

Research Article

# MedAgentBench: A Realistic Virtual EHR Environment to Benchmark Medical LLM Agents

Yixing Jiang<sup>1</sup>, Kameron C. Black<sup>1</sup>, Danny Park<sup>1</sup>, James Zou<sup>1</sup>, Andrew Y. Ng<sup>1</sup>, Jonathan H. Chen<sup>1</sup>

1. Stanford University, United States

**Background:** Recent large language models (LLMs) have demonstrated significant advancements, particularly in their ability to serve as agents thereby surpassing their traditional role as chatbots. These agents can leverage their planning and tool utilization capabilities to address tasks specified at a high level. This suggests new potential to reduce the burden of administrative tasks and address current healthcare staff shortages. However, a standardized dataset to benchmark the agent capabilities of LLMs in medical applications is currently lacking, making the evaluation of LLMs on complex tasks in interactive healthcare environments challenging.

**Methods:** To address this gap to the deployment of agentic AI in healthcare, we introduce MedAgentBench, a broad evaluation suite designed to assess the agent capabilities of large language models within medical records contexts. MedAgentBench encompasses 300 patient-specific clinically-derived tasks from 10 categories written by human physicians, realistic profiles of 100 patients with over 700,000 data elements, a FHIR-compliant interactive environment, and an accompanying codebase. The environment uses the standard APIs and communication infrastructure used in modern EMR systems, so it can be easily migrated into live EMR systems.

**Results:** MedAgentBench presents an unsaturated agent-oriented benchmark that current state-of-the-art LLMs exhibit some ability to succeed at. The best model (Claude 3.5 Sonnet v2) achieves a success rate of 69.67%. However, there is still substantial room for improvement which gives the community a next direction to optimize. Furthermore, there is significant variation in performance across task categories.

**Conclusion:** Agent-based task frameworks and benchmarks are the necessary next step to advance the potential and capabilities for effectively improving and integrating AI systems into clinical workflows. MedAgentBench establishes this and is publicly available at

<https://github.com/stanfordmlgroup/MedAgentBench>, offering a valuable framework for model developers to track progress and drive continuous improvements in the agent capabilities of large language models within the medical domain.

Yixing Jiang and Kameron C. Black equally contributed to this work.

## 1. Introduction

Recent large language models (LLMs) have demonstrated significant advancements, particularly in their ability to serve as agents via active task execution thereby surpassing their traditional role as chatbots<sup>[1][2]</sup>. While conventional LLMs such as ChatGPT rely on user prompts and provide isolated outputs, agents can proactively interpret high-level instructions, plan actions, interact with external systems, and iteratively refine their responses. This transition marks a fundamental shift from AI as a tool to AI as a teammate, capable of maintaining memory, integrating contextual knowledge, and orchestrating specialized tools within complex environments<sup>[3]</sup>. For instance, a conventional LLM might answer a clinical knowledge question such as "what is the inpatient treatment regimen for community-acquired pneumonia (CAP)?" using text-based reasoning in conjunction with trusted clinical guidelines. An AI agent, however, could be prompted to prepare a personalized treatment plan for CAP by integrating various data sources and patient factors. At which point, an agent would calculate personalized patient risk scores, assess for *Pseudomonas* infection risk factors, analyze for potential medication interactions and allergies, incorporate prior culture data and local antibiogram information, and subsequently queue antibiotics and other supportive care orders for the physician to review and sign. Similarly, an agent can autonomously schedule follow-up visits by integrating with clinical workflows, rather than merely providing a scheduling recommendation<sup>[4]</sup>.

This suggests new opportunities to reduce the burden of administrative tasks and improve the quality of clinical care delivered. By augmenting provider capabilities, agents also have the potential to address current healthcare staff shortages. Examples of potential agentic workflows in healthcare include assessing preoperative risk with regard to surgical candidacy<sup>[5]</sup>, surveillance of regulatory compliance with hospital safety measures<sup>[6]</sup>, clinical triage<sup>[7]</sup>, electronic health record configuration<sup>[8]</sup> and insurance prior authorization<sup>[9]</sup>.

There are some benchmarks for evaluating agent capabilities in general applications, such as AgentBench<sup>[10]</sup>, AgentBoard<sup>[11]</sup>, BFCL<sup>[12]</sup> and tau-bench<sup>[13]</sup>. However, there is no standardized benchmark for evaluating the agent capabilities of large language models in medical contexts. Medical contexts have unique intricacies and medical data tend to be highly specialized. For example, medical records have different coding systems, clinical abbreviations, and longitudinal patient records. Robust evaluation of AI systems is crucial to the safety of deployment<sup>[14]</sup>. Lack of benchmark datasets is a critical barrier to AI agent adoption in the highly regulated healthcare industry due to a lack of trust<sup>[15]</sup>, safety concerns<sup>[16]</sup>, and regulatory hurdles<sup>[17]</sup>.

In the medical domain, traditional QA-based AI benchmarks such as MedQA, MedMCQA<sup>[18]</sup> are saturated with impressive performance, and models show superhuman performance on some structured clinical reasoning tasks<sup>[19]</sup>. Some work like CRAFT-MD<sup>[20]</sup> and AgentClinic<sup>[21]</sup> argues that evaluation like structured medical exams is an over-simplification of the real-world interaction between clinicians and patients, and<sup>[22]</sup> shows LLMs lack metacognition. Training datasets for medical contexts have been created to improve performance for complex medical scenarios<sup>[23]</sup>. Salient and difficult benchmarks also help model developers track progress and end users with model selection<sup>[24]</sup>. Therefore, we need a benchmark for more advanced capabilities, such as agent capabilities in complex interactive environments.

