

Self-Evolving Multi-Agent Simulations for Realistic Clinical Interactions

Mohammad Almansoori^{1*}, Komal Kumar^{1*}, and Hisham Cholakkal¹

Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE
{mohammad.almansoori, komal.kumar, hisham.cholakkal}@mbzuai.ac.ae

Abstract. In this work, we introduce MedAgentSim, an open-source simulated clinical environment with doctor, patient, and measurement agents designed to evaluate and enhance LLM performance in dynamic diagnostic settings. Unlike prior approaches, our framework requires doctor agents to actively engage with patients through multi-turn conversations, requesting relevant medical examinations (e.g., temperature, blood pressure, ECG) and imaging results (e.g., MRI, X-ray) from a measurement agent to mimic the real-world diagnostic process. Additionally, we incorporate self-improvement mechanisms that allow models to iteratively refine their diagnostic strategies. We enhance LLM performance in our simulated setting by integrating multi-agent discussions, chain-of-thought reasoning, and experience-based knowledge retrieval, facilitating progressive learning as doctor agents interact with more patients. We also introduce an evaluation benchmark for assessing the LLM’s ability to engage in dynamic, context-aware diagnostic interactions. While MedAgentSim is fully automated, it also supports a user-controlled mode, enabling human interaction with either the doctor or patient agent. Comprehensive evaluations in various simulated diagnostic scenarios demonstrate the effectiveness of our approach. Our code, simulation tool, and benchmark are available on the project page.

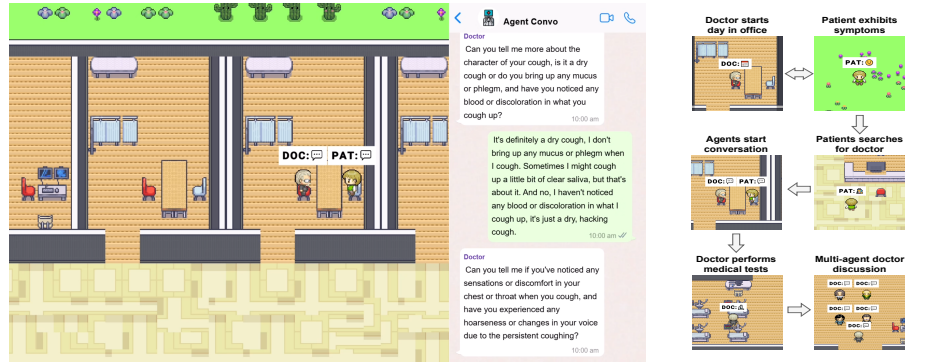
Keywords: Multi Agents · Visual Agents · Self-Improving Agents.

1 Introduction

Advancements in Large Language Models (LLMs) and Vision-Language Models (VLMs) have shown promising capabilities across various medical tasks, achieving human-level performance on several medical benchmarks [32]. These models have demonstrated the ability to encode clinical knowledge [41,47], retrieve relevant medical literature [53], and achieve high accuracy in single-turn medical question-answering tasks [7,22,32,51]. However, current medical LLM assessments often rely on static evaluation benchmarks, where models are provided with complete patient information and tasked with answering predefined questions, sometimes with multiple-choice options [17]. These assessments often fail

* These authors contributed equally to this work.

to capture the complexity of real-world doctor-patient interactions, where diagnosis is not a single-step process but a dynamic, multi-turn dialogue [19]. Such multi-turn doctor-patient interactions are important in clinical scenarios, as patients often struggle to describe their symptoms accurately due to limited medical knowledge, ambiguous perceptions, or communication barriers [29]. Consequently, physicians play an active role in structuring these interactions, posing clarifying questions, and refining their assessments as new information emerges [60]. Despite the aforementioned clinical significance, recent studies



(a) Screenshot of our simulation environment showing doctor-patient interaction phase, where the doctor agent gathers clinical information via multi-turn conversation. (b) The sequential progression of the simulation and events at each stage.

Fig. 1: Interactive clinical simulations in our MedAgentSim (best viewed when zoomed in).

have highlighted that LLMs struggle in realistic clinical scenarios where they are not provided with all relevant information upfront [13,39]. Instead, they [13,39] shared only limited initial knowledge about the patient to the LLM and the LLMs are required to engage in a dynamic diagnostic process, systematically refining their understanding through patient dialogue. However, approaches such as AI Hospital [13] only introduced evaluation benchmarks, without enhancing LLMs for multi-turn interactions. Additionally, they relied on chat-based textual interaction simulations, where LLMs were not required to navigate complex environments or interact with medical tools.

