learning to reason with llms</i> .	801	Xidong Wang, Guiming Chen, Song Dingjie, Zhang Zhiyi, Zhihong Chen, Qingying Xiao, Junying Chen, Feng Jiang, Jianquan Li, Xiang Wan, and 1 others. 2024. Cmb: A comprehensive medical benchmark in chinese. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6184–6205.	802
748	OpenAI. 2025a. <i>gpt-oss-120b & gpt-oss-20b model card</i> . <i>Preprint</i> , arXiv:2508.10925.	803		804
749		804		805
750	OpenAI. 2025b. <i>Introducing gpt-5</i> .	805		806
751	Ankit Pal, Logesh Kumar Umapathi, and Malaikanan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In <i>Conference on health, inference, and learning</i> , pages 248–260. PMLR.	806	Xidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao, Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, and 1 others. 2023. Cmb: A comprehensive medical benchmark in chinese. <i>arXiv preprint arXiv:2308.08833</i> .	807
752		807		808
753		808		809
754		809		810
755		810		811
756	Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th annual acm symposium on user interface software and technology</i> , pages 1–22.	811	Kevin Wu, Eric Wu, Rahul Thapa, Kevin Wei, Angela Zhang, Arvind Suresh, Jacqueline J Tao, Min Woo Sun, Alejandro Lozano, and James Zou. 2025. Medcasereasoning: Evaluating and learning diagnostic reasoning from clinical case reports. <i>arXiv preprint arXiv:2505.11733</i> .	812
757		812		813
758		813		814
759		814		815
760		815		816
761		816		817
762	Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein. 2024. Generative agent simulations of 1,000 people. <i>arXiv preprint arXiv:2411.10109</i> .	817		818
763		818		
764		819		
765		820		
766		821		
767	Jinghua Piao, Yuwei Yan, Jun Zhang, Nian Li, Junbo Yan, Xiaochong Lan, Zhihong Lu, Zhiheng Zheng, Jing Yi Wang, Di Zhou, and 1 others. 2025. Agentsociety: Large-scale simulation of llm-driven generative agents advances understanding of human behaviors and society. <i>arXiv preprint arXiv:2502.08691</i> .	822	Murong Yue, Wenhan Lyu, Wijdane Mifdal, Jennifer Suh, Yixuan Zhang, and Ziyu Yao. 2024. Mathvc: An llm-simulated multi-character virtual classroom for mathematics education. <i>arXiv preprint arXiv:2404.06711</i> .	823
768		823		
769		824		
770		825		
771				
772				

826 Kaiyuan Zhang, Jiaqi Li, Yueyue Wu, Haitao Li,
827 Cheng Luo, Shaokun Zou, Yujia Zhou, Weihang Su,
828 Qingyao Ai, and Yiqun Liu. 2025. Chinese court sim-
829 ulation with llm-based agent system. *arXiv preprint*
830 *arXiv:2508.17322*.

831 zhao zy15. 2024. Rarearena. Website. [https://](https://github.com/zhao-zy15/RareArena)
832 github.com/zhao-zy15/RareArena.

833 Yakun Zhu, Zhongzhen Huang, Linjie Mu, Yutong
834 Huang, Wei Nie, Jiaji Liu, Shaoting Zhang, Pengfei
835 Liu, and Xiaofan Zhang. 2025a. Diagnosisarena:
836 Benchmarking diagnostic reasoning for large lan-
837 guage models. *arXiv preprint arXiv:2505.14107*.

838 Yinghao Zhu, Ziyi He, Haoran Hu, Xiaochen Zheng,
839 Xichen Zhang, Zixiang Wang, Junyi Gao, Liantao
840 Ma, and Lequan Yu. 2025b. Medagentboard: Bench-
841 marking multi-agent collaboration with conventional
842 methods for diverse medical tasks. *arXiv preprint*
843 *arXiv:2505.12371*.

844 Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai
845 Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and
846 Bowen Zhou. 2025. Medxpertqa: Benchmarking
847 expert-level medical reasoning and understanding.
848 *arXiv preprint arXiv:2501.18362*.

A Further Analysis

A.1 Dialogue Efficiency

To evaluate computational efficiency and interaction patterns of different models in our agentic hospital environment, we analyze encounters per patient, dialogue turns per encounter, and token consumption across various LLMs. Table 3 presents comprehensive dialogue efficiency statistics.

The results reveal distinct consultation patterns across models despite modest overall differences. Qwen3-30B favors frequent consultations (3.54 encounters per patient), reflecting a cautious, step-by-step approach. In contrast, GLM-4.5-Air pursues fewer encounters but with deeper dialogue (3.51 turns per encounter), indicating comprehensive single-session exploration. This creates a complementary pattern: while Qwen3-30B shows the highest encounter frequency, it exhibits the lowest dialogue depth (2.22 turns per encounter), suggesting models maintain similar total information volumes per patient through different interaction strategies—either multiple focused consultations or fewer comprehensive sessions.

Token usage patterns strongly correlate with models’ reasoning characteristics. Reasoning models demonstrate substantially higher consumption, with Seed leading at the highest token usage, resulting in significantly elevated computational costs. Other reasoning models show more restrained consumption: GPT-5 uses 1,089.09 tokens per encounter while Gemini-2.5-Pro maintains conservative usage at 653.55 tokens, suggesting more efficient reasoning mechanisms or superior output control.

In our benchmark involving complex clinical pathways, extensive reasoning may introduce counterproductive effects. Models with advanced reasoning occasionally become myopically focused on immediate situational analysis rather than maintaining broader pathway awareness. This reasoning trap leads to inefficient resource allocation and potentially suboptimal clinical pathway decisions.

A.2 Departmental Characteristics Analysis

To investigate the performance characteristics of large language models across different medical departments, we conducted a Departmental Characteristics Analysis in this section. We systematically analyzed the Diagnosis Recall (DR) and Triage Precision (TP) of nine models across 24 departments. This analysis specifically encompasses

Model	Avg.E. Nums	Avg.E. Tokens	Avg.T. Nums	Avg.T. Tokens
Llama-3.3 _{70B}	3.06	454.36	3.11	145.86
Llama-4 _{Scout}	3.21	534.57	2.74	194.59
Qwen3 _{30B}	3.54	478.99	2.22	215.09
Qwen3-Next _{80B}	3.09	532.89	2.37	224.31
GLM-4.5-Air	3.02	774.35	3.51	220.33
GPT-OSS _{120B}	3.20	641.92	2.47	259.74
Gemini-2.5-Pro	3.14	653.55	2.67	244.52
GPT-5	3.33	1089.09	2.35	461.54
Seed-OSS _{36B}	3.35	2843.19	2.66	1066.52

Table 3: Average encounters, turns, and token consumption across models. Avg.E. means average encounters per patient; Avg.T. means average dialogue turns per encounter.

