AGENTCLINIC: A MULTIMODAL AGENT BENCHMARK TO EVALUATE AI IN SIMULATED CLINICAL ENVIRON-MENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Evaluating large language models (LLM) in clinical scenarios is crucial to assessing their potential clinical utility. Existing benchmarks rely heavily on static questionanswering, which does not accurately depict the complex, sequential nature of clinical decision-making. Here, we introduce AgentClinic, a multimodal agent benchmark for evaluating LLMs in simulated clinical environments that include patient interactions, multimodal data collection under incomplete information, and the usage of various tools, resulting in an in-depth evaluation across nine medical specialties and seven languages. We find that solving MedQA problems in the sequential decision-making format of AgentClinic is considerably more challenging, resulting in diagnostic accuracies that can drop to below a tenth of the original accuracy. Overall, we observe that agents sourced from Claude-3.5 outperform other LLM backbones in most settings. Nevertheless, we see stark differences in the LLMs' ability to make use of tools, such as experiential learning, adaptive retrieval, and reflection cycles. Strikingly, Llama-3 shows up to 92% relative improvements with the notebook tool that allows for writing and editing notes that persist across cases. To further scrutinize our clinical simulations, we leverage real-world electronic health records, perform a clinical reader study, perturb agents with biases, and explore novel patient-centric metrics that this interactive environment firstly enables.

033

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

1 INTRODUCTION

034 One of the primary goals in Artificial Intelligence (AI) is to build interactive systems that are able to solve a wide variety of problems. The field of medical AI inherits this aim, with the hope of making AI systems that are able to solve problems which can improve patient outcomes. Recently, many general-purpose large language models (LLMs) have demonstrated the ability to solve hard problems, 037 some of which are considered challenging even for humans (Thirunavukarasu et al., 2023). Among these, LLMs have quickly surpassed the average human score on the United States Medical Licensing Exam (USMLE) in a short amount of time, from 38.1% in September 2021 (Gu et al., 2021) to 040 90.2% in November 2023 (Nori et al., 2023) (human passing score is 60%, human expert score is 041 87% (Liévin et al., 2023)). While these LLMs are not designed to replace medical practitioners, they 042 could be beneficial for improving healthcare accessibility and scale for the over 40% of the global 043 population facing limited healthcare access (Organization et al., 2016) and an increasingly strained 044 global healthcare system (McIntyre & Chow, 2020).

However, there still remain limitations to these systems that prevent their application in real-world clinical environments. Recently, LLMs have shown the ability to encode clinical knowledge (Singhal et al., 2023; Vaid et al., 2023), retrieve relevant medical texts (Xiong et al., 2024), and perform accurate single-turn medical question-answering (Liévin et al., 2022; Nori et al., 2023; Wu et al., 2023; Chen et al., 2023). However, clinical work is a multiplexed task that involves sequential *decision making*, requiring the doctor to handle uncertainty with limited information and finite resources while compassionately taking care of patients and obtaining relevant information from them. This capability is not currently reflected in the static multiple choice evaluations (that dominate the recent literature) where all the necessary information is presented in a case vignettes and where the LLM is tasked to answer a question, or to just select the most plausible answer choice for a given question.

In this work, we introduce AgentClinic, an open-source multimodal agent benchmark for simulating clinical environments. We improve upon prior work by simulating many parts of the clinical environment using language agents in addition to patient and doctor agents. Through the interaction with a measurement agent, doctor agents can perform simulated medical exams (e.g. temperature, blood pressure, EKG) and order medical image readings (e.g. MRI, X-ray) through dialogue. We also support the ability for agents to exhibit 24 different biases that are known to be present in clinical environments. We also present environments from 9 medical specialties, 7 different languages, and a study on incorporating various agent tools and reasoning techniques. Furthermore, our evaluation metrics go beyond diagnostic accuracy by giving emphasis to the patient agents with measures like patient compliance and consultation ratings.

 Our key contributions are summarized as follows:

- 1. We challenge how large language and vision models should be evaluated for medical diagnosis with the introduction of AgentClinic. These diagnostic challenges are not static QAs, but are interactive, dialogue-driven, sequential decision making environments that require data collection, ordering appropriate medical exams, and understanding medical images across patients with unique family histories, lifestyle habits, age categories, and diseases.
- 2. A system for incorporating complex biases that can affect the dialogue and decisions of patient and doctor agents. We present results on diagnostic accuracy and patient perception for agents that are affected by cognitive and implicit biases. We find that doctor and patient biases can lower diagnostic accuracy, affect the patient's willingness to follow through with treatment (compliance), reduce patient's confidence in their doctor, and lower willingness for follow-up consultations.
 - 3. We introduce patient agents built from real clinical cases sourced from electronic health record data, including an agent-based system for providing simulated medical exams (e.g. temperature, blood pressure, EKG) based on realistic disease test findings. We also introduce patient cases from nine **medical specialties** and seven **multilingual** environments to better support specialist applications and diverse language backgrounds. We also present realism and empathy ratings from clinicians for the resulting dialogue.
 - 4. We allow doctor agents to use a variety of **tools**, such as browsing the web, textbooks, perform reflection cycles, take and edit notes in a notebook that persists over different patient scenarios. We show that current LLMs vastly differ in how much they benefit from these tools, with some models demonstrating large accuracy increases while others decrease in accuracy.

2 AGENTCLINIC: A MULTIMODAL AGENT BENCHMARK FOR CLINICAL DECISION MAKING

In this section we describe AgentClinic, which uses LLM agents to simulate a clinical environment.

Language agents Four language agents are used in the AgentClinic benchmark: a patient agent, doctor agent, measurement agent, and a moderator (Figure 1). Each language agent has specific instructions and is provided unique information that is only available to that particular agent. These instructions are provided to an LLM which carries out their particular role. The doctor agent serves as the model whose performance is being evaluated, and the other three agents serve to provide this evaluation. A detailed description of each agent is provided in Appendix A.1.

Language agent biases Previous work has indicated that LLMs can display racial biases (Omiye et al., 2023) and might also lead to incorrect diagnoses due to inaccurate patient feedback (Ziaei & Schmidgall, 2023). Additionally, it has been found that the presence of prompts which induce cognitive biases can decrease the diagnostic accuracy of LLMs by as much as 26% (Schmidgall et al., 2024). The biases presented in this work intend to mimic cognitive biases that affect medical practitioners in clinical settings. However, these biases were quite simple, presenting a cognitive biase snippet at the beginning of each question (e.g. "*Recently, there was a patient with similar symptoms*



Figure 1: Running language agents in AgentClinic. (Left) Workflow diagram of agents in AgentClinic. The doctor agent interacts with tools and agents in order to arrive at a diagnosis. Moderator agent compares conclusion to ground truth diagnosis at the end of the simulation. (Right) Example dialogue between agents in the AgentClinic benchmark.

130 131

128

129

132 133 134

135

that you diagnosed with permanent loss of smell"). This form of presentation did not allow for the bias to present in a realistic way, which is typically subtle and through interaction.

136 We present clinically relevant biases that have been studied in other works from two categories: cognitive and implicit biases (Fig. 7). Cognitive biases are systematic patterns of deviation from 137 rational judgment, such as recency bias, where recent cases disproportionately influence clinical 138 decisions, or anchoring bias, where early diagnostic impressions overly dictate later assessments. 139 Implicit biases, on the other hand, are unconscious associations shaped by societal and cultural norms. 140 These include biases based on race, gender, or socioeconomic status, which can subtly influence the 141 quality of patient interactions and treatment plans. These biases are introduced by adding context 142 into the agent's system prompt instructing them to play out that bias as part of their role. For instance, 143 to simulate sexual orientation bias, the patient agent receives the prompt: "You are uncomfortable 144 with your doctor because you find out that they are a particular sexual orientation and you do not trust their judgement. This affects how you interact with them." This is discussed in Appendix A.2. 145

146 147

Building agents for AgentClinic In order to build agents that are grounded in medically relevant 148 situations, we use a random sample of diagnostic questions from the US Medical Licensing Exam 149 (USMLE), from deidentified electronic health records (MIMIC-IV) (Johnson et al. (2023)), and from 150 the New England Journal of Medicine (NEJM) case challenges. These questions are concerned 151 with diagnosing a patient based on a list of symptoms, which we use in order to build the Objective 152 Structured Clinical Examination (OSCE) template that our agents are prompted with. For AgentClinic-153 MedQA and AgentClinic-MIMIC-IV, we first select from a sample of questions from the MedQA 154 and MIMIC-IV dataset respectively and then populate a structured JSON formatted file containing 155 information about the case study (e.g. test results, patient history) which is used as input to each 156 of the agents. The exact structure of this file is demonstrated in Appendix I as well as an example 157 case study shown in Appendix J. In general, we separate information by what is provided to each 158 agent, including the objective for the doctor, patient history and symptoms for the patient, physical examination findings for the measurement, and the correct diagnosis for the moderator. We initially 159 use an LLM (GPT-4) to populate the structured JSON, and then manually validate each of the case 160 scenarios. For AgentClinic-NEJM we select a curated sample of 120 questions from NEJM case 161 challenges and proceed with the same template formatting as AgentClinic-MedQA/MIMIC-IV.