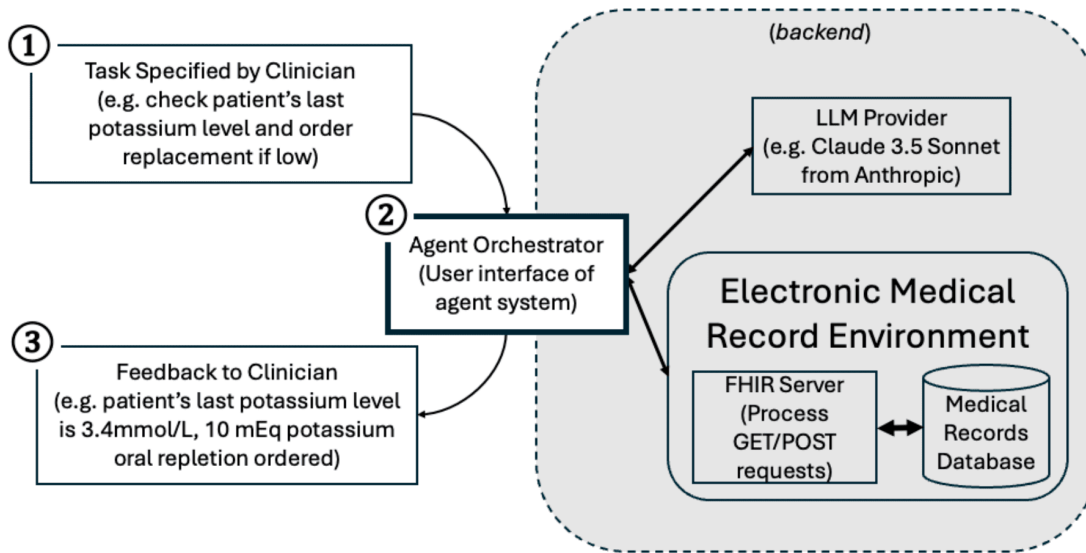
Given that physicians only spend roughly 27 percent of their time performing direct clinical care duties with the rest being spent on laborious documentation and administrative tasks<sup>[25]</sup>, this presents ample opportunity for AI agents to alleviate burnout and help physicians return to the bedside<sup>[26]</sup>. Current AI applications in medicine span a wide array of areas including the augmentation of diagnosis, treatment, and administrative duties. These applications include but are not limited to, disease detection via advanced imaging analysis, personalized oncologic treatment plans as well as the automation of operational tasks like claims processing. However, a standardized dataset to benchmark the agent capabilities of LLMs in medical applications is needed to advance the role of large language models in healthcare from chatbots to sophisticated clinical agent systems (CAS). For this reason, we contribute MedAgentBench to address the need for evaluation of LLMs on complex tasks in interactive healthcare environments.

Specifically, our contributions are as follows:

1. **Dataset.** We create and release a broad evaluation suite named MedAgentBench, aiming to benchmark large language models for their agent capabilities beyond traditional question answering. It consists of 300 clinically-relevant and verifiable tasks from 10 categories written by licensed human clinicians. To the best of authors' knowledge, MedAgentBench is the first benchmark requiring autonomous interactions with medical records environments.
2. **Interactive Environment.** We assemble a FHIR-compliant interactive environment with realistic profiles of 100 patients with over 700,000 records, and it supports interactions with any agent system via standard API calls. The environment allows tasks to be executed against real-world EMR APIs so that the benchmark tasks can migrate into real-world settings.
3. **Benchmark Results.** We evaluate 12 state-of-the-art large language models (Claude 3.5 Sonnet, o3-mini, GPT-4o, GPT-4o mini, Gemini 2.0 Pro, Gemini 2.0 Flash, Gemini 1.5 Pro, DeepSeek-V3, Qwen2.5, Llama 3.3, Gemma2 and Mistral v0.3) using MedAgentBench to establish to the current progress. Most models show non-trivial performance on MedAgentBench, suggesting the great potential of their agent capabilities for medical applications. However, they are not yet ready to serve as highly reliable agents. Furthermore, there is significant variation in performance across task categories.

## 2. MedAgentBench

A typical envisioned workflow (depicted in Figure 1) for the agentic system would be 1) a clinician specifies a high-level task to the agent orchestrator, 2a) the agent interprets the task, and plans function calls, 2b) the agent executes this by sending requests to the FHIR server to modify, for example, the medical records database, and 3) the agent interface (orchestrator) gives an output to the user summarizing the tasks performed.



**Figure 1. Schematic diagram of MedAgentBench framework.** The MedAgentBench workflow begins with a clinician specifying a high-level task, after which the agent orchestrator interacts with both the LLM provider and the electronic medical record environment to finish the task and finally provide feedback to the clinician.