Departmental Difficulty Stratification and Model Domain-Specific Analysis.

Departmental Difficulty Stratification We analyzed the performance of different models’ Diagnosis Recall (DR) and Triage Precision (TP) across various departments, with the results visualized in heatmaps of Figure 4. Our findings reveal distinct performance stratification patterns across different departments.

First, high-performing departments include Ophthalmology, Dermatology, and Stomatology & Maxillofacial Surgery, which demonstrate exceptional performance in both metrics. For department allocation tasks, all models achieve TP scores exceeding 0.90, with heatmaps displaying deep red high-value regions. This superior performance stems from these departments’ cases typically containing highly distinctive anatomical features and visual descriptive characteristics, enabling models to achieve accurate triage through straightforward semantic pattern matching. Moreover, their DR performance in these departments also stands out, generally surpassing other departments, suggesting that the clinical features for disease diagnosis in these specialties may possess strong discriminative properties. However, despite the remarkably high TP, the DR shows a notable decline, with dermatology’s mean DR primarily ranging between 0.40 and 0.50. This indicates that for LLMs, localizing lesions based on explicit features is considerably easier than confirming specific pathological types, and fine-grained disease differentiation remains a primary bottleneck for current models.

Second, moderately performing departments include traditional internal medicine specialties such

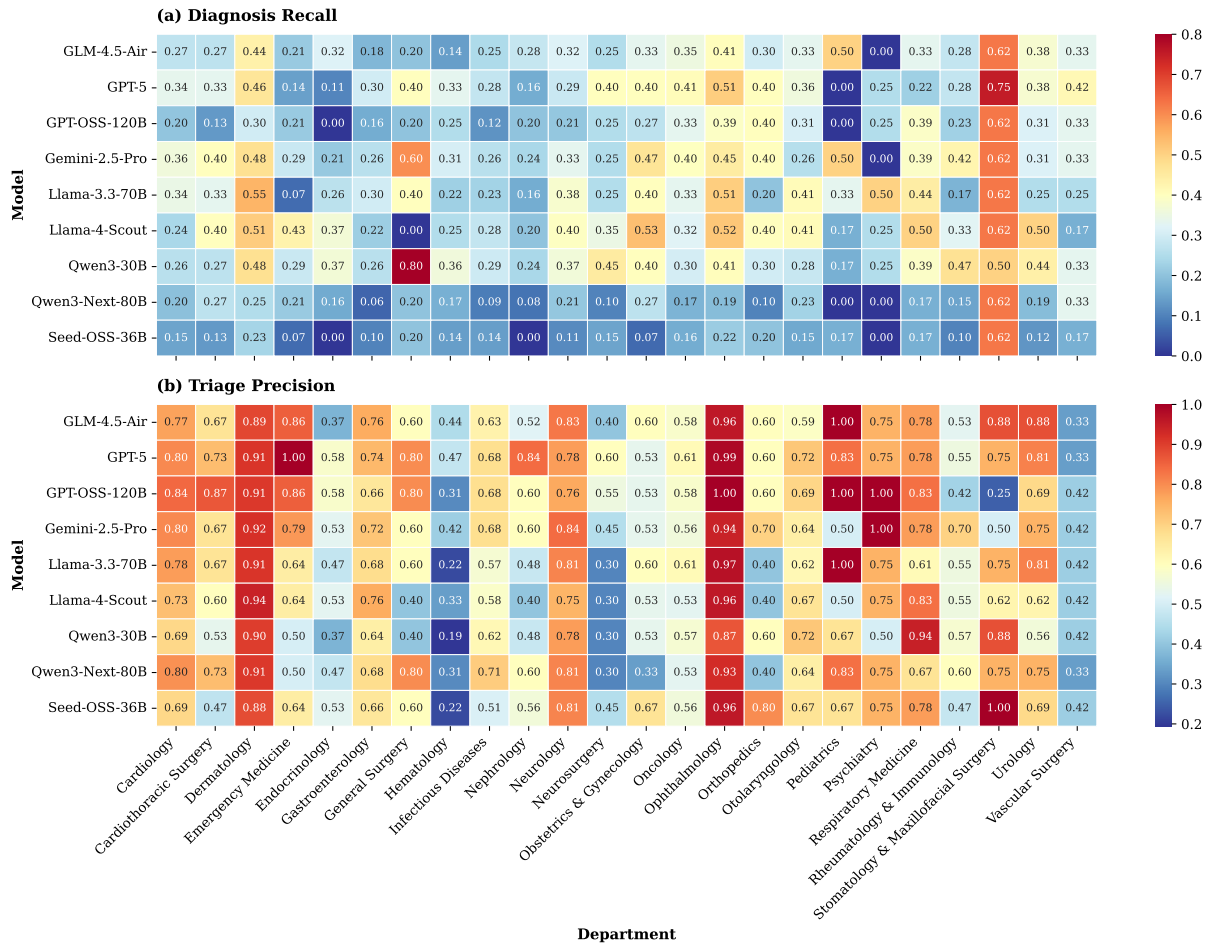


Figure 4: **Fine-grained performance analysis of different LLMs across medical departments.** The heatmaps illustrate the performance of various models across clinical specialties for Diagnosis Recall and Triage Precision. The color gradient represents the accuracy score, ranging from dark blue (lower performance) to deep red (higher performance).

934 as Cardiology, Gastroenterology, and Respiratory
 935 Medicine. These departments’ heatmaps predom-
 936 inantly display yellow-orange hues, with DR dis-
 937 tributed in the 0.30-0.50 range and TP fluctuating
 938 between 0.60-0.80. This tier is characterized by
 939 high symptom overlap and substantial clinical man-
 940 ifestation intersection. For instance, chest pain
 941 could indicate myocardial infarction or originate
 942 from respiratory or gastrointestinal disorders. This
 943 clinical ambiguity demands enhanced differential
 944 diagnostic capabilities from models, requiring them
 945 to move beyond simple keyword matching and in-
 946 tegrate patient history, accompanying symptoms,
 947 and exclusionary logic for multi-step hypothesis-
 948 deductive reasoning.