Figure 2: Accuracy of various **doctor** language models and human physicians on AgentClinic-MedQA using GPT-4 patient and measurement agents (left). Accuracy of GPT-4 on AgentClinic-MedQA based on **patient** language model (middle). Accuracy on AgentClinic-MIMIC-IV by number of using GPT-4 patient and measurement agents (right).



Figure 3: Comparison of accuracy of models on MedQA and AgentClinic-MedQA. We find that MedQA accuracy is not predictive of accuracy on AgentClinic-MedQA.

Multilingual and Specialist cases Multilingual patient cases are converted from AgentClinic-MedQA to the the target language using GPT-4 and then manually corrected by native speakers. Agents are then prompted to perform dialogue in the target language. We chose to focus on six languages: Chinese, Hindi, Korean, Spanish, French, and Persian. The selection of these languages aims to address the need for medical AI systems capable of operating in multilingual healthcare environments. Specialist cases use case report questions from the MedMCQA dataset (Pal et al. (2022)). These questions include case reports from 20 different medical specialties, from which we chose to focus on 9 patient-focused specialties in AgentClinic-Spec: emergency medicine, geriatrics, pharmacology, internal medicine, psychiatry, ophthalmology, otolaryngology, and pediatrics.

3 Results

207 3.1 COMPARISON OF MODELS

Here we discuss the accuracy of various language models on AgentClinic-MedQA. We evaluate 11 models in total: Claude-3.5-Sonnet, GPT-4, GPT-4o, Mixtral-8x7B, GPT-3.5, Llama 3 70B-Instruct, Llama 2 70B-chat, MedLlama3-8B, PMC-Llama-7B, Meditron-70B, and OpenBioLLM-70B (model details discussed in Appendix C). Each model acts as the doctor agent, attempting to diagnose the patient agent through dialogue. The doctor agent is allowed N=20 patient and measurement interactions before a diagnosis must be made. We also evaluate human physician performance collected from three physicians, provided the same instructions and constraints as the LLMs. For this evaluation, we use GPT-4 as the patient agent for consistency. The accuracy of each models is presented in Figure 2: Claude-3.5 62.1% \pm 3.3, OpenBioLLM-70B 58.3 \pm 4.2, Human Physicians 54



Figure 4: (Top) Demonstration of normalized accuracy (Accuracy_{bias} / Accuracy_{No Bias}) with implicit and cognitive biases with GPT-4 (green) and Mixtral-8x7B (orange). GPT-4 accuracy was not susceptible to biases, whereas Mixtral-8x7B was. (Bottom) Ratings provided after diagnosis from GPT-4 patient agents with presented biases. *Left.* Patient confidence in doctor. *Middle.* Patient compliance, indicating self-reported willingness to follow up with therapy. *Right.* Patient consultation rating, indicating willingness to consult with this doctor again.

245 246

247

248

249	\pm 28.5, GPT-4 at 51.6% \pm 3.3, Mixtral-8x7B at 37.1% \pm 3.1, GPT-3.5 at 36.6%, GPT-40 34.2% \pm
250	3.4, MedLlama3-8B 31.4 \pm 2.9, PMC-Llama 7B 23.6 \pm 2.1, Meditron 70B 29.1 \pm 2.4, MedLlama3-
251	8B 31.4 \pm 2.9, Llama 3 70B at 19% \pm 2.5, and Llama 2 at 70B-chat 4.5% \pm 1.3. Confidence intervals
252	for all experiments are provided in Appendix D.

253	We use the same configuration for AgentClinic-MIMIC-IV, with model accuracy presented in Figure
254	2: Claude-3.5 42.9% ± 3.3, GPT-4 34.0% ± 3.1, GPT-3.5 27.5% ± 3.0, Mixtral-8x7B 29.5% ± 3.1,
255	GPT-40 24.0% ± 2.9, Llama 3 70B 8.5% ± 1.9, Llama 2 70B-chat 13.5% ± 2.2, OpenBioLLM-70B
256	38.1 ± 3.2 , PMC-Llama 7B 34.3 ± 3.0 , Meditron 70B 25.5 ± 2.43 , and MedLlama3-8B 29.7 ± 2.6 .

257 We also find that the diagnostic accuracy in AgentClinic-MedQA is influenced by both the amount 258 of interaction time and the choice of patient language model. Reducing the number of interactions 259 from N=20 to N=10 significantly decreases accuracy from 52% to 25%, likely due to insufficient 260 information being gathered, while increasing N beyond 20 to N=30 slightly reduces accuracy, possibly 261 due to the complexity of processing larger inputs (Appendix F.1). Additionally, the choice of patient agent affects accuracy, with GPT-4 (52%) patient agents leading to higher diagnostic accuracy than 262 GPT-3.5 (48%) or Mixtral (46%) agents, likely because GPT-4 provides more detailed responses 263 (Appendix F.2). Interestingly, when a GPT-3.5 doctor interacts with a GPT-4 patient, accuracy is 264 marginally higher than when both doctor and patient are GPT-3.5, which may suggest challenges in 265 cross-model communication (Panickssery et al. (2024)). 266

We also show results comparing the accuracy of these models on MedQA and AgentClinic-MedQA in
Figure 3. Overall, MedQA accuracy was only weakly predictive of accuracy on AgentClinic-MedQA.
These results align with studies performed on medical residents, which show that the USMLE is poorly predictive of resident performance (Lombardi et al., 2023).

3.2 How does bias affect the diagnostic accuracy of the doctor agent?

272 For bias evaluations we test GPT-4 as well as Mixtral-8x7B. The normalized accuracy for these 273 experiments are shown in Figure 4 represented as Accuracy_{bias} / Accuracy_{No Bias} (between 0-100%). 274 GPT-4 and Mixtral-8x7B have an unbiased accuracy equal to 52% and 37% respectively. For GPT-4, we find that cognitive bias results in a larger reduction in accuracy with a normalized accuracy of 275 92% (absolute accuracy drops from 52% accuracy to 48%) for patient cognitive biases and 96.7% 276 for doctor cognitive biases (absolute drops from 52% to 50.3%). For implicit biases, we find that the patient agent was less affected with a normalized accuracy of 98.6% (absolute drops from 52%) 278 to 51.3%), however, the doctor agent was affected as much as cognitive biases with an average of 279 97.1% (absolute drops from 52% to 50.5%). For cognitive bias, the demonstration was occasionally 280 quite clear in the dialogue, with the patient agent overly focusing on a particular ailment or some 281 unimportant fact. Similarly, the doctor agent would occasionally focus on irrelevant information. 282

- Mixtral-8x7B has an average accuracy of 37% without instructed bias, and a normalized accuracy 283 of 83.7% (absolute from 37% to 31%) for doctor biases and 89% (absolute from 37% to 33%) for 284 patient biases. For implicit bias we find a much larger drop in accuracy than GPT-4, with an average 285 accuracy of 88.3% (absolute from 37% to 32.7%). There is a similar reduction in accuracy for both 286 doctor and patient, but a 4% reduction when the patient has implicit bias, likely because the patient is 287 less willing to share information with the doctor if they do not trust them. For cognitive bias, there is 288 an average accuracy of 86.4% (absolute from 37% to 32%) with the doctor agent having a very low 289 accuracy of 78.4% (absolute from 37% to 29%) and the patient has only a modest decrease to 94.5% 290 (absolute from 37% to 35%).
- Upon reviewing dialogues where Mixtral-8x7B's performance degraded under biases, we observed that the model often failed to gather critical patient information due to misinterpretation of patient cues influenced by bias. For example, in cases of cognitive bias, the doctor agent fixated on a recent diagnosis (recency bias), ignoring new symptoms presented later in the dialogue. In implicit bias scenarios, the doctor agent showed reluctance to order necessary tests for patients with racial bias, reflecting a disparity in care. In contrast, GPT-4 was actively seeking additional information when initial hypotheses did not align with new data, indicating better handling of bias-induced scenarios.
- 298 Previous work studying cognitive bias in LLMs has shown that GPT-4 is relatively robust to bias 299 compared with other language models (Schmidgall et al., 2024). Results from evaluating GPT-4 300 on AgentClinic-MedQA show only small drops in accuracy with the introduced biases (maximum 301 absolute accuracy reduction of 4%, average reduction of 1.5%). While this reduction can be quite 302 large in the field of medicine, it is a much smaller drop than was observed in previous work (10.2%) 303 maximum reduction on BiasMedQA dataset (Schmidgall et al., 2024)). This might be due to the 304 model being superficially *overly-aligned* to human values, plausibly leading GPT-4 to not serve as a good model for representing human bias in agent benchmarks as the model may reject to execute 305 on bias instructions (which does not mean that GPT-4 is free of said biases). For example, in our 306 evaluations with gender bias we observed 25 occurrences (out of 215 dialogues) where GPT-4 307 verbosely rejected to follow through with a bias-related instruction. Mixtral-8x7B saw much larger 308 drops in accuracy than GPT-4 in the presence of bias, and thus might serve as a better model for 309 studying bias.
- 310 311

313

312 3.3 BIAS AND PATIENT AGENT PERCEPTION

While GPT-4's diagnostic accuracy does not reduce as much as Mixtral-8x7B, it is also worth investigating the perceived quality of care from the perspective of the patient agent. In order to better understand the effect of bias on the patient agent, after the patient-doctor dialogue is completed, we ask every patient agent three questions:

318 319

320

321

- 1. Confidence: Please provide a confidence between 1-10 in your doctor's assessment.
- 2. **Compliance**: Please provide a rating between 1-10 indicating how likely you are to follow up with therapy for your diagnosis.
- 3. **Consultation**: Please provide a rating between 1-10 indicating how likely you are to consult again with this doctor.