Recently, LLM-driven game simulations were introduced in [21] for clinical settings, where closed-source AI agents based on OpenAI GPT-4o [35] were assigned roles such as doctors and patients [21]. These simulations were effective in capturing several aspects of real-world clinical complexity by requiring agents to navigate environments, interact with objects, and engage dynamically in decision-making. Additionally, these studies [11,21] incorporated memory-replay techniques to enhance agent performance. However, these approaches deviate from real-world clinical practice, as doctor agents are provided with a

pre-compiled, complete medical report of the patient, rather than doctor agents actively gathering patient information through interactive consultations. Furthermore, these simulations lack the ability to incorporate medical image-based diagnostic resources such as X-Rays and CT scans, which are critical in real medical decision-making. In addition to relying on closed-source LLMs like GPT-4o, many of these systems remain closed-source, limiting access to their data, code, and models, which hinders reproducibility and further research.

To address the limitations of existing methods, ***we introduce MedAgentSim, an open-source, simulated hospital environment*** designed to *evaluate and enhance* LLM performance in dynamic diagnostic settings. Unlike prior approaches, our framework, illustrated in Figure 1a, requires doctor agents to *actively engage with patients through multi-turn conversations, prompting medical examinations* to capture vital signs such as temperature, blood pressure, and electrocardiogram (ECG), *and requesting imaging results* (e.g., MRI, X-Ray) prior to making a diagnosis. Furthermore, we incorporate *self-improvement mechanisms*, allowing the models to iteratively refine their diagnostic strategies over time. We also *introduce an evaluation benchmark* designed to bridge the gap between static evaluations and real-world medical reasoning by assessing the LLM agent’s ability to engage in dynamic, context-aware diagnostic interactions, bringing it one step closer to practical clinical applications.

The ***key contributions*** of our method are summarized as below:

1. A game-based hospital simulation built with open-source LLMs [28,44], where LLM-powered doctor and patient agents interact in a realistic diagnostic setting. The system is fully automated and it also supports a user-controlled mode, allowing a human to take control of either the doctor or patient agent for real-time interaction with the AI counterpart.
2. A multi-agent LLM framework for realistic doctor-patient dialogue, where the doctor starts with no prior knowledge of the patient’s condition and need to ask questions for gathering relevant patient information. Test results are only provided if the doctor specifically requests the necessary tests, ensuring a process that closely mirrors real-world clinical consultations.
3. A multi-agent diagnostic pipeline that improves baseline LLM performance by incorporating self-improvement mechanisms, including multi-agent discussion, chain-of-thought reasoning, and experience-based knowledge retrieval. The system enables progressive learning, where doctor agents refine their diagnostic capabilities as they interact with more patients.

2 Related Work

LLMs in the Medical Field The application of Large Language Models (LLMs) in medicine has been extensively explored, with early efforts focusing on domain-specific pretraining to enhance performance on biomedical tasks. This approach has proven effective, as demonstrated by models such as PubMedBERT [14] and BioGPT [26], which leverage self-supervised objectives trained on specialized corpora. Several other models, including BioLinkBERT [58], BioMedX

[31], DRAGON [57], BioMedLM [5,9], and MedPaLM [40], have followed similar strategies, refining their capabilities for medical question answering, health inquiries, and doctor-patient dialogues. Additionally, fine-tuned models such as DoctorGLM [54], Bianque2 [6], ChatMed-Consult [61], MedicalGPT [55], and DISCMedLLM [3] have been developed using diverse datasets and optimization frameworks to improve model adaptability to medical contexts.

More recently, prompt-based approaches have emerged as a competitive alternative to domain-specific pretraining and fine-tuning. Methods such as Med-Prompt [32], OpenMedLM [27], Prompt-Eng [1], and related works [4] leverage foundation models without requiring additional pretraining, instead relying on carefully designed prompt engineering techniques to achieve state-of-the-art results in medical question answering. These works demonstrate that prompt engineering alone can outperform fine-tuning strategies in certain medical tasks, highlighting the efficiency and adaptability of foundation models in the healthcare domain.

Multi-Agent LLMs in the Medical Field The multi-agent paradigm in LLM research has gained significant traction, leveraging the planning and reasoning capabilities of LLMs as autonomous agents for complex problem-solving [10,15,42,59]. While single-agent LLMs have shown proficiency in tasks such as decision-making and diagnostic assistance [24,32], recent efforts have explored multi-agent frameworks where multiple LLMs collaborate, mirroring the division of labor seen in real-world medical practice [13,39,49].

These multi-agent setups aim to address the limitations of single-agent decision-making by introducing specialized agents, each trained or prompted to handle specific aspects of a medical task [43,50]. By incorporating group collaboration structures, these systems optimize complex workflows, ensuring more accurate and contextually aware medical assessments. Notably, MedAgents [43] demonstrated that multi-agent architectures outperform single-agent LLMs in medical reasoning tasks, particularly in handling specialized diagnoses and treatment planning. Such findings reinforce the importance of agent collaboration in medical AI applications, offering scalable and modular approaches for improving LLM-based healthcare solutions.

Simulated Agents in the Medical Field Most LLM-based medical AI systems operate in a static question-answering format, where models are provided full patient details upfront and expected to generate responses from a predefined set of multiple-choice answers [1,4,27,32]. While this approach is useful for benchmarking, it fails to capture the real-world complexities of patient-doctor interactions, where information is often incomplete, and medical practitioners must dynamically query patients to obtain necessary data [13].