949 Finally, low-performing departments include
 950 Hematology, Nephrology, and Infectious Diseases.
 951 Heatmaps reveal extensive deep blue low-score
 952 regions, with most models achieving DR below

0.20 and significantly reduced TP. The underly-
 953 ing cause of this phenomenon lies in these depart-
 954 ments’ heavy reliance on laboratory tests and quan-
 955 titative indicators. Unlike departments with explicit
 956 features such as Dermatology or Ophthalmology,
 957 Hematology and Nephrology require specific bio-
 958 chemical test data, while Infectious Diseases neces-
 959 sitates detailed epidemiological history screening.
 960 During consultation dialogues, patients typically
 961 can only describe non-specific symptoms like fever
 962 and fatigue, unable to directly provide crucial diag-
 963 nostic information. This requires models to possess
 964 proactive clinical information acquisition capabil-
 965 ities—identifying potential disease spectra under-
 966 lying non-specific symptoms and guiding patients
 967 to complete necessary examinations. The system-
 968 atic low scores in these departments indicate that
 969 current LLMs have not yet fully mastered proac-
 970 tive, goal-oriented clinical consultation strategies,
 971

972 representing a crucial direction for future model
973 optimization.

974 In summary, the current capabilities of LLMs
975 in medical consultation tasks demonstrate excel-
976 lent performance in specialties with distinctive fea-
977 tures, moderate performance in internal medicine
978 departments with symptom overlap, and poor per-
979 formance in departments heavily dependent on ob-
980 jective examinations.

981 **Model Domain-Specific Analysis** Beyond the
982 inherent difficulty of stratification across depart-
983 ments, our cross-model comparison within individ-
984 ual departments reveals domain-specific specializa-
985 tion differences. Notably, several models exhibit
986 performance spikes in particular departments that
987 exceed their average capabilities, revealing poten-
988 tial domain bias in general-purpose LLMs within
989 medical subspecialties.

990 In general surgery, the Qwen3-30B model
991 achieved a diagnosis recall of 0.80, significantly
992 outperforming other outstanding models. This ex-
993 ceptional performance in conversational diagnosis
994 suggests that the model may have acquired highly
995 efficient consultation strategies for handling acute
996 abdominal surgical or trauma cases. Given that
997 general surgery diagnosis heavily relies on precise
998 identification of specific physical signs (such as
999 rebound tenderness and muscle guarding), Qwen3-
1000 30B’s superior performance likely stems from its
1001 ability to rapidly identify critical surgical indica-
1002 tors during conversations without being distracted
1003 by irrelevant internal medicine symptoms. This
1004 single-domain breakthrough phenomenon demon-
1005 strates that small parameter models, through tar-
1006 geted knowledge and capability enhancement in
1007 specific domains, are fully capable of surpassing
1008 general-purpose large models in specialized medi-
1009 cal consultation tasks.

1010 Additionally, Gemini-2.5-Pro demonstrated ex-
1011 ceptional reasoning capabilities in gastroenterol-
1012 ogy, achieving a DR score of 0.60—significantly
1013 outperforming other models’ average of 0.30. Un-
1014 like surgical specialties that rely on specific physi-
1015 cal signs, gastroenterological diagnosis typically re-
1016 quires comprehensive evaluation of dietary history,
1017 pain patterns, and long-term medication use. This
1018 necessitates maintaining logical coherence through-
1019 out extended dialogues. In simulated consultations,
1020 Gemini-2.5-Pro effectively managed conversation
1021 flow, elucidating complex gastrointestinal symp-
1022 tom evolution through iterative questioning rather
1023 than being confounded by overlapping symptoms

1024 in patient presentations. This underscores that in
1025 internal medicine specialties requiring complex
1026 history-taking and temporal reasoning, a model’s
1027 contextual understanding and memory mechanisms
1028 are critical determinants of diagnostic accuracy.

1029 GPT-5 demonstrates the characteristics of a gen-
1030 eral practitioner, and its heatmap displays the most
1031 uniform color distribution with minimal deep-blue
1032 zero-score regions. This cross-specialty consis-
1033 tency validates its capabilities. In real-world pri-
1034 mary care settings where patients present with un-
1035 differentiated symptoms, an effective diagnostic
1036 AI must not only excel in its areas of strength but
1037 also avoid catastrophic errors in less familiar do-
1038 mains. GPT-5’s stable performance across all spe-
1039 cialties indicates it has developed a more gener-
1040 alizable medical diagnostic framework, enabling
1041 flexible transfer across different clinical contexts
1042 while maintaining consistently high baseline diag-
1043 nostic standards.

1044 B Case Example

1045 The following is an example of a patient’s medical
1046 record.

Example of Patient’s Medical Record

Recording Department: Registration Clinic Note

Solitary exophytic scalp nodule (~4 cm) on the vertex, present since birth with mild enlargement after minor trauma two months ago. No urgent symptoms reported. Recommend evaluation of a congenital scalp lesion and discussion of imaging and possible excision.

Suggestion: Dermatology

Recording Department: Specialist Consul- tation

Clinic Note

Chief complaint: Long-standing midline scalp nodule since birth with mild enlargement after minor trauma 2 months ago.

History of present illness: 72-year-old with a ~4 cm firm, exophytic scalp nodule centered at the vertex (midline). Surface waxy and normochromic. Surrounding plaque feels infiltrated. Mildly tender to pressure; no spontaneous pain. No discharge, bleeding, crusting, or pus. Not compressible or pulsatile; no change with Valsalva. Lesion feels fixed to the underlying skull, not mobile. No central

punctum or hair tuft. No headaches, nausea/vomiting, visual changes, seizures, fevers, or recurrent infections. No prior imaging or procedures.

Exam: Single ~4 cm firm exophytic nodule on the midline vertex with smooth/waxy skin-colored surface; surrounding plaque indurated. Lesion fixed to underlying calvarium. No visible sinus opening or hair tuft. No overlying erythema, warmth, or ulceration. No regional lymphadenopathy reported.

Assessment: Congenital midline scalp mass with bony fixation. Most likely diagnosis is a congenital dermoid cyst (or less likely epidermoid) with possible adherence to periosteum/skull and potential intracranial extension via a dermal sinus tract. Enlargement after trauma is plausible. Important to exclude connection to the superior sagittal sinus/dura before any biopsy or excision.

Differentials: (1) Nevus sebaceus with secondary adnexal neoplasm (e.g., trichoblastoma, syringocystadenoma papilliferum, less commonly BCC) presenting as a new nodule on a congenital plaque; (2) Proliferating trichilemmal tumor arising on the scalp; (3) Pilomatricoma; (4) Osteoma; (5) Encephalocele/sinus pericranii or vascular malformation are unlikely given non-pulsatility and lack of Valsalva change.