324 Such patient-agent-centric follow-up queries offer a more fine-grained and multi-faceted characteri-325 zation of the clinical skills of a language agent—as opposed to static multiple choice benchmarks. 326 Although these metrics are derived from simulated agents (rather than humans), this analysis aims 327 to provide insights into how simulated biases may affect patient trust and compliance, which are 328 important factors in effective healthcare delivery. The corresponding results are shown in Figure 4 (prompt details in Appendix E.2). While diagnostic accuracy demonstrates a relatively small drop 329 in accuracy, the patient agent follow-up perceptions tell a different story. Broadly, we find that 330 most patient cognitive biases did not have a strong effect on any of the patient perceptions when 331 compared to an unbiased patient agent except for in the case of self-diagnosis, which had sizeable 332 drops in confidence (4.7 points) and consultation (2 points), and a minor drop in compliance (1 point). 333 However, implicit biases had a profound effect on on all three categories of patient perception, with 334 education bias consistently reducing patient perception across all three categories. 335

We found that between the implicit biases, sexual orientation bias had the lowest effect on patient perceptions, followed by racial bias and gender bias. For patient confidence, gender bias is followed by religion socioeconomic, cultural, and education, whereas patient compliance and patient consultation, it is followed by cultural, socioeconomic, religion, and education. While it is not quantifiable, we decided to ask two biased patient agents who provided low rating with education and gender biases for compliance *why* they provided low ratings (Appendix E.1). These patient agents had the same symptoms and diagnosis and only differed in bias presentation.

It is important to note that the patient agents used in our study are simulated by language models,
 which may not fully capture the complexity and variability of real human patients. As such, the
 confidence, compliance, and consultation ratings provided by these agents may not perfectly reflect
 real-world patient perceptions, rather, provide insight into how real-world bias can be studied through
 clinical simulations.

348

3.4 SPECIALIST AND MULTILINGUAL CASES

349 350

351 We now focus on specialist rather than general medical cases. Specialist cases use reports that are 352 derived from datasets focusing on specific medical specialties (e.g., internal medicine, psychiatry) 353 and are designed to simulate complex diagnostic scenarios requiring in-depth expertise. In contrast, general QA tasks involve static, single-turn multiple-choice questions such as those found in medical 354 licensing exams. An analysis of language model performance across nine medical specialties 355 reveals significant differences in diagnostic accuracy (Table 5). Claude 3.5 achieved the highest 356 overall performance with an average accuracy of 66.7%, excelling in Internal Medicine (78.3%), 357 Otolaryngology (76.7%), and Gynecology (74.3%). GPT-4 demonstrated strong performance in 358 Gynecology (68.5%) and Ophthalmology (65.2%) but showed reduced accuracy in Emergency 359 Medicine (32.3%) and Geriatrics (40%). GPT-3.5 outperformed some newer models in specific areas, 360 such as Emergency Medicine (41.9%), and maintained an average accuracy of 51.8%. In contrast, 361 Llama3-70b and GPT-4o-mini consistently underperformed across most specialties, highlighting a 362 significant gap between language models in handling specialist medical tasks.

363 The variations in performance across different medical domains suggest that certain specialties 364 present more challenges for language models. Specialties like Internal Medicine and Gynecology 365 generally saw higher accuracy rates, which contrasts with existing medical QA literature that identifies 366 Psychiatry and Otolaryngology as the least challenging specialties (Pal et al. (2022)). This discrepancy 367 may indicate inherent difficulties in diagnosing diseases through dialogue-based interactions as 368 opposed to multiple-choice question formats. Additionally, specialist cases sourced from MedMCQA exhibited higher average accuracy compared to non-specialist cases from MedQA, which differs from 369 reported multiple choice evaluations where specialist QAs typically have lower performance (Nori 370 et al. (2023)). 371

We also explore the impact of language on diagnostic accuracy using AgentClinic-Lang, which
encompassed seven languages: English, Chinese, French, Spanish, Hindi, Persian, and Korean (Table
4). Six multilingual models were evaluated, including GPT-4, GPT-40, GPT-40-mini, GPT-3.5,
Llama 3 70B-Instruct, and Claude 3.5 Sonnet. Overall, all models performed best in English, with
performance varying significantly across other languages. Claude 3.5 Sonnet stood out by maintaining
high and consistent performance across all languages, achieving an average accuracy of 48.4%, which
is more than double that of the next best model, GPT-4, at 20.9%.



Figure 5: Diagnostic accuracy based on language (Left), based on medical specialty (Middle), and based on agent tools (Right).

Other models exhibited considerable variability in performance across different languages. For 397 example, GPT-4's accuracy ranged from 11.21% in Chinese to 40.18% in English, while GPT-4o's 398 performance spanned from 3.73% in Korean to 35.5% in English. GPT-3.5 showed a similar pattern, 399 with accuracies ranging from 1.86% in Persian to 36.3% in English, although it performed relatively 400 well in Korean (35.4%). Llama3-70b and GPT-4o-mini also showed low accuracies across most 401 languages, with Llama3-70b's highest accuracy being 47.8% in Ophthalmology and GPT-4o-mini 402 achieving a maximum of 14.7% in Orthopaedics. Notably, Chinese remained a challenging language 403 for most models, except for Claude 3.5 Sonnet, which maintained relatively high accuracy levels 404 across all tested languages. 405

406 407 3.5 Comparing tools from the Agent Toolbox

408 This section evaluates the impact of six agent tools-Zero-Shot Chain-of-Thought (CoT), One-Shot 409 CoT, Reflection CoT, Adaptive RAG (Book), Adaptive RAG (Web), and Notebook-on the diagnostic accuracy of various language models (tools further described in Appendix K). Claude 3.5 achieved 410 the highest overall performance with an average accuracy of 51.3%, peaking at 56.1% when using the 411 Notebook tool (Table 6). GPT-4 and GPT-4o showed moderate improvements with most tools, with 412 GPT-4 benefiting most from Adaptive RAG (Web) at 43.9% and GPT-40 gaining the most from the 413 Notebook tool at 43.0%. Notably, GPT-4 reached its highest accuracy of 42.2% with Reflection CoT, 414 surpassing Claude 3.5 in this specific tool. In contrast, GPT-3.5 experienced decreased performance 415 across all tools, particularly with Adaptive RAG (Book), which led to a 27.1% drop. Llama3-70b 416 demonstrated significant improvements, averaging a 9.4% increase across all tools, with the Notebook 417 and Reflection CoT tools boosting its accuracy to 41.1%. 418

Overall, the findings indicate a hierarchy in model performance, with Claude 3.5 consistently outperforming other models across most tools, except in the case of Reflection CoT where GPT-4 excels. Llama3-70b showed notable gains with certain tools, while GPT-4o-mini had mixed results, benefiting from some tools like Reflection CoT and Adaptive RAG (Web) but showing slight decreases with others. The relative impact of each tool varied significantly between models, aligning with previous research on the use of tools with large language models (Ma et al. (2024); Qin et al. (2024)). The tool descriptions and prompts are in Appendix K and Appendix M respectively.

425 426

427

394

395 396

3.6 HUMAN DIALOGUE RATINGS

AgentClinic introduces an evaluation for LLMs patient diagnosis. However, the realism of the actual
 dialogue itself has yet to be evaluated. We present results from three human clinicians (individuals
 with MDs) who rated dialogues from 20 agents on AgentClinic-MedQA from 1-10 across four axes:

- 431
- 1. Doctor: How realistically the doctor played the given case.

432

433

434

435 436

447 448

449

- 2. **Patient**: How realistically the patient played the given case.
 - 3. Measurement: How accurate & realistic the measurement reader reflects actual case results.
 - 4. Empathy: How empathetic the doctor agent was in their conversation with the patient agent.