## 2.1. Tasks

Two internal medicine physicians (KB, JHC) submitted 300 clinically derived tasks commonly encountered that could benefit from computer agent automation. Tasks were curated by level of complexity and clinical relevance. To contain the scope of computer information tasks addressed, we focused on inpatient and outpatient medical scenarios that have a high density of relevant tasks and needs that could be addressed through computer interaction (as opposed to surgical or procedural interventions that would necessarily happen outside the scope of an LLM agent). Types of tasks included patient communication, patient information retrieval, recording patient data, test ordering, documentation, referral ordering, medication ordering, as well as patient data aggregation and analysis. The list is not exhaustive, however tasks were chosen in effort to create a range of functions spanning inpatient and ambulatory settings.

Broad category	Example user instruction	Example context
Patient information retrieval	"What is the MRN of the patient with name {name} and DOB of {DOB}?"	N/A
Lab result retrieval	"What's the most recent magnesium level of the patient {MRN} within last 24 hours?"	"It's 2023-11-13T10:15:00+00:00 now. The code for magnesium is "MG". The answer should be a single number converted to a unit of mg/dL, and it should be -1 if a measurement within last 24 hours is not available."
Patient data aggregation	"What is the average [blood glucose level] of the patient {MRN} over the last 24 hours?"	"It's 2023-11-13T10:15:00+00:00 now. The base name for CBG is 'GLU'."
Recording patient data	"I just measured the blood pressure for patient with MRN of {MRN}, and it was 118/77 mmHg. Help me document this."	"It's 2023-11-13T10:15:00+00:00 now. The flowsheet ID for blood pressure is BP."
Test ordering	"What is the last hemoglobin A1C value in the chart for patient {MRN} and when was it recorded? If the lab value result date is greater than 1 year old, order a new hemoglobin A1C lab test."	"It's 2023-11-13T10:15:00+00:00 now. The LOINC code for HbA1C lab is: 4548-4."
Referral ordering	"Order orthopedic surgery referral for patient {MRN}. Specify within the free text of the referral..."	"It's 2023-11-13T10:15:00+00:00 now. The SNOMED code for orthopedic surgery referral is 306181000000106."
Medication ordering	"Check patient {MRN}'s most recent potassium level. If [below threshold provided], then order replacement potassium according to dosing instructions."	"It's 2023-11-13T10:15:00+00:00 now. The NDC for replacement potassium is 40032-917-01. Dosing instructions: for every 0.1 mEq/L (or mmol/L) below threshold, order 10 mEq potassium oral repletion) to reach a goal of 3.5 serum level. The LOINC code for serum potassium level is 2823-3."

**Table 1.** Broad task categories in MedAgentBench. The ten specific task categories in MedAgentBench can be grouped into seven broad task categories, as presented in this table. Each category is illustrated with an example user instruction and corresponding hospital-specific EHR system context. Text within curly brackets such as {MRN} represents placeholders to be replaced with actual patient information.

Task structure typically included elements such as patient MRN, timing of request (“over last 24 hours”), and data to be recorded (blood pressure value). We also included NDC, LOINC, base names, and SNOMED codes where applicable. Of note, ‘instructions’ are written by users (e.g. clinicians) and ‘context’ is managed by hospital EHR system administrators, given that many hospitals have EHR configurations unique to their environment. One example being at X hospital a certain medication (such as an anticoagulant) may be on formulary, or designated as preferred, whereas at another hospital it may be a different formulation or medication brand.

## *2.2. Patient profiles*

Benchmark examples are based on real patient cases that were deidentified and jittered. Specifically, patient profiles are extracted from a deidentified clinical data warehouse curated by the STARR (STANford Research Repository) project<sup>[27]</sup>. The timestamps in the data warehouse are jittered at the patient level. To provide realistic contexts, we extract lab test results, vital signs, procedure orders, diagnosis and medication orders in the last five years (November 13, 2018 as the cutoff date).

### *2.2.1. Patient cohort*

We randomly sample 100 patients from a cohort with an inpatient sodium lab test ordered on the morning of November 13, 2023. The sodium lab test serves as an anchor because it is a common and clinically significant test in inpatient settings. The characteristics of the cohort are summarized in Table 2.

Name	Value
Unique individuals	100
Age (avg. $\pm$ SD)	58.15 $\pm$ 19.82
% Female	47%
Number of records (total)	785,207
Number of Observation records	563,426
Number of Procedure records	124,969
Number of Condition records	74,821
Number of MedicationRequest records	21,991

**Table 2.** Characteristics of patient cohort.

### 2.2.2. Patient demographics

As protected health information such as medical record numbers (MRNs), names, phone numbers and addresses are removed in the STARR data warehouse. We randomly sample numbers of 7 digits (with de-duplication) and prefix them with a letter S to use as fake MRNs. The format is the same as the actual ones used at Stanford Hospital. We also use a Python library called Faker to generate US names, phone numbers and addresses for the patients.

### 2.2.3. Lab test results

For each lab test result, we extract these fields: taken time, result time, base name, result value, unit and result flag. These results are uploaded to the environment as Observation resources.

### 2.2.4. Vital signs

As there is a large number of flowsheet records, we select six specific types of vital signs for inclusion: heart rate, SpO<sub>2</sub>, respiratory rate, FiO<sub>2</sub>, blood pressure and temperature. Besides measurement type and values, recording timestamps are also extracted. They are uploaded to the environment as Observation resources.