Plan: Do not perform biopsy before imaging. Obtain cross-sectional imaging to define any intracranial extension and calvarial involvement: (a) MRI brain and scalp with and without contrast, including fat-suppressed sequences; (b) MR venography to evaluate relation to the superior sagittal sinus; (c) Non-contrast high-resolution CT of the skull to assess bony defects and calcifications.

Optional: targeted high-resolution ultrasound with Doppler to assess lesion composition and vascularity if immediately available, but not a substitute for MRI/CT. Multidisciplinary coordination with Neurosurgery (and Plastics) for surgical planning. Definitive management, if no intracranial extension, is complete surgical excision down to (or including) involved periosteum with histopathology. If there is intracranial extension, plan a combined neurosurgical approach.

Counseling: avoid manipulation/trauma; return urgently for new headache, fever, drainage, rapid enlargement, redness, or neurologic symptoms.

Suggestion

Leading diagnosis: congenital midline scalp dermoid cyst with possible cranial involvement. Differentials include epidermoid cyst; nevus sebaceus with secondary adnexal tumor; proliferating trichilemmal tumor; pilomatricoma; osteoma; encephalocele/sinus pericranii (less likely).

Recommended tests: MRI brain and scalp with and without contrast (include fat-suppressed sequences), MR venography, and non-contrast high-resolution CT of the skull. Consider targeted high-resolution ultrasound with Doppler if readily available. Avoid biopsy until imaging excludes intracranial connection. Based on imaging, proceed to complete excision with appropriate surgical team and send for histopathology.

Recording Department: Diagnostic Test Completed Tests

['MRI brain and scalp with and without contrast (fat-suppressed sequences)', 'MR venography', 'Non-contrast high-resolution CT of the skull', 'Targeted high-resolution ultrasound with Doppler']

Recording Department: Specialist Consultation

Clinic Note

Dermatology Consultation – Current Visit

Subjective:

- 72-year-old with a congenital midline scalp lesion present since birth. Mild enlargement after minor trauma 2 months ago. Mild tenderness to pressure only. No drainage, bleeding, redness, warmth, or rapid growth. No headaches, fevers, visual changes, weakness, seizures, or other neurologic symptoms. No prior procedures or radiation to the scalp. Family history notable for father with metastatic melanoma (not suggestive of relation to this lesion).

Objective (exam summary):

- Single 4 cm firm, exophytic nodule centered at the vertex (midline) with smooth/waxy skin-colored surface. Surrounding plaque feels indurated. Lesion is fixed to the underlying skull;

not pulsatile; no change with Valsalva. No central punctum or hair tuft. No overlying erythema, warmth, or ulceration. No regional lymphadenopathy.

Imaging (reviewed today):

- MRI brain and scalp with and without contrast (fat-suppressed sequences): 4.2 × 3.8 × 2.5 cm well-circumscribed subcutaneous mass at the vertex, T2 hyperintense, iso–mildly hypointense on T1, non-enhancing; no diffusion restriction. No calvarial involvement or intracranial extension.

- MR venography: Normal patent dural venous sinuses; no involvement by the scalp lesion.

- Non-contrast high-resolution CT skull: Well-defined subcutaneous mass without calcifications; intact calvarium with no erosion or bony defect; no intracranial abnormality.

- Targeted high-resolution ultrasound with Doppler: Well-defined, heterogeneously echogenic subcutaneous nodule abutting periosteum without bony invasion; mild peripheral vascularity; no aggressive flow pattern.

Assessment:

- Congenital midline scalp mass most consistent with cutaneous meningeal heterotopia (heterotopic meningotheial hamartoma; Type I cutaneous meningioma). Imaging shows a well-circumscribed subcutaneous lesion without bone involvement or intracranial extension, and normal MRV, which aligns with CMH and argues against dermoid (no fat signal), epidermoid (no diffusion restriction), proliferating trichilemmal tumor (no calcifications/keratinous cyst), pilomatricoma (typically calcified), or nevus sebaceus with secondary neoplasm. No features concerning for vascular malformation or encephalocele.

Plan:

-Surgical management: Recommend complete excision of the lesion en bloc down to (and including, if adherent) the periosteum, with submission for histopathologic confirmation. Avoid preoperative/incisional biopsy due to midline location and proximity to the superior sagittal sinus, although MRV shows no involvement.

- Multidisciplinary coordination: Refer to Plastic Surgery for closure/reconstruction planning and to Neurosurgery for standby/consul-

tation given the midline vertex location and periosteal adherence, despite imaging showing no intracranial connection.

- Perioperative considerations: Standard pre-op assessment and medication review; manage anticoagulants/antiplatelets per surgical protocol if applicable. Discuss risks (bleeding, infection, nerve injury, scarring, need for periosteal/outer table burring if required) and low recurrence risk with complete excision.

- Pathology: Request immunohistochemistry (e.g., EMA, vimentin; PR/SSTR2A) to support meningotheial differentiation and exclude adnexal neoplasms or cystic lesions.

- Follow-up: Routine postoperative wound care and pathology review. No further imaging is needed unless new symptoms or interval change occur.

- Return precautions: Seek urgent care for new or worsening headache, fever, drainage, redness, rapid enlargement, or neurologic symptoms.

Final Diagnosis:

- Congenital midline scalp lesion most consistent with cutaneous meningeal heterotopia (heterotopic meningotheial hamartoma), without calvarial involvement or intracranial extension on MRI/MRV/CT/ultrasound.

Suggestion

Proceed with definitive treatment.

Diagnosis: Congenital midline scalp lesion most consistent with cutaneous meningeal heterotopia, with no intracranial extension or bony involvement on MRI/MRV/CT and ultrasound.

Plan: Coordinate surgical excision (Dermatology/Plastic Surgery with Neurosurgery consult), excise en bloc down to periosteum as needed, and send for histopathology with immunohistochemistry to confirm. Routine pre-op assessment and postoperative follow-up. No additional preoperative imaging is required at this time.