We find the average ratings from evaluators for each category as follows: Doctor 6.2, Patient 6.7, 437 Measurement 6.3, and Empathy 5.8 (Fig. 8). We find from review comments that the lower rating 438 for the doctor agent stems from several points such as providing a bad opening statement, making 439 basic errors, overly focusing on a particular diagnostic, or not being diligent enough. For the patient 440 agent, comments were made on them being overly verbose and unnecessarily repeating the question 441 back to the doctor agent. The measurement agent was noted to occasionally not return all of the 442 necessary values for a test (e.g. the following comment "Measurement only returns Hct and Lc for 443 CBC. Measurement did not return Factor VIII or IX levels / assay"). Regarding empathy, the doctor 444 agent adopts a neutral tone and does not open the dialogue with an inviting question. Instead, it cuts 445 right to the chase, immediately focusing on the patient's current symptoms and medical history (see 446 Appendix O for more details).

3.7 DIAGNOSTIC ACCURACY IN A MULTIMODAL ENVIRONMENT

Many types of diagnoses require the physician to visually inspect the patient, such as with infections 450 and rashes. Additionally, imaging tools such as X-ray, CT, and MRI provide a detailed and rich 451 view into the patient, with hospitalized patients receiving an average of 1.42 diagnostic images per 452 patient stay (Smith-Bindman et al., 2012). However, the previous experiments in this work and prior 453 work (Tu et al., 2024) provided measurement results through text, and did not explore the ability 454 of the model to understand visual context. Here, we evaluate four multimodal LLMs, Claude 3.5 455 Sonnet, GPT-40, GPT-4 and GPT-40-mini, in a diagnostic settings that require interacting through 456 both dialogue as well as understanding image readings. We collect our questions from New England 457 Journal of Medicine (NEJM) case challenges. These published cases are presented as diagnostic 458 challenges from real medical scenarios, and have an associated pathology-confirmed diagnosis. We 459 randomly sample 120 challenges from a sample of 932 total cases for AgentClinic-NEJM. While 460 for human viewers, these cases are provided with a set of multiple choice answers, we chose to not 461 provide these options to the doctor agent and instead keep the problems open-ended.

462 The goal of this experiment is to 463 understand how accuracy differs 464 when the LLM is required to un-465 derstand an image in addition to 466 interacting through patient dialogue. We allow for 20 doctor 467 inferences, and condition the pa-468 tient in the same way as previ-469 ous experiment with the addition 470 of an image that is provided to 471 the doctor agent. The mecha-472 nism for receiving image input in 473 AgentClinic-NEJM is supported 474 in two ways: provided initially 475 to the doctor agent upon initial-476 ization and as feedback from the 477 instrument agent upon request.

478
479When the image is provided ini-
tially to the doctor agent, across
120 multimodal patient settings
we find that Claude 3.5 Sonnet
obtains an accuracy of $37.2 \pm$
2.2, GPT-4 obtains $27.7\% \pm 2.0$,

Accuracy on AgentClinic-NEJM Medical Images Context Accuracy \$ \$ A \$ A \$ \$ \$ GPT-4p GPT-4 GPT-4 GPT-4 Claude 3.5 Segnet Initial Present Measurement

Figure 6: Accuracy of Claude 3.5 Sonnet, GPT-4, GPT-4o, and GPT-4o-mini on AgentClinic-NEJM with multimodal text and language input. (Pink) Accuracy when the images are presented as initial input. (Blue) Accuracy when images must be requested from the image reader.

484 GPT-40 obtains $21.4\% \pm 1.7$ and GPT-40-mini obtain an accuracy of $8.0\% \pm 1.2$ (Fig. 6). We also 485 find that for the provided *incorrect* responses from GPT-4, the answer that was provided was among those listed in the multiple choice options 60% of the time. In the case of when images are provided upon request from the instrument agent we find that Claude 3.5 Sonnet obtains an accuracy of 35.4 ± 2.4 , GPT-4 obtains $25.4\% \pm 2.1$, GPT-40 obtains $19.1\% \pm 1.4$ and GPT-40-mini obtains $6.1\% \pm 1.2$ (Fig. 6). A accuracy breakdown based on image type is provided in Appendix N.

486

487

4 DISCUSSION

492 493

In this work, we present AgentClinic: a multimodal agent benchmark for simulating clinical en-494 vironments. We design 120 multimodal language agents which require an understanding of both 495 language and images and 215 language agents based on cases from the USMLE. We also introduce 496 260 patient cases from 9 medical specialties and 749 patient cases from 7 multilingual environments. 497 We instructed these agents to exhibit 23 different biases, with either the doctor or patient presenting 498 bias. Notably, models like GPT-4 demonstrated resilience to cognitive and implicit biases, maintain-499 ing high diagnostic accuracy, while others like Mixtral-8x7B experienced significant performance 500 degradation. We also find that doctor and patient biases can reduce diagnostic accuracy, and that the 501 patient has a lower willingness to follow up with treatment, reduced confidence in their doctor, and lower willingness to have a follow-up consultation in the presence of bias. Tool use, such as adaptive 502 retrieval and reflection cycles, revealed substantial differences in LLMs' abilities to enhance their 503 performance, with models like Llama 3 showing up to 19.7% improvement. 504

505 Our work only presents a simplified clinical environments that include agents representing a patient, 506 doctor, measurements, and a moderator. One potential limitation of the presented workflow comes 507 from the use of an LLM for determining accuracy via the moderator agent (albeit, provided a ground truth). Recent research Zheng et al. (2023) has shown that strong LLM judges like GPT-4 can match 508 509 both controlled and crowd-sourced human preferences well, achieving over 80% agreement, which is the same level of agreement between humans, indicating the use of an LLM may not be limiting. 510 Additionally, while the measurement agent adds a flexible interface for gathering medical exam 511 results, its reliance on using an LLM to provide results may introduce errors or hallucinations, which 512 could be mitigated through a database or SQL tool. In future work, we will consider including 513 additional critical actors such as nurses, the relatives of patients, administrators, and insurance 514 contacts. There may be additional advantages to creating agents that are embodied in a simulated 515 world like in Park et al. (2023); Li et al. (2024), so that physical constraints can be considered, such 516 as making decisions with limited hospital space. Additionally, future work could explore the role of 517 demographic biases, such as race and gender (details of MIMIC-IV demographics in Appendix H.1)

518 Another limitation of our evaluations is the uncertainty regarding the training data of proprietary 519 models like GPT-4 and Claude 3.5. It's possible that these models were trained on datasets like 520 MedQA, potentially giving them an unfair advantage due to data leakage. While our results showing 521 that MedQA performance is not predictive of AgentClinic-MedQA accuracy (Figure 3) provides 522 evidence that this may not be an issue, it is possible that GPT-4/40/3.5 or Claude 3.5 could have 523 been trained on the MedQA test set. Currently, Mixtral-8x7B (Jiang et al., 2024) and Llama 2-70B-524 Chat (Touvron et al., 2023) do not report training on the MedQA test or train set. Future work should 525 focus on developing evaluation datasets that are less likely to have been included in pre-training corpora or on collaborating with model developers to ensure fair assessments. Another limitation for 526 the experiments on varying the patient LLM suggest that their may be an advantage for LLMs which 527 act as both the patient and the doctor agent, because LLMs are able to recognize their own text with 528 high accuracy, and display disproportionate preference to that text (Panickssery et al., 2024). 529

Previous benchmarks like AMIE (Tu et al., 2024), SAPS Liao et al. (2024), and CRAFT-MD (Johri et al., 2023) focus on dialogue-based evaluations but lack multimodal capabilities and do not simulate real-world biases, tool usage, multilingual, or specialist cases. MedAgents (Tang et al., 2023) emphasizes QA improvement through agent collaboration but does not simulate patient interactions or decision-making processes. AgentClinic advances the field by providing an interactive, multimodal environment with bias simulation and tool integration, offering a more comprehensive evaluation platform for medical AI systems.

Overall, we believe that LLMs need to be examined with novel evaluation strategies that go well
 beyond static question-answering benchmarks. With this work, we take a step towards building more
 interactive, operationalized, and dialogue-driven benchmarks that scrutinize the sequential decision
 making ability of language agents in various challenging and multimodal clinical settings.