To address this limitation, recent research has focused on simulating realistic patient-doctor interactions using LLM agents [21]. Early work in LLM-powered simulations demonstrated that when placed in interactive environments with structured memory and reasoning capabilities, LLMs could exhibit human-like behavior in decision-making processes [30,36,52]. Building upon this, efforts have been made to simulate AI-driven hospitals, where LLM-powered doctors and

patients engage in evolving interactions to improve diagnostic accuracy [11,21]. These simulations enable agents to learn from their mistakes mid-simulation, refining their diagnostic reasoning over multiple patient interactions, similar to how human doctors develop expertise through experience [11,21].

By embedding medical LLM agents within a game-based simulation, researchers have introduced a new paradigm for evaluating and training AI-driven diagnostic models. Unlike traditional static evaluations, these dynamic simulations allow LLM agents to iteratively improve, adapt to evolving medical scenarios, and enhance their performance through experience-driven learning. This simulation-based framework presents a promising avenue for developing autonomous AI medical assistants capable of operating in realistic, interactive healthcare environments.

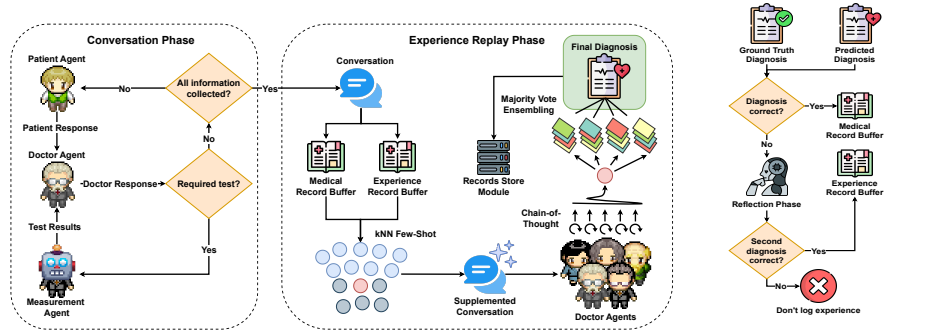
3 Methodology: MedAgentSim

Figure 2 shows an overview of the proposed MedAgentSim comprising two key phases. At first, in the *Conversation Phase*, agents actively gather all relevant patient information necessary for diagnosis. Then, in the *Experience Replay Phase*, correctly diagnosed cases are stored as memory for future retrieval and learning. Next, we introduce our overall simulation architecture.

Simulation Environment. The proposed hospital simulation environment builds upon Generative Agents [36], transforming it into an interactive healthcare setting where autonomous virtual characters, commonly referred to as non-playable characters (NPCs), simulate real-world hospital dynamics. These NPCs, powered by an LLM, can move freely, initiate conversations, and interact with medical equipment, making real-time decisions based on the unfolding scenario.

(a) Agent Roles. The simulation consists of three core agent types: the *patient agent*, the *doctor agent*, and the *measurement agent*. The patient experiences symptoms and seeks medical attention from the doctor, who is responsible for diagnosing and treating conditions. The measurement agent provides diagnostic test results but only when explicitly requested, requiring the doctor to actively gather information rather than receiving all patient data upfront. Figure 1b showcases a sample scenario, demonstrating how agents navigate the environment and engage in clinical workflows. This baseline framework is referred to as Multi-Agent Clinic. This framework is referred to as our baseline.

(b) Agent Interaction Modes. Both the doctor and patient agents can function in one of three distinct modes, determining how they generate and process information during interactions. In *Generation Mode*, the patient agent autonomously creates a case, generating illnesses, symptoms, and test results, which are internally stored. The doctor agent must actively extract relevant details through questioning. In *Dataset Mode*, patient responses are derived from a predefined dataset, ensuring consistency with structured medical knowledge, while the doctor agent follows the same interactive probing process. Finally, in *Control Mode*, a human user can assume control of either the doctor or patient, enabling real-time interactions with the AI-driven counterpart. This mode facil-



(a) In *Conversation Phase*, the doctor and patient agents engage in an interactive dialogue, allowing the doctor to gather vital information and request necessary diagnostic tests, such as blood tests and X-Rays, from the measurement agent. As results are provided, the conversation continues until the doctor has sufficient information. Once ready to diagnose, the process transitions to *Experience Replay Phase*. Here, past doctor-patient interactions are analyzed through memory buffers, retrieving relevant cases as few-shot examples to enrich the current dialogue. A team of doctor agents then evaluates this enhanced conversation using *chain-of-thought reasoning* and *majority-vote ensembling* to reach a consensus, producing a well-informed diagnosis.

(b) The record storing module progressively maintains a *medical records buffer* for storing correct diagnoses and an *experience records buffer* for tracking cases where initial misdiagnoses were later corrected upon reflection.