Treatment Planing

Final Diagnosis ['Cutaneous meningeal heterotopia (heterotopic meningotheial hamartoma; primary cutaneous meningioma, Type I)', 'Congenital dermoid cyst', 'Epidermoid cyst', 'Proliferating trichilemmal tumor', 'Nevus sebaceus with secondary adnexal neo-

plasm (e.g., trichoblastoma or syringocystadenoma papilliferum)']

Treatment Plan

Definitive management is complete surgical excision. Coordinate Dermatology/Plastic Surgery with Neurosurgery consultation given the midline vertex location. Perform en bloc excision down to the periosteum (including adherent periosteum and outer-table burring if necessary) with careful hemostasis to avoid superior sagittal sinus injury. Submit the specimen for histopathology and immunohistochemistry (e.g., EMA, vimentin ± PR/SSTR2A) to confirm meningotheial differentiation and exclude adnexal or cystic lesions. Standard preoperative assessment and perioperative medication management (hold anticoagulants/antiplatelets as appropriate). Routine postoperative wound care and follow-up; no additional imaging unless new symptoms or interval change occur. Provide return precautions for headache, fever, drainage, redness, rapid enlargement, or neurologic symptoms.

The following is an example of the doctor-patient conversation.

Example

Encounter Turn 1 - Registration

Doctor: Hi, I'm here to help direct you to the right department. What is your main symptom and where is it located? How long has this been going on? Do you have any urgent signs like chest pain, trouble breathing, severe weakness, high fever, or heavy bleeding?

Patient: I have a solitary lump on the top of my scalp. It's an exophytic nodule about 4cm in size. The lesion has been present since I was born, but I noticed that it grew a little after a minor bump to the area a couple of months ago. I don't have chest pain, shortness of breath, severe weakness, fever, heavy bleeding, or any other urgent symptoms.

Doctor: Please go to the Dermatology Department for further examination.

[Makes a note in the patient's record]

Encounter Turn 2 - Specialist Consultation

Doctor: How old are you, and can you describe the lump's features: its color, whether it's tender or painful, any discharge or bleed-

ing, whether it feels soft or firm, and if it seems compressible or pulsatile? Does its size change when you cough, strain, or cry?

Is the lump exactly in the midline of the scalp at the vertex? Does it move over the bone or feel fixed to the skull? Is there a central opening/punctum or a tuft of coarse hair over it?

Have you had headaches, nausea/vomiting, vision changes, seizures, fevers, or recurrent infections? Have you ever had any imaging (ultrasound/CT/MRI) or procedures on this lump before?

Patient: I am in my early's (about 72 years old). The lesion is a firm, exophytic nodule on the scalp with a waxy, normochromic surface; the surrounding plaque feels infiltrated.

It is mildly tender to pressure but not painful. There is no discharge, bleeding, crusting, or pus. The nodule is not compressible, not pulsatile, and its size does not change when I cough, strain, or cry.

The lesion is centered on the vertex of the scalp, essentially mid-line. It feels adherent to the underlying bone and does not glide over the skull; it is fixed in place. There is no visible central opening, punctum, or a tuft of coarse hair over the lesion.

I have had no headaches, nausea, vomiting, visual changes, seizures, fevers, or recurrent scalp infections. I have never had any imaging studies (ultrasound, CT, MRI) or any procedures performed on this lump before this visit.

Doctor: Leading diagnosis: congenital midline scalp dermoid cyst with possible cranial involvement. Differentiels include epidermoid cyst; nevus sebaceus with secondary adnexal tumor; proliferating trichilemmal tumor; pilomatricoma; osteoma; encephalocele/sinus pericranii (less likely). Recommended tests: MRI brain and scalp with and without contrast (include fat-suppressed sequences), MR venography, and non-contrast high-resolution CT of the skull. Consider targeted high-resolution ultrasound with Doppler if readily available. Avoid biopsy until imaging excludes intracranial connection. Based on imaging, proceed to complete excision with appropriate surgical team and send for histopathology.

[Makes a note in the patient's record]

1052

1053

1054

1055

1056

Encounter Turn 3 - Specialist Consultation

Doctor: Thanks for coming in. I have a few quick questions, and I'd like to review your imaging reports to make sure we plan the safest next steps.

Since your scans were done, have you noticed any new symptoms such as headaches, fever, drainage from the lesion, rapid growth, redness, or any neurological issues (vision changes, weakness, seizures)?

Have you ever had any prior scalp surgeries, radiation, or similar lesions elsewhere? Any family history of similar congenital lumps?

Please provide the official reports for the completed tests: MRI brain and scalp with and without contrast (including fat-suppressed sequences); MR venography; Non-contrast high-resolution CT of the skull; Targeted high-resolution ultrasound with Doppler

Patient: I have not noticed any new problems. There are no headaches, fevers, or any drainage from the lesion. It has not grown rapidly, there is no redness or warmth around it, and I have had no visual changes, weakness, or seizures.

I have never had any surgery on the scalp, nor have I received radiation to the head. I do not have any other skin or subcutaneous lumps elsewhere on my body.

In my family the only notable skin history is that my father died of metastatic cutaneous melanoma at age; there are no congenital scalp or "birth-mark" type lesions reported in relatives.

At this point I have not had any of the studies you listed (MRI brain+scalp with and without contrast, MR venography, high-resolution non-contrast CT of the skull, or a targeted high-resolution ultrasound with Doppler). Therefore I do not have any official reports to give you yet. If you think imaging is indicated, I am happy to arrange for them.

[Using get_info tool]

Doctor: Proceed with definitive treatment. Diagnosis: Congenital midline scalp lesion most consistent with cutaneous meningeal heterotopia, with no intracranial extension or bony involvement on MRI/MRV/CT and ultrasound. Plan: Coordinate surgical excision (Dermatology/Plastic Surgery with Neu-

rosurgery consult), excise en bloc down to periosteum as needed, and send for histopathology with immunohistochemistry to confirm. Routine pre-op assessment and postoperative follow-up. No additional preoperative imaging is required at this time.

[Makes a note in the patient's record]

C Implementation Details

C.1 Experimental Setup

In our benchmark experiments, we deployed GPT-OSS-120B as the patient agent model using 2 H100 GPUs, and subsequently evaluated the performance of various doctor models. For the doctor models, we utilized official APIs for the proprietary models Gemini-2.5-Pro and GPT-5. For open-source models, we deployed them using vLLM across 2 to 8 H100 GPUs. Qwen3-Next was an exception due to unresolved issues with its vLLM implementation; therefore, we deployed it using SGLang with 8 H100 GPUs.

C.2 Roles Simulation

We deployed GPT-OSS-120B as the patient agent model using 2 H100 GPUs, with the model under evaluation serving as the physician model.