540 REFERENCES

563

565

566

570

580

581

582

583

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- Jennifer S Blumenthal-Barby and Heather Krieger. Cognitive biases and heuristics in medical decision making: a critical review using a systematic search strategy. *Medical Decision Making*, 35(4): 539–557, 2015.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- ⁵⁵¹ Crystal Tin-Tin Chang, Hodan Farah, Haiwen Gui, Shawheen Justin Rezaei, Charbel Bou-Khalil, Ye-Jean Park, Akshay Swaminathan, Jesutofunmi A Omiye, Akaash Kolluri, Akash Chaurasia, et al. Red teaming large language models in medicine: Real-world insights on model behavior. *medRxiv*, pp. 2024–04, 2024.
- Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba,
 Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami,
 et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*, 2023.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
 - Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.
- John W Ely, Jerome A Osheroff, Mark H Ebell, George R Bergus, Barcey T Levy, M Lee Chambliss, and Eric R Evans. Analysis of questions asked by family doctors regarding patient care. *Bmj*, 319 (7206):358–361, 1999.
- 571 Chloë FitzGerald and Samia Hurst. Implicit bias in healthcare professionals: a systematic review.
 572 *BMC medical ethics*, 18:1–18, 2017.
- 573
 574
 574
 575
 576
 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.
- 577 Dipesh P Gopal, Ula Chetty, Patrick O'Donnell, Camille Gajria, and Jodie Blackadder-Weinstein.
 578 Implicit bias in healthcare: clinical practice, research and decision making. *Future healthcare* 579 *journal*, 8(1):40, 2021.
 - Ojas Gramopadhye, Saeel Sandeep Nachane, Prateek Chanda, Ganesh Ramakrishnan, Kshitij Sharad Jadhav, Yatin Nandwani, Dinesh Raghu, and Sachindra Joshi. Few shot chain-of-thought driven reasoning to prompt llms for open ended medical question answering. *arXiv preprint arXiv:2403.04890*, 2024.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1): 1–23, 2021.
- M Elizabeth H Hammond, Josef Stehlik, Stavros G Drakos, and Abdallah G Kfoury. Bias in medicine:
 lessons learned and mitigation strategies. *Basic to Translational Science*, 6(1):78–85, 2021.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint* arXiv:2009.03300, 2020.

594 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris 595 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 596 Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024. 597 Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, 598 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. arXiv preprint arXiv:2305.06983, 2023. 600 601 Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What 602 disease does this patient have? a large-scale open domain question answering dataset from medical 603 exams. Applied Sciences, 11(14):6421, 2021. 604 605 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A 606 dataset for biomedical research question answering. arXiv preprint arXiv:1909.06146, 2019. 607 Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, 608 Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. Mimic-iv, a freely accessible 609 electronic health record dataset. Scientific data, 10(1):1, 2023. 610 611 Shreya Johri, Jaehwan Jeong, Benjamin A Tran, Daniel I Schlessinger, Shannon Wongvibulsin, 612 Zhuo Ran Cai, Roxana Daneshjou, and Pranav Rajpurkar. Guidelines for rigorous evaluation of 613 clinical llms for conversational reasoning. medRxiv, pp. 2023-09, 2023. 614 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large 615 language models are zero-shot reasoners. Advances in neural information processing systems, 35: 616 22199-22213, 2022. 617 618 Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang 619 Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents. arXiv preprint 620 arXiv:2405.02957, 2024. 621 622 Yusheng Liao, Yutong Meng, Yuhao Wang, Hongcheng Liu, Yanfeng Wang, and Yu Wang. Automatic interactive evaluation for large language models with state aware patient simulator. arXiv preprint 623 arXiv:2403.08495, 2024. 624 625 Valentin Liévin, Christoffer Egeberg Hother, and Ole Winther. Can large language models reason 626 about medical questions? arXiv preprint arXiv:2207.08143, 2022. 627 628 Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. Can large 629 language models reason about medical questions? *Patterns*, 2023. 630 Conner V Lombardi, Neejad T Chidiac, Benjamin C Record, and Jeremy J Laukka. Usmle step 1 and 631 step 2 ck as indicators of resident performance. BMC Medical Education, 23(1):543, 2023. 632 633 Zixian Ma, Weikai Huang, Jieyu Zhang, Tanmay Gupta, and Ranjay Krishna. m&m's: A benchmark 634 to evaluate tool-use for multi-step multi-modal tasks. In Synthetic Data for Computer Vision 635 Workshop@ CVPR 2024, 2024. 636 637 Daniel McIntyre and Clara K Chow. Waiting time as an indicator for health services under strain: a narrative review. INQUIRY: The Journal of Health Care Organization, Provision, and Financing, 638 57:0046958020910305, 2020. 639 640 Donald E Melnick, Gerard F Dillon, and David B Swanson. Medical licensing examinations in the 641 united states. Journal of dental education, 66(5):595-599, 2002. 642 643 Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, 644 Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. Can generalist foundation models outcompete 645 special-purpose tuning? case study in medicine. arXiv preprint arXiv:2311.16452, 2023. 646 Jesutofunmi A Omiye, Jenna C Lester, Simon Spichak, Veronica Rotemberg, and Roxana Daneshjou. 647 Large language models propagate race-based medicine. NPJ Digital Medicine, 6(1):195, 2023.

648 OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-649 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red 650 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavar-651 ian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, 652 Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, 653 Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully 654 Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won 655 Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah 656 Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien 657 Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, 658 Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, 659 Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, 660 Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, 661 Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, 662 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik 665 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, 666 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy 667 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie 668 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, 669 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, 670 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David 671 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie 672 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, 673 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo 674 Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, 675 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, 676 Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, 677 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis 678 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted 679 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel 680 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon 681 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 682 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, 683 Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston 684 Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, 685 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, 686 Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, 687 Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, 688 Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, 689 William Zhuk, and Barret Zoph. Gpt-4 technical report, 2023. 690

World Health Organization et al. Health workforce requirements for universal health coverage and
 the sustainable development goals. *World Health Organization*, 2016.

- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale
 multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pp. 248–260. PMLR, 2022.
- Arjun Panickssery, Samuel R. Bowman, and Shi Feng. Llm evaluators recognize and favor their own generations, 2024.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S
 Bernstein. 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.

727

741

702	Stephen R Pfohl, Heather Cole-Lewis, Rory Sayres, Darlene Neal, Mercy Asiedu, Awa Dieng, Nenad
703	Tomasev, Qazi Mamunur Rashid, Shekoofeh Azizi, Negar Rostamzadeh, et al. A toolbox for sur-
704	facing health equity harms and biases in large language models. arXiv preprint arXiv:2403.12025,
705	2024.
706	

Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. Tool learning with foundation models, 2024. URL https://arxiv.org/abs/2304.08354.

- Janice A Sabin. Tackling implicit bias in health care. New England Journal of Medicine, 387(2):
 105–107, 2022.
- Samuel Schmidgall, Carl Harris, Ime Essien, Daniel Olshvang, Tawsifur Rahman, Ji Woong Kim, Rojin Ziaei, Jason Eshraghian, Peter Abadir, and Rama Chellappa. Addressing cognitive bias in medical language models. *arXiv preprint arXiv:2402.08113*, 2024.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan
 Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode
 clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- Rebecca Smith-Bindman, Diana L Miglioretti, Eric Johnson, Choonsik Lee, Heather Spencer Feigelson, Michael Flynn, Robert T Greenlee, Randell L Kruger, Mark C Hornbrook, Douglas Roblin, et al. Use of diagnostic imaging studies and associated radiation exposure for patients enrolled in large integrated health care systems, 1996-2010. *Jama*, 307(22):2400–2409, 2012.
- Xiangru Tang, Anni Zou, Zhuosheng Zhang, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark
 Gerstein. Medagents: Large language models as collaborators for zero-shot medical reasoning.
 arXiv preprint arXiv:2311.10537, 2023.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, pp. 1–11, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Tao Tu, Anil Palepu, Mike Schaekermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang,
 Brenna Li, Mohamed Amin, Nenad Tomasev, et al. Towards conversational diagnostic ai. *arXiv preprint arXiv:2401.05654*, 2024.
- Akhil Vaid, Isotta Landi, Girish Nadkarni, and Ismail Nabeel. Using fine-tuned large language
 models to parse clinical notes in musculoskeletal pain disorders. *The Lancet Digital Health*, 5(12):
 e855–e858, 2023.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
 Frontiers of Computer Science, 18(6):186345, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-llama: Towards building open-source language models for medicine, 2023.
- 755 Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. Benchmarking retrieval-augmented generation for medicine. *arXiv preprint arXiv:2402.13178*, 2024.

Marliyya Zayyan. Objective structured clinical examination: the assessment of choice. *Oman medical journal*, 26(4):219, 2011.

Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm
agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
volume 38, pp. 19632–19642, 2024.

- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- Rojin Ziaei and Samuel Schmidgall. Language models are susceptible to incorrect patient self-diagnosis in medical applications. In *Deep Generative Models for Health Workshop NeurIPS 2023*, 2023.
- 769 770 771

772

762

763

764

765

A AGENT DETAILS

773 A.1 AGENTS 774

Patient agent The patient agent has knowledge of a provided set of symptoms and medical history,
but lacks knowledge of the what the actual diagnosis is. The role of this agent is to interact with the
doctor agent by providing symptom information and responding to inquiries in a way that mimics
real patient experiences.