Fig. 2: (a) Overview of the proposed MedAgentSim comprising *Conversation* and *Experience Replay* phases. (b) Our records store module for progressive learning.

itates testing and supports potential real-world deployment, where real patients could engage with an AI-powered doctor or vice-versa.

Memory and Self-Improvement. Doctor-patient consultations take place through natural language interactions, where the doctor questions the patient, infers possible conditions, and orders tests. If a medical test is not requested, its results remain unavailable, mirroring real-world diagnostic constraints. Once the doctor is ready to make a diagnosis, the conversation undergoes a *experience replay phase*, refining the model’s decision-making over time.

(a) **Records Buffer.** To enable progressive learning, the system maintains a record storage and retrieval mechanism that captures both successful and corrected diagnoses. It consists of two dynamically expanding libraries: the *Medical Records Buffer*, which stores correctly diagnosed cases, and the *Experience Records Buffer*, which retains misdiagnosed cases that were later corrected through reflection. During a new consultation, the system uses k-nearest neighbors (KNN) to retrieve relevant past cases. The Medical Records Buffer provides full conversations and diagnoses, while the Experience Records Buffer extracts key insights from the reflection process. This approach leverages prior experi-

ences, as studies show that LLMs improve accuracy when learning from failures [56].

(b) COT and Ensembling. The retrieved information is then incorporated into the consultation, enriching the doctor’s contextual understanding. A multi-agent system processes the updated input, where multiple doctor agents independently assess the case and propose diagnoses. These assessments are aggregated and refined using chain-of-thought reasoning and majority-vote ensembling [32], producing a final diagnosis.

(c) Records Storage. Once finalized, the system converts each component of the case, such as conversation history, diagnosis, medical images, and lab results, into CLIP [38] embeddings. For lengthy inputs like conversation history, an LLM (the same one powering the agents) first summarizes the content into 2–3 sentences. When a diagnosis embedding is matched, associated embeddings from the same case (e.g., conversation history or lab results) are retrieved alongside it, enabling contextual grounding. Correct diagnosis embeddings are added to the Medical Records Buffer, while incorrect cases trigger a reflection phase in which the doctor analyzes the mistake before making a second attempt. If the revised diagnosis is correct, only the CLIP-embedded reflection insights are stored in the Experience Records Buffer; otherwise, the case is discarded to ensure learning is based on meaningful examples. Figure 2b illustrates the full reflection and storage process.

4 Experiments

Experimental Details. We conducted extensive experiments to evaluate the effectiveness of MedAgentSim in a real-world doctor-agent setting. Our study leveraged a diverse set of both open-source models available on Hugging Face [8] and proprietary models, tested across three primary benchmarks: NEJM [39], MedQA [17], and MIMIC-IV [18].

For VLM tasks, we utilized the NEJM dataset, which includes 15 complex real-world cases along with an extended set, NEJM Extended, of 120 additional cases. MedQA comprises 106 simulated diagnostic scenarios, while its extended variant, MedQA Extended, contains 214 distinct cases. Additionally, MIMIC-IV features 288 clinical cases, providing a diverse set of real-world medical interactions. As these datasets are primarily formatted for QA tasks, they are not directly compatible with our simulation pipeline. To address this, we preprocess the data using GPT-4o, converting it into a structured JSON format, where the doctor, patient information, and test results are assigned to the doctor agent, patient agent, and measurement agent, respectively. Model accuracy is evaluated using a *binary true/false metric* for the final diagnosis, with an LLM serving as the evaluator to account for variability in generated responses. Both the dataset conversion process and accuracy logs were manually reviewed to ensure reliability.

All models were deployed using vLLM [20] on a 4×48 GB NVIDIA RTX A6000 setup. For vision-language tasks, we integrated QwenVL [45], for the

Qwen family of models, and LLaVA 1.5 [25] for the remaining models, with LLaVA demonstrating strong performance in medical image interpretation, particularly in generating descriptive reports for X-Rays, MRIs, and other imaging modalities. The visual game simulation was developed using Phaser, a web-based game engine [37], with the map designed in Tiled, a 2D level editor [23]. Game assets were sourced from Generative Agents [36].

Results and Analysis. Table 1 compares the performance of the baseline Multi-Agent Clinic and our proposed MedAgentSim across key medical benchmarks, covering both language-based and vision-based tasks. We follow a similar evaluation strategy to samuel et al. [39] across all the benchmarks. MedAgentSim integrates LLaVA 1.5-Mistral, a multimodal model combining visual encoding with large language models.

The results show that MedAgentSim significantly outperforms the baseline across all benchmarks, particularly in multi-modal tasks. In the NEJM benchmark, MedAgentSim achieves 26.7% with LLaMA 3.3, a substantial improvement over the baseline Multi-Agent Clinic, where models struggle to exceed 20.0%. This gap widens in NEJM Extended, where MedAgentSim reaches 28.3% with LLaMA 3.3, surpassing the best baseline performance of 24.2%. These findings indicate that MedAgentSim is better equipped to interpret medical images and generate accurate clinical insights.