The LLM was instructed to simulate a patient through the specific task description and prompt. The simulated patient was provided only with medical history and physical examination findings, without access to actual laboratory results or final diagnoses, thereby reflecting the information asymmetry typical of real-world patient encounters.

You are a simulated patient, intended to test the hospital's medical procedures and the doctor's diagnostic skills. You are currently role-playing as a patient at a hospital, where you will interact with various individuals and engage in limited communication with them.

Below is the simulated case provided to you:
{INSERT_CASE_HERE}

Please remember the following:

1. When the doctor inquires about your medical condition, you should respond based on the provided simulated case.
2. You only need to answer the questions the doctor asks you. If a question is not asked, you do not need to provide any information.

The simulated patient is then navigated through different clinical scenarios. Since patient behavior

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104	may vary depending on the specific context or setting, we provide tailored prompts for each scenario to ensure appropriate and realistic responses.	You have the tool to get completed test reports mentioned in the medical record. If a test you need is not available, you should list the required tests in your final response.	1167 1168 1169 1170
1107	Registration:		1171
1108	You have just arrived at the hospital. Your first step is to provide basic personal information to register. A guidance nurse will speak with you to get a general understanding of your condition and recommend an appropriate department.	You may ask up to {INSERT_QUERYNUMS_HERE} questions per turn, with a maximum of {INSERT_TURNS_HERE} rounds of dialogue. You must deliver your final diagnosis before the dialogue ends.	1172 1173 1174 1175 1176
1111			1177
1112		Once you have reached a conclusion, respond in the following JSON format, enclosed by ```json and ```:	1178 1179
1113		```json	1180
1114	Now the guidance nurse says: {INSERT_QUERY_HERE}	{	1181
1115		"clinic_note": "A comprehensive clinic note for the patient's current visit. This should include your clinical assessment, the final diagnosis, and the proposed management or treatment plan.",	1182 1183 1184 1185 1186 1187
1116	Specialist Consultation:	"suggestion": "Your professional recommendation. If further tests are required, list them and set `next_step` to `diagnostic_test`. If the final diagnosis is confirmed, outline the diagnosis and treatment plan, and set `next_step` to `end_of_diagnosis`.",	1188 1189 1190 1191 1192 1193
1117	You have now arrived at the specialist consultation department of {INSERT_DEPARTMENT_HERE}.	"next_step": "Specify one of the following options: 'diagnostic_test' or 'end_of_diagnosis'."	1194 1195 1196 1197
1118		}	1198
1119		```	1199
1120			1200
1121	Now the physician asks: {INSERT_QUERY_HERE}	Treatment Planing:	1201
1122	Begin role-playing as the patient!	Based on your final analysis, enumerate the top 5 most likely diagnoses for this patient, ordered from most to least probable. In addition, provide the definitive treatment plan.	1202 1203 1204 1205 1206 1207
1123		Output in JSON format, enclosed by ```json and ```:	1208 1209
1124	Furthermore, different physicians assume distinct roles and responsibilities within the clinical workflow.	```json	1210
1125		{	1211
1126	Registration:	"final_diagnosis": ["Disease 1", "Disease 2", "Disease 3", "Disease 4", "Disease 5"],	1212 1213 1214
1127	You are a hospital guidance assistant stationed in the main lobby. Your job is to briefly assess each patient's general symptoms and recommend the appropriate department for consultation. Remember, you are just a guide, so keep the inquiry simple and focused on directing the patient efficiently.	"treatment_plan": "A treatment plan for the patient"	1215 1216
1128		}	1217
1129		```	1218
1130			
1131	You may ask up to {INSERT_QUERYNUMS_HERE} questions per turn, with a maximum of {INSERT_TURNS_HERE} rounds of dialogue. You must provide your final recommendation before the dialogue ends.	C.3 Evaluation Metrics	1219
1132		Work Completion. WC is calculated based on whether the final task is successfully completed, with a value of 1 for success and 0 for failure. Since the CP-Env involves multiple interactions and scenario transitions controlled by LLMs at each phase, WC assesses their capacity to handle complex task demands through long-sequence agent processes in real-world hospital environments.	1220 1221 1222 1223 1224 1225 1226
1133	Once you have reached a conclusion, respond in the following JSON format, enclosed by ```json and ```:	Diagnosis Recall. It evaluates the LLM's diagnostic capability by measuring alignment between its	1227 1228 1229
1134	```json		
1135	{		
1136	"clinic_note": "A guide note of the patient's reported symptoms and the reason for the referral.",		
1137	"suggestion": "The single, most appropriate department for the patient to visit. Must be one department name only.",		
1138	"next_step": "specialist_consultation"		
1139	}		
1140	```		
1141			
1142			
1143	Now the guidance nurse says: {INSERT_QUERY_HERE}		
1144			
1145			
1146			
1147			
1148			
1149			
1150			
1151			
1152			
1153			
1154			
1155			
1156			
1157	Specialist Consultation:		
1158	You are a specialist physician in the {INSERT_DEPARTMENT_HERE} department, responsible for conducting hospital consultations. Your task is to evaluate the patient's condition through dialogue and ultimately provide a diagnosis or recommend the necessary diagnostic tests.		
1159			
1160			
1161			
1162			
1163			
1164			
1165			
1166	Patient's Medical Record: {INSERT_RECORD_HERE}.		