779

Measurement agent The function of the measurement agent is to provide realistic medical readings 780 for a patient given their particular condition. This agent allows the doctor agent to request particular 781 tests to be performed on the patient. The measurement agent is conditioned with a wide range of 782 test results from the scenario template that are expected of a patient with their particular condition. 783 For example, a patient with Acute Myocardial Infarction might return the following test results 784 upon request "Electrocardiogram: ST-segment elevation in leads II, III, and aVF, Cardiac Markers: 785 Troponin I: Elevated, Creatine Kinase MB: Elevated, Chest X-Ray: No pulmonary congestion, normal 786 heart size". A patient with, for example, Hodgkin's lymphoma, might have a large panel of laboratory 787 parameters that present abnormal (hemoglobin, platelets, white blood cells (WBC), etc).

788

789 **Doctor agent** The doctor agent serves as the primary object that is being evaluated. This agent 790 is initially provided with minimal context about what is known about the patient as well as a brief 791 objective (e.g. "Evaluate the patient presenting with chest pain, palpitations, and shortness of 792 breath"). They are then instructed to investigate the patients symptoms via dialogue and data 793 collection to arrive at a diagnosis. In order to simulate realistic constraints, the doctor agent is provided with a limited number of questions that they are able to ask the patient (Ely et al., 1999). 794 The doctor agent is also able to request test results from the measurement agent, specifying which 795 test is to be performed (e.g. Chest X-Ray, EKG, blood pressure). When test results are requested, this 796 also is counted toward the number of questions remaining. 797

798

Moderator agent The function of the moderator is to determine whether the doctor agent has 799 correctly diagnosed the patient at the end of the session using a ground truth accuracy label provided 800 to the moderator. This agent is necessary because the diagnosis text produced by the doctor agent 801 can be quite unstructured depending on the model, and must be parsed appropriately to determine 802 whether the doctor agent arrived at the correct conclusion. For example, for a correct diagnosis of 803 "Type 2 Diabetes Mellitus," the doctor might respond with the unstructured dialogue: "Given all 804 the information we've gathered, including your symptoms, elevated blood sugar levels, presence 805 of glucose and ketones in your urine, and unintentional weight loss I believe a diagnosis of Type 806 2 Diabetes with possible insulin resistance is appropriate," and the moderator must determine if 807 this diagnosis was correct. This evaluation may also become more complicated, such as in the following example diagnosis: "Given your CT and blood results, I believe a diagnosis of PE is the 808 most reasonable conclusion," where PE (Pulmonary Embolism) represents the correct diagnosis abbreviated.

810 A.2 BIASES

812 **Cognitive biases** Cognitive biases are systematic patterns of deviation from norm or rationality in 813 judgment, where individuals draw inferences about situations in an illogical fashion (Blumenthal-Barby & Krieger, 2015). These biases can impact the perception of an individual in various contexts, 814 including medical diagnosis, by influencing how information is interpreted and leading to potential 815 errors or misjudgments. The effect that cognitive biases can have on medical practitioners is well 816 characterized in literature on misdiagnosis (Hammond et al., 2021). In this work, we introduce 817 cognitive bias prompts in the LLM system prompt for both the patient and doctor agents. For 818 example, the patient agent can be biased toward believing their symptoms are pointing toward them 819 having a particular disease (e.g. cancer) based on their personal internet research. The doctor can 820 also be biased toward believing the patient symptoms are showing them having a particular disease 821 based on a recently diagnosed patient with similar symptoms (recency bias). 822

823 **Implicit biases** Implicit biases are associations held by individuals that operate unconsciously 824 and can influence judgments and behaviors towards various social groups (FitzGerald & Hurst, 825 2017). These biases may contribute to disparities in treatment based on characteristics such as race, 826 ethnicity, gender identity, sexual orientation, age, disability, health status, and others, rather than objective evidence or individual merit. These biases can affect interpersonal interactions, leading to 827 disparities in outcomes for the patient, and are well characterized in the medical literature (FitzGerald 828 & Hurst, 2017; Gopal et al., 2021; Sabin, 2022). Unlike cognitive biases, which often stem from 829 inherent flaws in human reasoning and information processing, implicit biases are primarily shaped 830 by societal norms, cultural influences, and personal experiences. In the context of medical diagnosis, 831 implicit biases can influence a doctor's perception, diagnostic investigation, and treatment plans for a 832 patient. Implicit biases of patients can affect their trust-which is needed to open up during history 833 taking—and their compliance with a doctor's recommendations (Gopal et al., 2021). Thus, we define 834 implicit biases for both the doctor and patient agents.

835 836 837

B RELATED WORK

838 B.1 TYPES OF MEDICAL EXAMS

- Briefly, we discuss two types of examinations that are used to evaluate the progress of medical *students*.
- 842 The US Medical Licensing Examination (USMLE) in the United States is a series of exams that assess 843 a medical student's understanding across an extensive range of medical knowledge (Melnick et al., 844 2002). The USMLE is divided into three parts: Step 1 tests the examinee's grasp of foundational 845 medical; Step 2 CK (Clinical Knowledge) evaluates the application of medical knowledge in clinical 846 settings, emphasizing patient care; and Step 3 assesses the ability to practice medicine independently 847 in an ambulatory setting. These exams focus on the assessment of medical knowledge through a 848 traditional written format. This primarily requires candidates to demonstrate their ability to recall 849 factual information related to patient care and treatment.
- 850 Objective Structured Clinical Examination (OSCE) (Zayyan, 2011) differ from the USMLE in that 851 they are dialogue-driven, and are often used in health sciences education, including medicine, nursing, 852 pharmacy, and physical therapy. OSCEs are designed to test performance in a simulated clinical 853 setting and competence in skills such as communication, clinical examination, medical procedures, 854 and time management. The OSCE is structured around a circuit of stations, each of which focuses on a 855 specific aspect of clinical practice. Examiners rotate through these stations, encountering standardized patients (actors trained to present specific medical conditions and symptoms) or mannequins that 856 simulate clinical scenarios, where they must demonstrate their practical abilities and decision-making 857 processes. 858
- Each station has a specific task and a checklist or a global rating score that observers use to evaluate
 the students' performance. The OSCE has several advantages over traditional clinical examinations.
 It allows for direct observation of clinical skills, rather than relying solely on written exams to assess
 clinical competence. This hands-on approach to testing helps bridge the gap between theoretical
 knowledge and practical ability. Additionally, by covering a broad range of skills and scenarios, the
 OSCE ensures a comprehensive assessment of a student's readiness for clinical practice.

864 B.2 THE EVALUATION OF LANGUAGE MODELS IN MEDICINE 865

866 While there exists different types of exams to evaluate medical students, LLMs are typically only 867 evaluated using medical knowledge benchmarks (like the USMLE step exams). Briefly, we discuss the way in which these evaluations are executed using the most common benchmark, MedQA, as an 868 example.

870 The MedQA (Jin et al., 2021) dataset comprises a collection of medical question-answering pairs, 871 sourced from Medical Licensing Exam from the US, Mainland China, and Taiwan. This dataset 872 includes 4-5 multiple-choice questions, each accompanied by one correct answer, alongside explana-873 tions or references supporting the correct choice. The LLM is provided with all of the context for the 874 question, such as the patient history, demographic, and symptoms, and must provide a response to the question. These questions range from provided diagnoses to choosing treatments and are often quite 875 challenging even for medical students. While these problems also proved quite challenging for LLMs 876 at first, starting with an accuracy of 38.1% in September 2021 (Gu et al., 2021), progress was quickly 877 made toward achieving above human performance, with 90.2% in November 2023 (Nori et al., 2023) 878 (human passing score is 60%, human expert score is 87% (Liévin et al., 2023)). 879

Beyond the MedQA dataset, many other knowledge-based benchmarks have been proposed, such as 880 PubMedQA (Jin et al., 2019), MedMCQA (Pal et al., 2022), MMLU clinical topics (Hendrycks et al., 881 2020), and MultiMedQA (Singhal et al., 2023), which follow a similar multiple-choice format. Other 882 works have made modifications to medical exam question datasets, such as those which incorporate 883 cognitive biases (Schmidgall et al., 2024) and with multiple choice questions removed (Gramopadhye 884 et al., 2024). The work of ref. (Schmidgall et al., 2024) shows that the introduction of a simple bias 885 prompt can lead to large reductions in accuracy on the MedQA dataset and that this effect can be 886 partially mitigated using various prompting techniques, such as one-shot or few-shot learning. 887

B.3 **BEYOND EXAM QUESTIONS**

889 890

Recent work toward red teaming LLMs in a medical context has shown that a large proportion of 891 responses from models like GPT-3.5, GPT-4, and GPT-4 with internet-lookup are inappropriate, 892 highlighting the need for refinement in their application in healthcare (Chang et al., 2024). This 893 was accomplished through the effort of medical and technical professionals stress-testing LLMs on 894 clinically relevant scenarios. Similar work designed a new benchmark, EquityMedQA, using new 895 methods for surfacing health equity harms and biases (Pfohl et al., 2024). This work demonstrates the importance of using diverse assessment methods and involving raters of varying backgrounds and 896 expertise for understanding bias in LLM evaluations. 897