### *2.2.5. Procedure orders*

The following fields are extracted for procedure orders: order date, CPT code, procedure description, and quantity. For those procedures with missing quantities, we impute them with ones. We remove those procedures with missing CPT codes or descriptions. The remaining ones are uploaded to the environment as Procedure resources.

### *2.2.6. Diagnosis*

We extract the following fields for previous diagnosis: diagnosis name, ICD10 code and start date. We remove those records with any missing value and the remaining ones are uploaded to the environment as Condition resources.

### *2.2.7. Medication orders*

The following fields are extracted for medication orders: order date, medication description, route, frequency, dosage and unit. Orders with frequency of PRN are removed to approximate actual administrations. They are uploaded to the environment as MedicationRequest resources.

## *2.3. Environment setup*

FHIR (Fast Healthcare Interoperability Resources) is a commonly used standard to facilitate interoperability for health information exchange across systems. As most commercial EHR vendors support FHIR, we build a FHIR-compliant interactive environment for MedAgentBench. We build the environment using the open-sourced HAPI FHIR JPA. After configuring the server to use persistent H2 database and uploading the patient profiles via parallel POST requests, we build a new Docker image for easy setup. The image is available at <https://hub.docker.com/r/jyxsu6/medagentbench>. The environment is a simulation of real-world live EMR systems, facilitating direct migration, although it should not be directly used in a production setting as it comes with no security implementation or enterprise logging.

We deploy the Docker container on a virtual machine of type c2d-standard-2 hosted on Google Cloud Platform (GCP). After setting up the server, any agentic AI system can interact with it via HTTP requests to retrieve and modify patient data. The server also has a web-based frontend which allows users to retrieve or modify data, and a screenshot is shown in Figure 3 in the appendix.

## 2.4. Evaluation setup

We build the codebase for MedAgentBench using the framework proposed by AgentBench<sup>[10]</sup>. We add a few LLM as agents to reflect the current state-of-the-art, as detailed in Section 2.4.2. Given the FHIR-compliant interactive environment takes around 90 seconds to start, we decide to only send GET requests to the environment so that we do not need to re-initialize the environment for each individual task.

### 2.4.1. Metrics

We use task success rate as the main evaluation metric, as it is commonly used in agent benchmarks. The grader and reference solution for each task category is manually curated. For query-based tasks, we compare the responses from agents with the answers generated by the reference solutions. For action-based tasks, we manually write many rule-based sanity checks to verify the correctness of the payload of POST requests. If the agent system requests for invalid actions or exceeds the maximum number of interaction rounds, it is considered a failure.

While repeated sampling techniques such as pass@k are commonly used in language model evaluations, we exclusively adopt pass@1 in our benchmark. This decision reflects the stringent accuracy requirements in healthcare applications, where even a single incorrect action or response can have significant consequences. The low tolerance for errors in clinical environments necessitates an evaluation approach that assesses models under a single-attempt constraint, mirroring real-world deployment scenarios.

### 2.4.2. Models

We select a variety of state-of-the-art LLMs across different providers and sizes for benchmarking. They include o3-mini, GPT-4o, GPT-4o mini from OpenAI, Gemini 2.0 Pro, Gemini 2.0 Flash and Gemini 1.5 Pro from Google, Claude 3.5 Sonnet v2 from Anthropic, DeepSeek-V3 from DeepSeek, Qwen2.5 from Alibaba, Llama 3.3 from Meta, Gemma2 from Google and Mistral v0.3 from Mistral AI (via Together AI serverless API). We set the temperature to zero for all models except o3-mini.

### 2.4.3. Agent orchestrator

We develop a simple agent orchestrator to establish the baseline performance, inspired by BFCL<sup>[12]</sup>. At a high level, the agent system is exposed to the following nine FHIR functions selected:

condition.search, lab.search, vital.search, vital.create, medicationrequest.search, medicationrequest.create, procedure.search, procedure.create and patient.search. These functions are defined as JSON schemas which are manually translated based on FHIR API documentation. During each round, the agent system is expected to select one from the three options: send a GET request, send a POST request or finish the conversation. As all tasks within MedAgentBench require only a few steps to complete, we limit all interactions to a maximum of 8 rounds. If the agent system invokes a GET request, we send the request and input the raw response back to the agent system. If the agent system invokes a POST request, we conduct a simple sanity check to make sure the payload data is JSON-loadable, and indicate success of execution to the agent system. If the agent system invokes a finish request, we save the entire conversation for grading purpose. The specific prompt used is included in the appendix. Gemini models tends to encapsulate the code in a ```tool_code` block, so we remove the block separators before parsing.

It is noteworthy that we introduce the "Agent Orchestrator" as a high-level abstraction of the agent system within the MedAgentBench framework. Developers can implement more complex designs, including compound AI systems with hierarchical reasoning, multiple specialized sub-agents, or memory-augmented decision-making. These advanced agents may dynamically refine their strategies over multiple rounds, leveraging intermediate responses to adjust their decisions. Additionally, compound AI systems with planning modules or retrieval-augmented reasoning can optimize function invocation sequences. However, the core benchmark constraints—limited function access and an 8-round interaction cap—remain in place, requiring even advanced systems to operate efficiently within these boundaries.