1230	final diagnosis and the ground truth, where DR@k	Inquiry Sufficiency. IS measures the extent to	1294
1231	indicates whether the top k predictions contain the	which physicians successfully identify, elicit, and	1295
1232	correct diagnosis. The judge categorizes diagnosis	document essential diagnostic information during	1296
1233	pairs as identical, relevant, or irrelevant, with only	patient encounters. It evaluates clinicians' profi-	1297
1234	identical matches scoring 1 and all others scoring	ciency in capturing critical diagnostic elements nec-	1298
1235	0 (Zhu et al., 2025a). Additionally, LLMs pro-	essary for accurate case assessment. It is calculated	1299
1236	vide 5 diagnoses ranked by confidence level, where	by comparing core information points extracted	1300
1237	DR@k indicates whether the top k predictions con-	from the original case with those documented in	1301
1238	tain the correct diagnosis. The prompt is as follows:	the clinical records. The prompt for extracting core	1302
1239	You are an expert in diagnosing challenging	information is as follows:	1303
1240	cases. You will receive a student's answer		
1241	containing 5 differential diagnoses, as	You are an experienced clinical expert familiar	1304
1242	well as the reference diagnosis. You need	with medical diagnoses. Given a medical	1305
1243	to score each diagnosis from the student's	case: {case_info}, and its confirmed final	1306
1244	answer according to the following rules:	diagnosis: {right_diagnosis}, perform the	1307
1245	2 = The student's diagnosis exactly matches the	following tasks:	1308
1246	reference diagnosis;	1. List all the core and most important	1309
1247	1 = A broader or narrower disease category that	information a doctor must ask the patient	1310
1248	includes or is included in the reference	before making a correct diagnosis (e.g.,	1311
1249	diagnosis (e.g., 'pneumonia' vs 'bacterial	past medical history, family history).	1312
1250	pneumonia');	2. From the actual doctor communication record:	1313
1251	0 = Unrelated or incorrect.	{consult_communiaction}, identify which of	1314
1252		the core information you mentioned has	1315
1253	Here is the student's answer:	actually been asked by the doctor.	1316
1254	{','.join(diagnosis)}. Here is the	Only consider information that is a direct	1317
1255	reference diagnosis: {gt}.	match to your core list, only include items	1318
1256		that are present in your core information	1319
1257	Output Format: Output the scores in the	list.	1320
1258	following format. 1. Disease 1 Name: <The	Output format:	1321
1259	Score of Disease 1>; 2. Disease 2 name:	@ Core information needed: info 1, info 2, ...	1322
1260	<The Score of Disease 2>; ...	@ Inquired information in practice: info 1,	1323
1261		info 2, ...	1324
1262	Triage Precision. It measures the model's ability	- No additional text, explanation, or	1325
1263	to accurately route patients to appropriate medi-	punctuation is allowed. - Use consistent	1326
1264	cal departments based on initial assessment dia-	terminology to ensure exact matching.	1327
1265	logues. The judge validates the correctness of de-		
1266	partment assignments against case presentations.	Logic Coherence. LC quantifies the complete-	1328
1267	The prompt is as follows:	ness and consistency of diagnostic reasoning chains	1329
1268	You are an experienced clinical expert familiar	throughout the healthcare continuum, spanning	1330
1269	with medical diagnoses. Given the medical	from inquiry and examination to diagnosis and	1331
1270	case: {case_info} and the correct	treatment. It evaluates LLMs' ability to maintain	1332
1271	diagnosis: {right_diagnosis}, the doctor in	logical coherence across the entire diagnostic pro-	1333
1272	registration reception leads the patient to	cess. It is assessed by extracting and scoring the	1334
1273	the {department_suggestion} department of	cognitive processes demonstrated during conversa-	1335
1274	the hospital and gives a reason: {reason}.	tional patient encounters. The prompt is as follows:	1336
1275	Do you think this initial department		
1276	suggestion is correct?	You are an experienced clinical expert familiar	1337
1277		with medical diagnoses. Based on the doctor	1338
1278	Score the suggestion based on following rules:	consultation record:	1339
1279	2 = The suggested department fully matches the	{consult_communiacion} and the reference	1340
1280	correct diagnosis pathway and is clearly	correct diagnosis: {right_diagnosis}, you	1341
1281	the most appropriate destination (e.g.,	need to identify the Clinical Logic	1342
1282	chest pain -> Cardiology for myocardial	Coherence of doctor.	1343
1283	infarction).	Standard:	1344
1284	1 = The suggested department is somewhat	- 10 points: Complete and coherent clinical	1345
1285	reasonable but not optimal (e.g., dizziness	reasoning chain - 'inquiry -> tests ->	1346
1286	-> Neurology, but should ideally go to ENT	diagnosis -> treatment' - with consistent	1347
1287	based on final diagnosis).	causal logic (e.g., inquiry about dysphagia	1348
1288	0 = The suggested department is completely	-> esophageal test -> diagnosis of	1349
1289	inappropriate or inconsistent with the	esophageal injury -> mucosal protective	1350
1290	diagnosis (e.g., appendicitis ->	drugs).	1351
1291	Dermatology).	- 5 points: Partially coherent reasoning; some	1352
1292	Only response the score <0 or 1 or 2>, and no	steps logically related but missing one	1353
1293	other content is allowed.	link or minor mismatch.	1354
		- 0 points: Major logical break or	1355