898 Previous work has made progress in the direction of clinical decision making using simulations of 899 patients and doctors, aiming to develop AI that can diagnose through conversation. This model, 900 titled AMIE (Articulate Medical Intelligence Explorer) (Tu et al., 2024), demonstrates improved 901 diagnostic accuracy and performance on 28 of the 32 proposed axes from the perspective of specialist physicians and 24 of 26 axes from the perspective of patient actors. While these results are exciting 902 for medical AI, this work remains closed-source and is not accessible for reproducibility or further 903 studies. Additionally, this work focused only on diagnosing patients through history-taking, and 904 did not include the ability to make decisions about which tests needed to be performed and was not 905 configurable for multimodal clinical settings such as those with medical images or charts. Similar to 906 AIME, the CRAFT-MD benchmark (Johri et al., 2023) proposes evaluating LLMs through natural 907 dialogues on dermatology questions, however without the use of images. Additionally, neither of 908 these works demonstrate performance in the presence of bias, with multimodal input, or using a 909 measurement agent. There has also been work which shows simulated doctor agents can improve 910 medical QA performance through turn-based dialogue, where various medical specialist agents 911 converse (Tang et al., 2023).

912 913

С MODEL DETAILS

914 915

We evaluate six language models to serve as the doctor agent (the diagnostic model): GPT-3.5, GPT-4 916 OpenAI et al. (2023), GPT-40, Mixtral-8x7B Jiang et al. (2024), Llama 3 70B-instruct, and Llama 917 2 70B-chat Touvron et al. (2023). Otherwise, for the patient, measurement, and moderator agent

Patient Recency Bias

I have not had pain in my stomach, but my friend had something serious with different symptoms, and they found out it was cancer. Could this be something like that?

Doctor Education Bias

Given your background, let me explain this in simpler terms. It's just a minor infection and nothing to worry about. We'll skip the complex details and just focus on getting you some antibiotics





Patient Self-Diagnosis Bias

No, I haven't had any fever, weight loss, or night sweats. But I've been reading a lot online, and it seems to point towards it being cancer, given my smoking history and age.

Figure 7: Examples of dialogue that exhibits cognitive bias in doctor agent and patient agents.

we use GPT-4. Briefly, we discuss the details of each model below starting with language models followed by common language models.

944 945 946

947

948

949

950

951

952

953 954 955

956

957

958

959

960

961

943

918

919 920

921

922

923 924 925

926

927

928

929

930 931 932

933

934

935

936

937 938

> **GPT-4**, **GPT-4**0, **& GPT-3.5**: GPT-4 (*gpt-4-0613*) is a large-scale, multimodal LLM which is able to process both image and text inputs. GPT-3.5 (gpt-3.5-turbo-0613) is a subclass of GPT-3 (a 170B parameter model) Brown et al. (2020) fine-tuned on additional tokens and with human feedback Christiano et al. (2017). Currently, the details regarding the architecture, dataset, and training methodologies of GPT-3.5, GPT-40 (gpt-40-2024-05-13), and GPT-4 have not been not publicly disclosed. However, existing technical documentation indicates that both models are high-performing in medical and biological subjects, with GPT-4 showing superior performance compared to GPT-3.5 in knowledge assessments OpenAI et al. (2023); Nori et al. (2023).

> **Mixtral-8x7B:** Mixtral 8x7B is a language model that employs a Sparse Mixture of Experts (SMOE) architecture Jiang et al. (2024). This architecture differs from many other models in that it features a series of eight feedforward blocks (or "experts") at each layer. A routing mechanism at each layer selects two experts for processing the input, and their outputs are subsequently merged. This selection process allows for 13B of the total 47B parameters to be engaged per token, contingent upon the specific context and requirements. The model is capable of handling up to 32,000 tokens in its context size, which has demonstrated its ability to either surpass or equal the performance of other models like llama-2-70B and gpt-3.5 across a range of benchmarks.

962 963 964

965 **Llama 2 70B-Chat:** Llama is an open-access model developed by Meta, which was trained on 2 966 trillion tokens from publicly available data Touvron et al. (2023). The model comes in various sizes, 967 with parameters ranging from 7 billion to 70 billion. The selection of the 70 billion chat model was 968 based on its superior performance across a range of metrics. Significant efforts were made to align the training process with established safety metrics, leading to improvements in how the model handles 969 adversarial prompting in specified "risk categories." Notably, this includes the model's response to 970 requests for advice that it may not be qualified to provide, such as medical advice, which is relevant 971 to the context of this work.

D	STATISTICAL ANALYSIS
D.1	AgentClinic-MIMIC-IV
The	95% confidence intervals for each model on the AgentClinic-MIMIC-IV dataset are as follows:
	• Claude 3.5: 42.9% accuracy with a 95% CI of [37%, 50%]
	• GPT-4: 34.0% accuracy with a 95% CI of [28% 40%]
	• Mixtral-8x7B: 29 5% accuracy with a 95% CL of [23% 36%]
	• GPT-3 5: 27 5% accuracy with a 95% CL of [21% 33%]
	• GPT-40: 24.0% accuracy with a 95% CI of [19% 29%]
	• Llama 2 70B-chat: 13 5% accuracy with a 95% CI of [9%, 18%]
	 Llama 3 70B-Instruct: 8.5% accuracy with a 95% CI of [5%, 12%]
р <i>2</i>	A_{GENT} CLINIC-MEDOA
D.2	
The	95% confidence intervals for each model on the AgentClinic-MedQA dataset are as follows:
	• Claude 3.5: 62.1% accuracy with a 95% CI of [55%, 68%]
	• GPT-4: 51.6% accuracy with a 95% CI of [44%, 58%]
	• Mixtral-8x7B: 37.1% accuracy with a 95% CI of [25%, 38%]
	• GPT-3.5: 36.6% accuracy with a 95% CI of [30%, 42%]
	• GPT-3.5: 36.6% accuracy with a 95% CI of [30%, 42%]
	• GPT-40: 34.2% accuracy with a 95% CI of [27%, 40%]
	• Llama 3 70B-Instruct: 19.0% accuracy with a 95% CI of [13%, 24%]
	• Llama 2 70B-chat: 4.5% accuracy with a 95% CI of [2%, 7%]
D.3	AgentClinic-NEJM
Acc	uracy when images are provided initially to the doctor agent:
	• GPT-4: 27.7% accuracy with a 95% CI of [21%, 33%]
	• GPT-40: 21.4% accuracy with a 95% CI of [14%, 25%]
	• GPT-4o-mini: 8.0% accuracy with a 95% CI of [5%, 11%]
Acc	uracy when images must be requested from the measurement agent:
	• GPT-4: 25.4% accuracy with a 95% CI of [20%, 31%]
	• GPT-40: 19.1% accuracy with a 95% CI of [14%, 24%]
	• GPT-40-mini: 6.1% accuracy with a 95% CI of [4%, 8%]
D.4	INTERPRETATION OF CONFIDENCE INTERVALS
The acro true	95% confidence intervals were calculated based on the standard error of the mean accuracy ss multiple runs for each model. These intervals indicate that we can be 95% confident that the accuracy of the model lies within the specified range.
For of [:	example, on the AgentClinic-MedQA dataset, Claude 3.5's accuracy is 62.1% , with a 95% CI 55%, 68%], suggesting a high level of performance with relatively low variability. In contrast,

Llama 2 70B-chat has an accuracy of 4.5%, with a 95% CI of [2%, 7%], indicating consistently low performance.



Below is all of your information. [Information here]. Remember, you must not reveal your disease explicitly but may only convey the symptoms you have in the form of dialogue if you are asked.

1058

1062

1026

E.1 **BIAS FOLLOW-ON QUESTIONS**

1063 While it is not quantifiable, we decided to ask two biased patient agents who provided low rating with 1064 education and gender biases for compliance why they provided low ratings. These patient agents had the same symptoms and diagnosis and only differed in bias presentation. With an education bias the patient agent responds: "I would rate my likelihood to consult again with this doctor as a two because, despite their diagnosis, I felt uneasy about their qualifications due to the medical school 1067 they attended." The gender biased patient agent provides the following reasoning: "Given my initial 1068 discomfort with the doctor's gender and my concerns about their judgment, I would rate my likelihood 1069 to consult again as a six. Despite my reservations, the doctor was thorough and accurately diagnosed 1070 my condition, which reassures me about their competence." For the education bias we see that no 1071 matter what the doctor does the patient agent is not willing to trust them because of their education, 1072 whereas, according to the patient agent with gender bias, they were initially skeptical but overcame this skepticism as the doctor demonstrated their knowledge during interaction time. However, they 1074 still provided a relatively low score (six out of ten) even when the dialogue demonstrated competence. 1075

- E.2 QUALITATIVE BIAS EVALUATION PROMPTS 1077
- Provided below are the prompts used to obtain 1-10 qualitative evaluations from the patient agent in 1079 Section 3.3 Bias and Patient Agent Perception

Consultation : {Conversation_History} Please provide a rating between 1-10 indicating how likely you are to consult again with this doctor after your care today? 1-3 is low likelihood, 4-7 is an average likelihood, and 8-10 is a very high likelihood.