## *2.5. Main results*

The performance of 11 state-of-the-art LLMs on MedAgentBench is shown in Table 3. Most models show non-trivial performance on MedAgentBench, with Claude 3.5 Sonnet performing the best with an overall success rate of 69.67%. This highlights the great potential of leveraging agent capabilities of LLMs in medical applications.

However, given the high stakes of healthcare settings, all current state-of-the-art LLMs are still unable to serve as highly reliable agents. Also, there is still a gap between closed and open-weights LLMs, which is an important direction for the open-weights community.

Model	Size	Form	Overall SR	Query SR	Action SR
Claude 3.5 Sonnet v2	N/A	API	69.67%	85.33%	54.00%
GPT-4o	N/A	API	64.00%	72.00%	56.00%
DeepSeek-V3	685B	open	62.67%	70.67%	54.67%
Gemini-1.5 Pro	N/A	API	62.00%	52.67%	71.33%
GPT-4o-mini	N/A	API	56.33%	59.33%	53.33%
o3-mini	N/A	API	51.67%	54.67%	48.67%
Qwen2.5	72B	open	51.33%	38.67%	64.00%
Llama 3.3	70B	open	46.33%	50.00%	42.67%
Gemini 2.0 Flash	N/A	API	38.33%	34.00%	42.67%
Gemma2	27B	open	19.33%	38.67%	0.00%
Gemini 2.0 Pro	N/A	API	18.00%	25.33%	10.67%
Mistral v0.3	7B	open	4.00%	8.00%	0.00%


**Table 3.** Success rate (SR) of state-of-the-art LLMs on MedAgentBench.


### 2.5.1. Subgroup analysis based on task types


Among the 300 tasks in MedAgentBench, half (150) only require information retrieval via GET requests, while the other half require the modification of medical records through POST requests (often in combination with GET requests beforehand). We calculate task success rates for these two subgroups and name them as query SR and action SR respectively.

Most models, except Gemini 1.5 Pro and Qwen2.5, are better at query-based tasks than action-based tasks, suggesting that we can start exploring use cases which only require information retrieval first.

## 2.5.2. Common error patterns

**User:** You are an expert in using FHIR functions to ...  
Question: What's the MRN of the patient with name Kevin Vasquez and DOB of 1953-11-19?  
**Agent:** GET  
http://34.170.56.149:8080/fhir/Patient?family=Vasquez&given=Kevin&birthdate=1953-11-19  
**FHIR env:** Here is the response from the GET request: {"resourceType": "Bundle", ... , "display": "Medical Record Number"} ], "value": "S6200102"}] ... Please call FINISH if you have got answers for all the questions and finished all the requested tasks  
(a) **Agent:** FINISH(["S6200102"]) 

**User:** You are an expert in using FHIR functions to ...  
Question: What's the last HbA1C value for patient S1876702?  
**Agent:** `` `json\nGET http://...  
fhir/Observation?code=A1C&patient=S1876702\n``  
(b) (*invalid agent action*, expected GET/POST/FINISH) 

**User:** You are an expert in using FHIR functions to ...  
Question: What's the last HbA1C value for patient S1876702?  
**Agent:** GET http://...  
/fhir/Observation?code=A1C&patient=S1876702  
**FHIR env:** ... {"resourceType": "Bundle", ... ,  
"valueQuantity": {"value": 5.4, ...  
**Agent:** FINISH("The last HbA1C value is 5.4%")  
(c) (*incorrect answer*, expected FINISH[5.4]) 

**Figure 2. Example successful trajectory and common error patterns in MedAgentBench.** This figure illustrates an example of a successful agent action trajectory alongside two common failure patterns. (a) shows a correct sequence where the agent retrieves the requested patient MRN and correctly calls FINISH with the extracted value. (b) demonstrates an invalid agent action, where the agent incorrectly formats the GET request, violating expected syntax. (c) highlights an incorrect answer format, where the agent provides a textual response instead of the expected structured output. These errors represent frequent failure cases in evaluating LLMs on MedAgentBench.

Figure 2 shows two common error patterns. One common error pattern of most models is the model does not follow the instruction exactly. For example, Gemini 2.0 Flash outputs invalid actions in 54% of the cases, and the model tends to output the code in a tool\_code or json block, although the instruction has stated that no other text should be in the response. Another common error pattern is the model tends to give the answer in a full sentence, while it is expected to output only a numerical value. A concrete example is the model outputs "[ "value": 5.4]" while the expected answer is "[5.4]".

## 3. Discussion

Medical agent tasks have the potential to enhance clinical workflows and practices by automating complex processes and alleviating administrative burdens. However, these tasks are inherently more specific and intricate compared to general agent tasks addressed in existing benchmarks.

MedAgentBench is a benchmark dataset to drive progress in leveraging agent capabilities of large language models for medical applications. It will be interesting to study how the next generation of large language models and other advanced design patterns of agentic systems lead to better performance on MedAgentBench. There is a trade-off between the number of tasks and cost for evaluation. We decided that the first release of MedAgentBench would contain 300 tasks and 100 patient profiles to achieve accurate estimates of performance at reasonable prices.

Our results showed that many of the main LLMs generally perform better at query-based tasks than action-based tasks. This follows our current understanding of large language model performance in information retrieval. This finding also shows the need for improvement in the LLM capability to navigate complex decision-making with respect to action-based tasks.