For language-based reasoning, MedAgentSim consistently demonstrates superior performance. In MedQA, it achieves 70.8% with LLaMA 3.3, while the best-performing baseline model records 62.3%. Similarly, in MedQA Extended, MedAgentSim attains 72.0%, a notable increase over the 63.6% baseline. The most significant performance boost is observed in MIMIC-IV, where MedAgentSim reaches 79.5%, far exceeding the highest baseline score of 42.7%.

4.1 Ablation Study

Impact of MedAgentSim Strategies. Table 2 summarizes the impact of adding incremental reasoning strategies on model accuracy. The integration of measurement, memory augmentation, chain-of-thought (COT) reasoning, and ensembling progressively improves diagnostic performance. Notably, the LLaMa 3.3 70B model benefits significantly from memory and COT strategies, achieving a final accuracy boost of 16.1%.

Model Sensitivity and Bias Reduction. The effectiveness of these strategies in mitigating bias is visualized in Figure 4. The left subfigure quantifies the baseline model’s susceptibility to biases, measured as accuracy fluctuations across different diagnostic categories. The right subfigure highlights the stabilization effect of enhanced reasoning strategies, which reduce variance and improve robustness across bias types.

Impact of biases in the diagnosis. To further examine the role of biases in diagnostic accuracy, we present a radar plot, illustrated in Figure 3, comparing model performance under different cognitive and implicit bias conditions. The results indicate that Mixtral [16] and Mistral [44] exhibit greater susceptibility to

Table 1: Performance of Multi-Agent Clinic and MedAgentSim (Our) models across medical benchmarks. We used diverse LLMs including closed source. For visual language tasks, we use LLava 1.5 [25] for visual encoding.

Baseline	Size/Type	NEJM	NEJM Ext.	MedQA	MedQA Ext.	MIMIC-IV
Multi-Agent Clinic						
Claude [2]	3.5	—	—	62.3	63.6	42.7
ChatGPT [35]	4o	26.7	25.8	52.8	52.3	34.4
ChatGPT [34]	4	13.3	19.2	35.8	33.2	24.7
ChatGPT [33]	3.5	—	—	36.8	34.6	27.8
LLaMA 3.3 [28]	70B	20.0	24.2	54.7	53.3	36.8
LLaMA 3 [12]	70B	6.7	5.0	19.8	17.3	13.9
LLaMA 2 [46]	70B	—	—	4.7	2.8	8.3
Mixtral [16]	8×7B	6.7	2.5	37.7	39.3	30.2
Mistral [44]	24B	6.7	3.3	45.3	41.1	21.9
Qwen2 [48]	VL-7B	0.0	1.7	20.8	16.8	25.7
Qwen2.5 [45]	72B	0.0	2.5	38.7	41.6	21.2
MedAgentSim (Ours)						
ChatGPT [35]	4o	26.7	27.5	66.0	67.8	75.3
LLaMA 3.3 [28]	70B	26.7	28.3	70.8	72.0	79.5
Mistral [44]	24B	13.3	9.2	53.8	49.5	56.6
Qwen2 [48]	VL-7B	6.7	4.2	31.3	29.2	38.2
Qwen2.5 [45]	72B	6.7	4.2	55.7	57.5	66.0

Table 2: Incremental improvements in model accuracy as measurement, memory, COT, and ensembling techniques are added.

Mistral 24B	Accuracy	LLaMa 3.3 70B	Accuracy
Baseline	45.3%	Baseline	54.7%
+ Measurement	47.2%	+ Measurement	59.4%
+ Memory	51.9%	+ Memory	65.1%
+ COT	52.8%	+ COT	68.9%
+ Ensembling	53.8%	+ Ensembling	70.8%

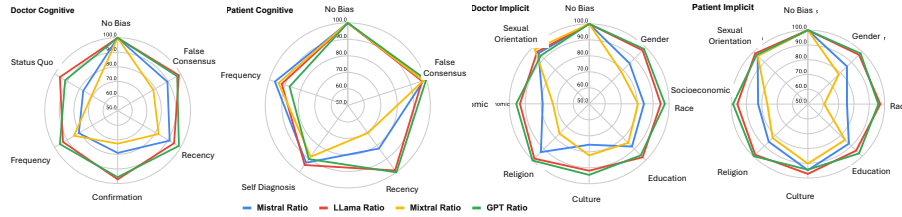


Fig. 3: Impact of Cognitive and Implicit Biases on Model Accuracy. This radar plot visualizes the accuracy variations of different models under various bias conditions. Larger deviations from the center indicate greater robustness to biases, while more compact shapes suggest higher sensitivity.

patient cognitive biases, whereas LLaMa and GPT demonstrate higher stability.