1356	inconsistency (e.g., inquiry about	{','}.join(suggested_tests)), perform the	1418
1357	dysphagia -> no esophageal test ->	following steps:	1419
1358	unrelated diagnosis).	1. Deduplicate the tests and merge items that	1420
1359	Notes:	refer to the same test.	1421
1360	- Do not assume missing steps; rely only on the	2. Identify which suggested tests are	1422
1361	given record.	equivalent to the reference tests (same	1423
1362	- Planned but not yet performed steps count if	clinical purpose, even if named	1424
1363	explicitly reasoned.	differently).	1425
1364	- Apply strict and responsible judgment; avoid	Output format strictly:	1426
1365	0 or 10 unless strongly justified.	@ Reference tests: item1, item2, ...	1427
1366	Output in the following format:	@ Suggested tests in practice: item1, item2, ...	1428
1367	@Rating: <integer score 0-10>	@ Equivalent tests: item1, item2, ...	1429
1368	@Reason: <one-sentence explanation of the	- Use consistent terminology.	1430
1369	rating>	- Only output in the above format.	1431
		- Do not include extra text.	1432
1370	Record Compliance. RC assesses whether physi-	Result Utilization. The IC and RU assess both	1433
1371	cians accurately document core symptoms, adhere	the physician's diagnostic test prescription com-	1434
1372	to standardized documentation structures, and uti-	petency and the agent's proficiency in workflow	1435
1373	lize appropriate medical terminology. It evaluates	management and tool utilization.	1436
1374	the medical professionalism and clinical documen-		
1375	tation compliance of LLMs. We evaluate the clini-	You are an experienced clinical expert familiar	1437
1376	cal documentation quality at each stage by examin-	with medical diagnoses. Given the doctor's	1438
1377	ing whether physicians' clinical notes contain suffi-	suggested medical tests: {suggest_tests}	1439
1378	cient patient information, maintain proper structure,	and the retrieved tests from knowledge	1440
1379	and demonstrate appropriate use of professional	tools: {','}.join(re_tests)), identify which	1441
1380	medical terminology. The prompt is as follows:	suggested tests are equivalent to the	1442
		retrieved tests (same clinical purpose).	1443
1381	You are an experienced clinical expert familiar	Output format strictly:	1444
1382	with medical diagnoses. Evaluate the	@ Retrieved Suggested tests: item1, item2, ...	1445
1383	compliance of the doctor's clinical note	- Use consistent terminology.	1446
1384	({record}) based on the doctor-patient	- Only output in the above format.	1447
1385	communication record	- Do not include extra text.	1448
1386	({consult_communiaction}). You need to	Privacy Safeguard. Effective clinical practice re-	1449
1387	Evaluate whether the clinical note	quires physicians to balance comprehensive infor-	1450
1388	accurately reflects patient-reported	mation gathering with appropriate privacy bound-	1451
1389	symptoms and history, is well-organized,	aries, avoiding unnecessary inquiries into patients'	1452
1390	and uses precise medical terminology.	personal matters during clinical interviews. The	1453
1391	Scoring standard:	prompt is as follows:	1454
1392	- 10 points: Fully cites patient-reported	You are an experienced clinical expert familiar	1455
1393	information; structure is clear;	with medical diagnoses. Based on the doctor	1456
1394	terminology is precise and professional	consultation record:	1457
1395	(e.g., 'stridor' instead of 'throat noise').	{consult_communiaction}, evaluate the	1458
1396	- 5 points: Captures part of the patient's	degree to which the doctor adheres to	1459
1397	report but lacks completeness or uses some	Privacy Safeguard.	1460
1398	non-standard wording.	Scoring standard:	1461
1399	- 0 points: Misses key patient information,	- 10 points: Inquiries are strictly limited to	1462
1400	disorganized, or vague (e.g., 'Patient	condition-related information, with no	1463
1401	unwell, order tests').	irrelevant personal questions (e.g.,	1464
1402	Apply strict and responsible judgment when	income, relationship status).	1465
1403	scoring, be careful when scoring 0 or 10.	- 5 points: Generally appropriate but includes	1466
1404	Output in the following format:	one or two mildly irrelevant or borderline	1467
1405	@Rating: <integer score 0-10>	personal questions without persistent	1468
1406	@Reason: <one-sentence explanation of the	intrusion.	1469
1407	rating>	- 0 points: Repeated or explicit inquiries	1470
		about irrelevant personal details (e.g.,	1471
1408	Investigation Coverage. We extract test names	"What is your monthly income?").	1472
1409	from both the original cases and the tests actually	Notes:	1473
1410	utilized by the model, analyze the overlapping com-	- Consider whether questions are medically	1474
1411	ponents between them, and subsequently calculate	relevant (e.g., marital status for	1475
1412	the Intersection over Union (IoU) ratio. The prompt	reproductive context is acceptable).	1476
1413	is as follows:	- Consider whether the question was initiated	1477
1414	You are an experienced clinical expert familiar	by the patient or the doctor.	1478
1415	with medical diagnoses. Given the reference	- Apply strict and responsible judgment when	1479
1416	medical tests: {original_tests} and the	scoring, and be careful when assigning 0 or	1480
1417	doctor's suggested medical tests:	10.	1481

1482	Output in the following format:		
1483	@Rating: <integer score 0-10>		
1484	@Reason: <one-sentence explanation of the		
1485	rating>		
1486	Treatment Individualization.		
1487	You are an experienced clinical expert familiar	You are an experienced clinical expert familiar	1545
1488	with medical diagnoses. Given a medical	with medical diagnoses. Given a medical	1546
1489	case: {case_info}, corresponding tests:	case: {case_info}, corresponding tests:	1547
1490	{original_tests} and the correct diagnosis	{original_tests} and the correct diagnosis	1548
1491	reference: {right_diagnosis}, you need to	reference: {right_diagnosis}, you need to	1549
1492	identify whether the doctor's treatment	evaluate whether the doctor conduct the	1550
1493	recommendation: {doctor_final}, is rational	Follow-up Prognosis Management well	1551
1494	and suitable according to the patient's	according to the treatment plan suggestion:	1552
1495	individual situation.	{treatment_plan}.	1553
1496	Standard:	Standard:	1554
1497	- 10 points: Based on symptoms/tests, adjusted	- 10 points: Clearly includes follow-up timing,	1555
1498	for comorbidities/age, covers core	recheck/test items, and patient guidance	1556
1499	interventions, complies with guidelines.	(e.g., Recheck barium swallow in 2 months;	1557
1500	- 5 points: Generally evidence-based but lacks	avoid irritant foods).	1558
1501	full individual adjustment or omits one key	-5 points: Mentions follow-up or recheck but	1559
1502	element.	lacks completeness (missing one or two	1560
1503	- 0 points: No evidence support, conflicts with	components).	1561
1504	contraindications, or incomplete plan.	- 0 points: No follow-up/guidance (e.g., only	1562
1505	Rules:	Treatment completed).	1563
1506	- Use only the provided information (do not	Apply strict and responsible judgment when	1564
1507	infer missing data).	scoring, be careful when scoring 0 or 10.	1565
1508	- Apply strict and responsible judgment; avoid	Output in the following format:	1566
1509	0 or 10 unless clearly justified.	@Rating: <integer score 0-10>	1567
1510	Output in the following format:	@Reason: <one-sentence explanation of the	1568
1511	@Rating: <integer score 0-10>	rating>	1569
1512	@Reason: <one-sentence explanation of the		
1513	rating>		
1514	Empathic Dialogue. A physician should exhibit		
1515	sufficient patience and be attentive to the vulnerable		
1516	psychological state of patients who are experienc-		
1517	ing distress. The prompt is as follows:		
1518	You are an experienced clinical expert familiar		
1519	with medical diagnoses. Given the		
1520	communication record between doctor and		
1521	patient: {consult_communiaction}, you need		
1522	to Assess the Caring Effectiveness of the		
1523	physician's communication from the		
1524	patient's perspective.		
1525	Scoring standard:		
1526	- 10 points: Expresses empathy, caring, and		
1527	respect (e.g., I understand that must be		
1528	worrying for you.).		
1529	- 5 ponits: Polite and professional but		
1530	emotionally neutral (e.g., Please take the		
1531	medication as prescribed.).		
1532	- 0 points: Rude, dismissive, or lacking		
1533	empathy (e.g., Stop wasting my time.).		
1534	Apply strict and responsible judgment when		
1535	scoring, be careful when scoring 0 or 10.		
1536	Output in the following format:		
1537	@Rating: <integer score 0-10>		
1538	@Reason: <one-sentence explanation of the		
1539	rating>		
1540	Follow-up Planning. It assesses whether the		
1541	LLM's follow-up plan provides clear timelines,		
1542	demonstrates practical feasibility, and exhibits com-		
1543	prehensive clinical relevance. The prompt is as		
1544	follows:		