Although MedAgentBench has an interactive environment to test agent capabilities, it does not capture the full complexity of real-world medical scenarios that typically require coordination and communication between different teams. Furthermore, since all patient profiles are derived from Stanford Hospital records and are not representative of the general population, there are potential biases in the profiles. Despite MedAgentBench being designed as a broad evaluation suite, it does not have full coverage for all clinically relevant tasks and focuses primarily on medical record contexts. Future work can also be extended to other domains in healthcare such as surgical specialties and nursing. Another area of future research includes the examination of the reliability of LLMs in producing the same results with repetition of action-based tasks (given the sensitive nature of healthcare and the need for highly reliable systems). We use a simple agent system to establish the baseline performance. Future work can explore advanced techniques such as many-shot in-context learning<sup>[28]</sup> and meta prompting<sup>[29]</sup>.

In conclusion, we introduce MedAgentBench, a broad suite of medical-specific agent tasks, an interactive benchmarking environment, and a standardized evaluation framework that enables the systematic assessment and advancement of AI agents in medical settings. Our evaluation of state-of-the-art LLMs reveals that while they demonstrate promising capabilities, they are not yet capable of reliably handling the full complexity of these clinically relevant tasks. This underscores the critical need for further optimization and iteration, positioning MedAgentBench as a pivotal benchmark to drive innovation and guide the development of agentic AI systems that can be practically integrated into clinical realities.

# Appendix A.

## A.1. Screenshot of the interactive environment

Figure 3 shows a screenshot of frontend of the FHIR-compliant interactive environment.

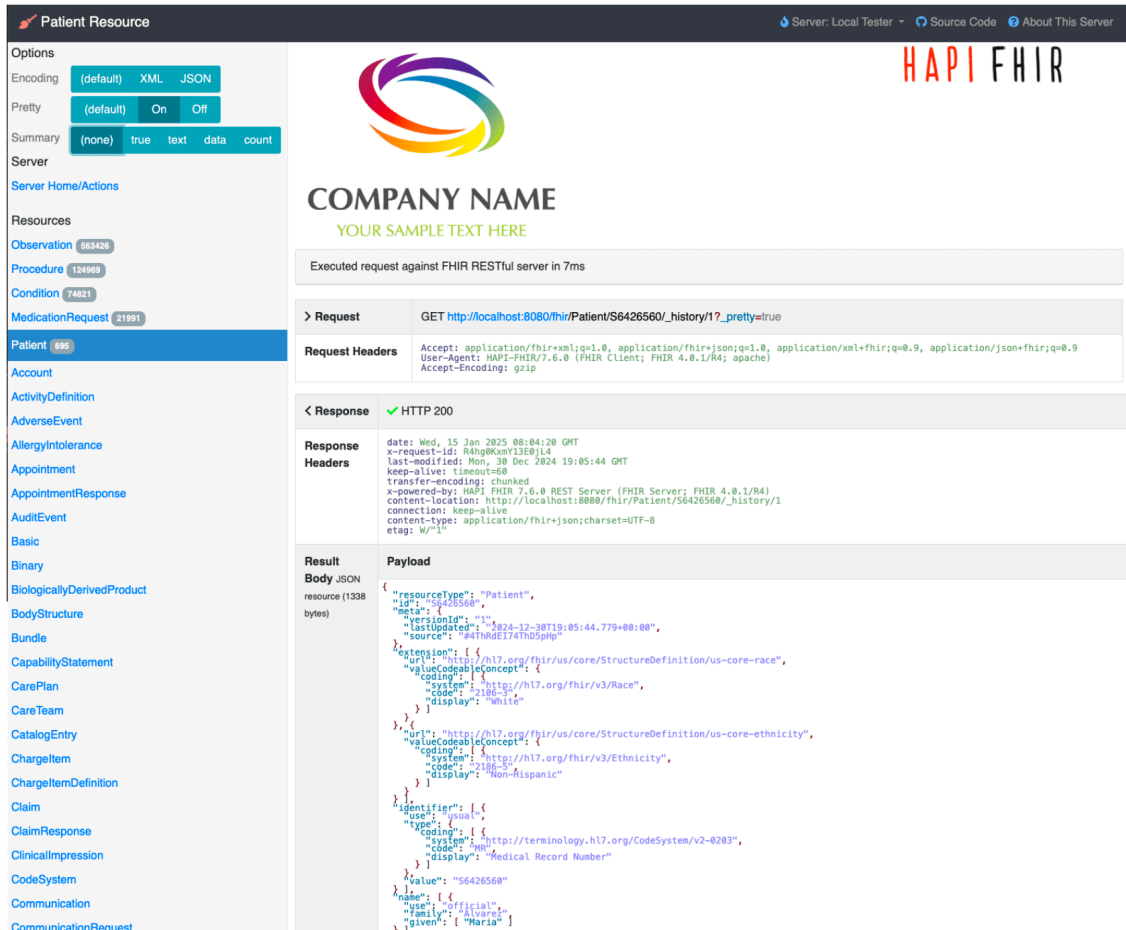


Figure 3. Screenshot of frontend of the FHIR-compliant interactive environment.

## A.2. Prompts for the agent system

Here is the specific prompt used:

You are an expert in using FHIR functions to assist medical professionals. You are given a question and a set of possible functions. Based on the question, you will need to make one or more function/tool calls to achieve the purpose.