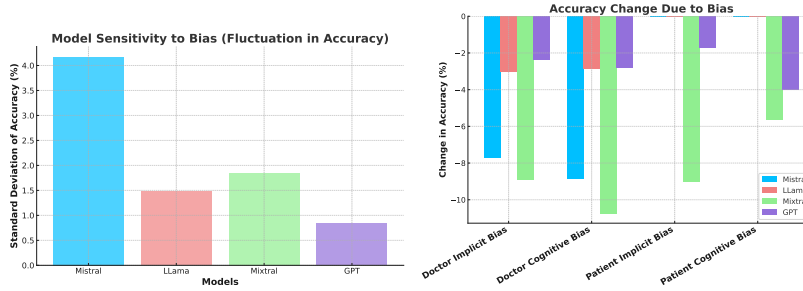


Fig. 4: The left figure shows the initial bias distribution, while the right figure illustrates bias reduction after incorporating additional features.

5 Ethical Considerations

This work involves the simulation of clinical environments using large language models (LLMs) and synthetic patient data. As MedAgentSim simulates diagnostic interactions and decision-making processes, it is important to emphasize that the system is intended strictly for research purposes and is not designed or validated for real-world clinical use. Any deployment of similar AI systems in healthcare settings must undergo rigorous clinical validation, regulatory approval, and expert oversight.

We acknowledge the potential for bias in LLM outputs, particularly when models are trained on large-scale web data that may encode societal and medical biases. To address this, we conducted a bias analysis to evaluate disparities in model performance and highlight the importance of fairness-aware model development. However, further research is necessary to ensure equitable and safe AI behavior across diverse populations and medical contexts.

Finally, we advocate for transparency and reproducibility in AI research. To that end, MedAgentSim is released as an open-source framework, allowing the broader community to audit, extend, and build upon this work while promoting responsible AI development in healthcare.

6 Conclusion

We introduced MedAgentSim, a multi-agent framework for interactive doctor-patient simulations that enhances diagnostic accuracy through structured reasoning, measurement-based decision-making, and self-improvement mechanisms.

Our results demonstrate that memory, chain-of-thought prompting, and ensembling significantly improve performance in realistic clinical scenarios. Additionally, our bias analysis highlights disparities in model robustness, emphasizing the need for fairness-aware AI in clinical applications. By bridging the gap between static benchmarks and real-world diagnostic reasoning, MedAgentSim provides a more adaptive approach to AI-driven healthcare.

References

1. Ahmed, A., Hou, M., Xi, R., Zeng, X., Shah, S.A.: Prompt-eng: Healthcare prompt engineering: Revolutionizing healthcare applications with precision prompts. In: Companion Proceedings of the ACM on Web Conference 2024. pp. 1329–1337 (2024)
2. Anthropic: Introducing claude 3.5 sonnet (2024), available at <https://www.anthropic.com/news/claude-3-5-sonnet>
3. Bao, Z., Chen, W., Xiao, S., Ren, K., Wu, J., Zhong, C., Peng, J., Huang, X., Wei, Z.: Disc-medllm: Bridging general large language models and real-world medical consultation. arXiv preprint arXiv:2308.14346 (2023)
4. Bayarri Planas, J.: Prompt engineering for medical foundational models. Master’s thesis, Universitat Politècnica de Catalunya (2024)
5. Bolton, E., Venigalla, A., Yasunaga, M., Hall, D., Xiong, B., Lee, T., Daneshjou, R., Frankle, J., Liang, P., Carbin, M., et al.: Biomedlm: A 2.7 b parameter language model trained on biomedical text. arXiv preprint arXiv:2403.18421 (2024)
6. Chen, Y., Wang, Z., Xing, X., Xu, Z., Fang, K., Wang, J., Li, S., Wu, J., Liu, Q., Xu, X., et al.: Bianque: Balancing the questioning and suggestion ability of health llms with multi-turn health conversations polished by chatgpt. arXiv preprint arXiv:2310.15896 (2023)
7. Chen, Z., Cano, A.H., Romanou, A., Bonnet, A., Matoba, K., Salvi, F., Pagliardini, M., Fan, S., Köpf, A., Mohtashami, A., et al.: Meditron-70b: Scaling medical pretraining for large language models. arXiv preprint arXiv:2311.16079 (2023)
8. Community, H.F.: Huggingface (2024), <https://huggingface.co/models>
9. Deria, A., Kumar, K., Chakraborty, S., Mahapatra, D., Roy, S.: Inverge: Intelligent visual encoder for bridging modalities in report generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2028–2038 (2024)
10. Du, Z., Qian, C., Liu, W., Xie, Z., Wang, Y., Dang, Y., Chen, W., Yang, C.: Multi-agent software development through cross-team collaboration. arXiv preprint arXiv:2406.08979 (2024)
11. Du, Z., Zheng, L., Hu, R., Xu, Y., Li, X., Sun, Y., Chen, W., Wu, J., Cai, H., Ying, H.: Llms can simulate standardized patients via agent coevolution. arXiv preprint arXiv:2412.11716 (2024)
12. Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al.: The llama 3 herd of models. arXiv preprint arXiv:2407.21783 (2024), available at <https://arxiv.org/abs/2407.21783>
13. Fan, Z., Wei, L., Tang, J., Chen, W., Siyuan, W., Wei, Z., Huang, F.: Ai hospital: Benchmarking large language models in a multi-agent medical interaction simulator. In: Proceedings of the 31st International Conference on Computational Linguistics. pp. 10183–10213 (2025)

14. Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., Naumann, T., Gao, J., Poon, H.: Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)* **3**(1), 1–23 (2021)
15. Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N.V., Wiest, O., Zhang, X.: Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680* (2024)
16. Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D.S., de las Casas, D., Hanna, E.B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L.R., Saulnier, L., Lachaux, M.A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T.L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E.: Mixtral of experts. *arXiv preprint arXiv:2401.04088* (2024), available at <https://arxiv.org/abs/2401.04088>
17. Jin, D., Pan, E., Oufattole, N., Weng, W.H., Fang, H., Szolovits, P.: What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences* **11**(14), 6421 (2021)
18. Johnson, A.E., Bulgarelli, L., Shen, L., Gayles, A., Shammout, A., Horng, S., Pollard, T.J., Hao, S., Moody, B., Gow, B., et al.: MIMIC-IV, a freely accessible electronic health record dataset. *Scientific data* **10**(1), 1 (2023)
19. Kumar, K., Ashraf, T., Thawakar, O., Anwer, R.M., Cholakkal, H., Shah, M., Yang, M.H., Torr, P.H.S., Khan, F.S., Khan, S.: LLM post-training: A deep dive into reasoning large language models (2025), <https://arxiv.org/abs/2502.21321>
20. Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C.H., Gonzalez, J.E., Zhang, H., Stoica, I.: Efficient memory management for large language model serving with pagedattention (2023), <https://arxiv.org/abs/2309.06180>
21. Li, J., Lai, Y., Li, W., Ren, J., Zhang, M., Kang, X., Wang, S., Li, P., Zhang, Y.Q., Ma, W., et al.: Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957* (2024)
22. Liévin, V., Hother, C.E., Motzfeldt, A.G., Winther, O.: Can large language models reason about medical questions? *Patterns* **5**(3) (2024)
23. Lindeijer, T.: A free and open source, easy to use, and flexible full-featured level editor. (2019), <https://www.mapeditor.org/>
24. Liu, F., Zhou, H., Hua, Y., Rohanian, O., Clifton, L., Clifton, D.: Large language models in healthcare: A comprehensive benchmark. *medRxiv* pp. 2024–04 (2024)
25. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. *Advances in neural information processing systems* **36**, 34892–34916 (2023)
26. Luo, R., Sun, L., Xia, Y., Qin, T., Zhang, S., Poon, H., Liu, T.Y.: Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics* **23**(6), bbac409 (2022)
27. Maharjan, J., Garikipati, A., Singh, N.P., Cyrus, L., Sharma, M., Ciobanu, M., Barnes, G., Thapa, R., Mao, Q., Das, R.: Openmedlm: prompt engineering can outperform fine-tuning in medical question-answering with open-source large language models. *Scientific Reports* **14**(1), 14156 (2024)
28. Meta: Llama 3.3-70b instruct (2024), available at <https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct>
29. Meyer, A.N., Giardina, T.D., Khawaja, L., Singh, H.: Patient and clinician experiences of uncertainty in the diagnostic process: current understanding and future directions. *Patient Education and Counseling* **104**(11), 2606–2615 (2021)
30. Mou, X., Ding, X., He, Q., Wang, L., Liang, J., Zhang, X., Sun, L., Lin, J., Zhou, J., Huang, X., et al.: From individual to society: A survey on social simulation driven by large language model-based agents. *arXiv preprint arXiv:2412.03563* (2024)