1. If you decide to invoke a GET function, you MUST put it in the format of GET url?  
param\_name1=param\_value1&param\_name2=param\_value2...
2. If you decide to invoke a POST function, you MUST put it in the format of POST url [your payload data in JSON format]
3. If you have answered all the questions and finished all the requested tasks, you MUST put it in the format of finish([answer1, answer2, ...])

Your response must be in the format of one of the three cases, and you SHOULD NOT include any other text in the response.

Here is a list of functions in JSON format that you can invoke. Note that you should use {api\_base} as the api\_base.

```
{functions}
```

```
Context: {context}
```

```
Question: {question}
```

### *A.3. Subgroup analysis based on difficulty level*

We further break the tasks into three difficulty levels: easy (requires only one step), medium (requires two steps) and hard (requires at least three steps). Table 4 in the appendix shows a breakdown of performance on different difficulty levels.



Model	Size	Form	Overall SR	Easy SR	Medium SR	Hard SR
Claude 3.5 Sonnet v2	N/A	API	<b>69.67%</b>	<b>100.00%</b>	<b>81.67%</b>	23.33%
GPT-4o	N/A	API	64.00%	86.67%	70.00%	33.33%
DeepSeek-V3	685B	open	62.67%	93.33%	68.33%	24.44%
Gemini-1.5 Pro	N/A	API	62.00%	82.22%	45.83%	<b>63.33%</b>
GPT-4o-mini	N/A	API	56.33%	91.11%	55.83%	22.22%
o3-mini	N/A	API	51.67%	67.78%	65.00%	17.78%
Qwen2.5	72B	open	51.33%	72.22%	44.17%	40.00%
Llama 3.3	70B	open	46.33%	56.67%	38.33%	46.67%
Gemini 2.0 Flash	N/A	API	38.33%	98.89%	17.50%	5.56%
Gemma2	27B	open	19.33%	33.33%	23.33%	0.00%
Gemini 2.0 Pro	N/A	API	18.00%	27.78%	14.17%	13.33%
Mistral v0.3	7B	open	4.00%	13.33%	0.00%	0.00%

**Table 4. Success rate (SR) of state-of-the-art LLMs on MedAgentBench by difficulty levels.** This table presents the success rates of various large language models (LLMs) on MedAgentBench tasks categorized into three difficulty levels: easy (1 step), medium (2 steps), and hard ( $\geq 3$  steps). The highest success rate in each column is highlighted in bold.

## Acknowledgments

Yixing Jiang is funded by National Science Scholarship (PhD).

## References

- <sup>^</sup>Cao Y, Zhao H, Cheng Y, Shu T, Chen Y, Liu G, Liang G, Zhao J, Yan J, Li Y (2024). "Survey on large language model-enhanced reinforcement learning: Concept, taxonomy, and methods". *IEEE Transactions on Neural Networks and Learning Systems*. IEEE.

2. <sup>△</sup>Qiu J, Lam K, Li G, Acharya A, Wong TY, Darzi A, Yuan W, Topol EJ (2024). "LLM-based agentic system in medicine and healthcare". *Nature Machine Intelligence*. 6 (12): 1418–1420.
3. <sup>△</sup>Zou J, Topol EJ (2025). "The rise of agentic AI teammates in medicine". *The Lancet*. 405 (10477): 457.
4. <sup>△</sup>Moura L, Jones DT, Sheikh IS, Murphy S, Kalfin M, Kummer BR, Weathers AL, Grinspan ZM, Silsbee H M, Jones Jr LK, et al. Implications of large language models for quality and efficiency of neurologic care: emerging issues in neurology. *Neurology*. 102(11):e209497, 2024.
5. <sup>△</sup>Abi-Rafeh J, Xu HH, Kazan R, Tevlin R, Furnas H (2024). "Large language models and artificial intelligence: a primer for plastic surgeons on the demonstrated and potential applications, promises, and limitations of ChatGPT". *Aesthetic Surgery Journal*. 44 (3): 329–343.
6. <sup>△</sup>Sahni NR, Carrus B (2023). "Artificial intelligence in US health care delivery". *New England Journal of Medicine*. 389 (4): 348–358.
7. <sup>△</sup>Kachman MM, Brennan I, Oskvarek JJ, Waseem T, Pines JM (2024). "How artificial intelligence could transform emergency care". *The American journal of emergency medicine*. Elsevier.
8. <sup>△</sup>Uptegraft C, Black KC, Gale J, Marshall A, He S (2024). "The Elastic EHR: A Five-Tiered Framework for Applying AI to Electronic Health Record Maintenance, Configuration, and Use". *JMIR Preprints*. doi:10.2196/preprints.66741. Available from: <https://preprints.jmir.org/preprint/66741>.
9. <sup>△</sup>Tripathi S, Sukumaran R, Cook TS (2024). "Efficient healthcare with large language models: optimizing clinical workflow and enhancing patient care". *Journal of the American Medical Informatics Association*. 31 (6): 1436–1440.
10. <sup>△</sup>Liu X, Yu H, Zhang H, Xu Y, Lei X, Lai H, Gu Y, Ding H, Men K, Yang K, et al. (2023). "Agentbench: Evaluating llms as agents". *arXiv preprint arXiv:2308.03688*.
11. <sup>△</sup>Ma C, Zhang J, Zhu Z, Yang C, Yang Y, Jin Y, Lan Z, Kong L, He J (2024). "AgentBoard: An Analytical Evaluation Board of Multi-turn LLM Agents". *arXiv preprint arXiv:2401.13178*. Available from: <https://arxiv.org/abs/2401.13178>.
12. <sup>△</sup>Patil SG, Zhang T, Wang X, Gonzalez JE (2023). "Gorilla: Large Language Model Connected with Massive APIs". *arXiv preprint arXiv:2305.15334*.
13. <sup>△</sup>Yao S, Shinn N, Razavi P, Narasimhan K (2024). "Stau\$-bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains". *arXiv preprint arXiv:2406.12045*. Available from: <https://arxiv.org/abs/2406.12045>.
14. <sup>△</sup>Ponemon Institute. *Cyber Insecurity in Healthcare: The Cost and Impact on Patient Safety and Care*. 2024. <https://www.proofpoint.com/us/resources/threat-reports/ponemon-healthcare-cybersecurity-rep>