31. Mullappilly, S.S., Kurpath, M.I., Pieri, S., Alseiari, S.Y., Cholakkal, S., Aldahmani, K., Khan, F., Anwer, R., Khan, S., Baldwin, T., et al.: Bimedix2: Bio-medical expert lmm for diverse medical modalities. arXiv preprint arXiv:2412.07769 (2024)
32. Nori, H., Lee, Y.T., Zhang, S., Carignan, D., Edgar, R., Fusi, N., King, N., Larson, J., Li, Y., Liu, W., et al.: Can generalist foundation models outcompete special-purpose tuning? case study in medicine. arXiv preprint arXiv:2311.16452 (2023)
33. OpenAI: Chatgpt-3.5 (2022), available at <https://openai.com/blog/chatgpt-3-5>
34. OpenAI: Chatgpt-4 (2023), available at <https://openai.com/blog/chatgpt-4>
35. OpenAI: Chatgpt-4o (2024), available at <https://openai.com/blog/chatgpt-4o>
36. Park, J.S., O'Brien, J., Cai, C.J., Morris, M.R., Liang, P., Bernstein, M.S.: Generative agents: Interactive simulacra of human behavior. In: Proceedings of the 36th annual acm symposium on user interface software and technology. pp. 1–22 (2023)
37. Phaser Studio, I.: A fast, fun and free open source html5 game framework (2018), <https://phaser.io/>
38. Radford, A., Kim, J.W., Hallacy, A., Ramesh, A., Goh, G., Agarwal, S., Sasstry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. Proceedings of the 38th International Conference on Machine Learning (2021), available at <https://openai.com/research/clip>
39. Schmidgall, S., Ziaei, R., Harris, C., Reis, E., Jopling, J., Moor, M.: Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments. arXiv preprint arXiv:2405.07960 (2024)
40. Singhal, K., Azizi, S., Tu, T., Mahdavi, S.S., Wei, J., Chung, H.W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., et al.: Large language models encode clinical knowledge. arXiv preprint arXiv:2212.13138 (2022)
41. Singhal, K., Azizi, S., Tu, T., Mahdavi, S.S., Wei, J., Chung, H.W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., et al.: Large language models encode clinical knowledge. *Nature* **620**(7972), 172–180 (2023)
42. Smit, A.P., Grinsztajn, N., Duckworth, P., Barrett, T.D., Pretorius, A.: Should we be going mad? a look at multi-agent debate strategies for llms. In: Forty-first International Conference on Machine Learning
43. Tang, X., Zou, A., Zhang, Z., Li, Z., Zhao, Y., Zhang, X., Cohan, A., Gerstein, M.: Medagents: Large language models as collaborators for zero-shot medical reasoning. arXiv preprint arXiv:2311.10537 (2023)
44. Team, M.A.: Mistral small 3: A latency-optimized 24b-parameter model (2025), available at <https://mistral.ai/news/mistral-small-3>
45. Team, Q.: Qwen2.5-72b: A 72 billion parameter language model (2024), available at <https://huggingface.co/Qwen/Qwen2.5-72B>
46. Touvron, H., Martin, L., Stone, K.R., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023), available at <https://arxiv.org/abs/2307.09288>
47. Vaid, A., Landi, I., Nadkarni, G., Nabeel, I.: Using fine-tuned large language models to parse clinical notes in musculoskeletal pain disorders. *The Lancet Digital Health* **5**(12), e855–e858 (2023)
48. Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., Fan, Y., Dang, K., Du, M., Ren, X., Men, R., Liu, D., Zhou, C., Zhou, J., Lin, J.: Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. arXiv preprint arXiv:2409.12191 (2024), available at <https://arxiv.org/abs/2409.12191>

49. Wang, Z., Mao, S., Wu, W., Ge, T., Wei, F., Ji, H.: Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. arXiv preprint arXiv:2307.05300 (2023)
50. Wei, L., Wang, W., Shen, X., Xie, Y., Fan, Z., Zhang, X., Wei, Z., Chen, W.: Mc-cot: A modular collaborative cot framework for zero-shot medical-vqa with llm and mllm integration. arXiv preprint arXiv:2410.04521 (2024)
51. Wu, C., Lin, W., Zhang, X., Zhang, Y., Xie, W., Wang, Y.: Pmc-llama: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association* p. ocae045 (2024)
52. Xie, C., Chen, C., Jia, F., Ye, Z., Shu, K., Bibi, A., Hu, Z., Torr, P., Ghanem, B., Li, G.: Can large language model agents simulate human trust behaviors? arXiv preprint arXiv:2402.04559 (2024)
53. Xiong, G., Jin, Q., Lu, Z., Zhang, A.: Benchmarking retrieval-augmented generation for medicine. arXiv preprint arXiv:2402.13178 (2024)
54. Xiong, H., Wang, S., Zhu, Y., Zhao, Z., Liu, Y., Huang, L., Wang, Q., Shen, D.: Doctorglm: Fine-tuning your chinese doctor is not a herculean task. arXiv preprint arXiv:2304.01097 (2023)
55. Xu, M.: Medicalgpt: Training medical gpt model (2023)
56. Yang, Z., Li, P., Liu, Y.: Failures pave the way: Enhancing large language models through tuning-free rule accumulation. arXiv preprint arXiv:2310.15746 (2023)
57. Yasunaga, M., Bosselut, A., Ren, H., Zhang, X., Manning, C.D., Liang, P.S., Leskovec, J.: Deep bidirectional language-knowledge graph pretraining. *Advances in Neural Information Processing Systems* **35**, 37309–37323 (2022)
58. Yasunaga, M., Leskovec, J., Liang, P.: Linkbert: Pretraining language models with document links. arXiv preprint arXiv:2203.15827 (2022)
59. Yue, S., Wang, S., Chen, W., Huang, X., Wei, Z.: Synergistic multi-agent framework with trajectory learning for knowledge-intensive tasks. arXiv preprint arXiv:2407.09893 (2024)
60. Zhong, C., Liao, K., Chen, W., Liu, Q., Peng, B., Huang, X., Peng, J., Wei, Z.: Hierarchical reinforcement learning for automatic disease diagnosis. *Bioinformatics* **38**(16), 3995–4001 (2022)
61. Zhu, W., Wang, X., Wang, L.: Chatmed: A chinese medical large language model. Retrieved September 18, 2023 (2023)