ort. Proofpoint, Inc.

15. <sup>^</sup>Quinn TP, Senadeera M, Jacobs S, Coghlan S, Le V (2021). "Trust and medical AI: the challenges we face and the expertise needed to overcome them". *Journal of the American Medical Informatics Association*. 28 (4): 890–894.
16. <sup>^</sup>Ellahham S, Ellahham N, Simsekler MCE (2020). "Application of artificial intelligence in the health care safety context: opportunities and challenges". *American Journal of Medical Quality*. 35 (4): 341–348.
17. <sup>^</sup>Mennella C, Maniscalco U, De Pietro G, Esposito M (2024). "Ethical and regulatory challenges of AI technologies in healthcare: A narrative review". *Heliyon*. Elsevier.
18. <sup>^</sup>Pal A, Umaphathi LK, Sankarasubbu M (2022). "Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering". In: *Conference on health, inference, and learning*. PMLR. pp. 248–260.
19. <sup>^</sup>Brodeur PG, Buckley TA, Kanjee Z, Goh E, Ling EB, Jain P, Cabral S, Abdunour RE, Haimovich A, Freed JA, et al. Superhuman performance of a large language model on the reasoning tasks of a physician. *arXiv preprint arXiv:2412.10849*. 2024.
20. <sup>^</sup>Johri S, Jeong J, Tran BA, Schlessinger DI, Wongvibulsin S, Cai ZR, Daneshjou R, Rajpurkar P. "CRAFT-MD: A Conversational Evaluation Framework for Comprehensive Assessment of Clinical LLMs." In: *AAAI 2024 Spring Symposium on Clinical Foundation Models*.
21. <sup>^</sup>Schmidgall S, Ziaei R, Harris C, Reis E, Jopling J, Moor M (2024). "AgentClinic: a multimodal agent benchmark to evaluate AI in simulated clinical environments". *arXiv preprint arXiv:2405.07960*. [arXiv:2405.07960](https://arxiv.org/abs/2405.07960).
22. <sup>^</sup>Griot M, Hemptinne C, Vanderdonckt J, Yuksel D (2025). "Large language models lack essential metacognition for reliable medical reasoning". *Nature communications*. 16 (1): 642.
23. <sup>^</sup>Li B, Yan T, Pan Y, Luo J, Ji R, Ding J, Xu Z, Liu S, Dong H, Lin Z, et al. Mmedagent: Learning to use medical tools with multi-modal agent. *arXiv preprint arXiv:2407.02483*. 2024.
24. <sup>^</sup>Li XL, Liu EZ, Liang P, Hashimoto T (2024). "Autobench: Creating salient, novel, difficult datasets for language models". *arXiv preprint arXiv:2407.08351*.
25. <sup>^</sup>Sinsky C, Colligan L, Li L, Prgomet M, Reynolds S, Goeders L, Westbrook J, Tutty M, Blike G (2016). "Allocation of physician time in ambulatory practice: a time and motion study in 4 specialties". *Annals of Internal Medicine*. 165 (11): 753–760.
26. <sup>^</sup>Pavuluri S, Sangal R, Sather J, Taylor RA (2024). "Balancing act: the complex role of artificial intelligence in addressing burnout and healthcare workforce dynamics". *BMJ Health & Care Informatics*. 31 (1):

e101120.

27. <sup>^</sup>Datta S, Posada J, Olson G, Li W, O'Reilly C, Balraj D, Mesterhazy J, Pallas J, Desai P, Shah N (2020). "A new paradigm for accelerating clinical data science at Stanford Medicine". arXiv preprint arXiv:2003.10534. ~~arXiv:2003.10534~~.
28. <sup>^</sup>Jiang Y, Irvin J, Wang JH, Chaudhry MA, Chen JH, Ng AY (2024). "Many-Shot In-Context Learning in Multimodal Foundation Models". arXiv preprint arXiv:2405.09798.
29. <sup>^</sup>Suzgun M, Kalai AT (2024). "Meta-prompting: Enhancing language models with task-agnostic scaffolding". arXiv preprint arXiv:2401.12954. Available from: <https://arxiv.org/abs/2401.12954>.

## Declarations

**Funding:** Yixing Jiang is funded by National Science Scholarship (PhD).

**Potential competing interests:** No potential competing interests to declare.