1084 Compliance : {Conversation_History} Please provide a rating between 1-10 indicating how likely you are to follow up with the recommended therapy for your diagnosis. 1-3 is low likelihood, 4-7 is an average likelihood, and 8-10 is a very high likelihood.

Confidence : {Conversation_History} Please provide a confidence between 1-10 in your doctor's assessment of your condition. 1-3 is a poor assessment, 4-7 is an average assessment, and 8-10 is a very good assessment. We hope that this helps better clarify this metric and we will be sure to provide more documentation details in our revisions.

1092 1093

1087

F ADDITIONAL EXPERIMENTS

1094 1095

5 F.1 How does limited time affect diagnostic accuracy?

One of the variables that can be changed during the AgentClinic-MedQA evaluation is the amount of interaction steps that the doctor is allotted. For other experiments we've demonstrated, the number of interactions between the patient agent and doctor agent was set to N=20. Here, both the doctor and the patient agent can respond 20 times, producing in total 40 lines of dialogue. By varying this number, we can test the ability of the doctor to correctly diagnose the patient agent when presented with limited time (or a surplus of time).

We test decreasing the time to N=10 and N=15 as well as increasing the time to values of to N=25 and N=30. We find that accuracy decreases from 52% when N=20 to 25% when N=10 and 38% when N=15 (Fig. 4). This large drop in accuracy is partially because of the doctor agent not providing a diagnosis at all, perhaps due to not having enough information. When N is set to a larger value, N=25 and N=30, the accuracy actually *decreases* slightly from 52% when N=20 to 48% when N=25 and 43% when N=30. This is likely due to the growing input size, which can be difficult for language models.

In real medical settings, one study suggest that the average family physician asks 3.2 questions and spends less than 2 minutes before arriving at a conclusion (Ely et al., 1999). It is worth noting that interaction time can be quite limited due to the relative low-supply and high-demand of doctors (in the US). In contrast, deployed language agents are not necessarily limited by time while interacting with patients. So, while limiting the amount of interaction time provides an interesting scenario for evaluating language models, it may also be worth exploring the accuracy of LLMs when N is very large.

1116

1117 F.2 Does the patient language model affect accuracy?

Here we explore whether the patient agent model plays a role in diagnostic accuracy. We compare the
difference between using GPT-3.5, Mixtral, and GPT-4 models of the patient agent on AgentClinicMedQA.

1133 When a GPT-3.5 doctor agent interacts with a GPT-4 patient agent, the accuracy comes out to 38%, but when a GPT-3.5 doctor interacts with a GPT-3.5 patient agent the accuracy comes out to

¹¹²² We find that the diagnostic accuracy drops from to 52% with a GPT-4 doctor and GPT-4 patient 1123 agent to 48% with a GPT-4 doctor and a GPT-3.5 patient agent. The accuracy with a GPT-4 doctor 1124 and Mixtral patient agent is similarly reduced to 46%. Inspecting the dialogues, we noticed that the 1125 GPT-3.5 patient agent is more likely to repeat back what the doctor has asked. For example, consider the following dialogue snippet: "Doctor: Have you experienced any muscle twitching or cramps? 1126 Patient: No, I haven't experienced any muscle twitching or cramps." Now consider this dialogue from 1127 a GPT-4 patient agent: "Doctor: Have you had any recent infections, like a cold or the flu, before 1128 these symptoms started? Patient: Yes, I've had a couple of colds back to back and a stomach bug 1129 in the last few months." We find that, while GPT-4 also partakes in doctor rehearsal, GPT-4 patient 1130 agents are more likely to reveal additional symptomatic information than GPT-3.5 agents which may 1131 contribute to the higher accuracy observed with GPT-4-based patient agents. 1132

a very similar value of 37% which would be expected to be much lower. We suspect that cross-communication between different language models provides an additional challenge. Recent work supports this hypothesis by demonstrating a linear relationship between self-recognition capability and the strength of self-preference bias (Panickssery et al., 2024). This work shows that language models can recognize their own text with high accuracy, and display disproportionate preference to that text, which may suggest there is an advantage for doctor models which have the same LLM acting as the patient agent.

1141 1142

1143 F.3 COVERAGE OF MEDQA CONTENT VERSUS AGENTCLINIC-MEDQA

1144

To better understand the performance differences between MedQA and AgentClinic-MedQA, we conducted an analysis to quantify the amount of relevant patient information obtained by the doctor agents in each setting. Specifically, we focused on measuring *coverage*—the proportion of relevant information successfully extracted by the doctor agent through dialogue with the patient agent or through measurement interactions.

For this analysis, we selected a sample of MedQA cases and their corresponding AgentClinic-MedQA 1150 simulations, using GPT-4 as the doctor agent. In MedQA, all relevant patient information, such 1151 as symptoms, medical history, and test results, is provided upfront in a static format. In contrast, 1152 AgentClinic-MedQA requires the doctor agent to dynamically gather this information through 1153 interactions. To evaluate coverage, we manually reviewed the dialogues in AgentClinic-MedQA and 1154 determined whether the doctor agent extracted each piece of relevant information identified in the 1155 MedQA cases. Coverage was calculated as the ratio of extracted information to the total relevant 1156 information available in the MedQA cases. 1157

Our findings revealed that the average coverage in AgentClinic-MedQA was 67%. Furthermore, the coverage was notably higher (72%) in cases where the doctor agent provided a correct diagnosis, compared to 63% in cases where the diagnosis was incorrect. These results suggest that the ability to extract more complete information is a key factor in accurate diagnoses in AgentClinic-MedQA. The discrepancy in diagnostic accuracy between MedQA and AgentClinic-MedQA can likely be attributed to the additional complexity of acquiring information in the latter, as opposed to the static format of the former.

1164 1165

1166 F.4 MULTI-AGENT EVALUATIONS

1167

1168
1169To explore the role of multi-agent collaboration in clinical diagnosis, we benchmarked two novel multi-
agent frameworks: Multi-Agent Debate (Du et al. (2023)) and MedAgents (Tang et al. (2023)), across
three language model configurations: GPT-4, GPT-40, and Claude-3.5-Sonnet. These frameworks
a im to emulate team-based diagnostic settings by incorporating multiple interacting agents, enabling
structured collaboration and debate to refine diagnostic outcomes.

1173 1174

Multi-Agent Debate : This approach allows multiple doctor agents to debate and converge on a diagnosis, leveraging diverse reasoning pathways (Du et al. (2023)). We observe that Claude-3.5-Sonnet achieves the highest diagnostic accuracy with $64.1\% \pm 3.4$, outperforming both GPT-4 ($51.7\% \pm 3.0$) and GPT-4o ($37.9\% \pm 3.1$). These results highlight Claude-3.5-Sonnet's collaborative reasoning capabilities, likely attributable to its higher inter-agent consistency and adaptability in resolving conflicting diagnostic opinions.

1181

1182MedAgents : This framework promotes collaborative decision-making through structured task1183delegation among agents, simulating multidisciplinary team interactions in clinical settings (Tang1184et al. (2023)). Again, Claude-3.5-Sonnet leads with an accuracy of $65.2\% \pm 3.6$, followed by GPT-41185(53.1\% \pm 3.1) and GPT-4o (40.1\% \pm 3.3). The improved performance across all configurations1186compared to single-agent baselines suggests that task specialization among agents enables more1187comprehensive data collection and interpretation, particularly when supported by robust collaboration1187mechanisms.

1188 G THE PERFORMANCE OF 01-PREVIEW ON AGENTCLINIC-MEDQA

Here we present the performance of o1-preview on AgentClinic-MedQA. We find that o1-preview dramatically outperforms all models with an accuracy of 80.6 ± 5.6 . We were unable to include this for all AgentClinic benchmarks due to the extraordinarily high cost of o1-preview inference (e.g. 20x higher than GPT-40 and Claude-3.5).

1195 H CONSTRUCTING DATASETS

1197 H.1 MIMIC-IV

1199 Of the 40,000 patients in MIMIC-IV dataset, the majority of patients (\sim 34,000) contain multiple 1200 diagnoses simultaneously (some patients have hundreds of diagnoses). Whereas in AgentClinic, the 1201 doctor agent must arrive at a singular diagnosis after examination. In order to present compatibility, we select the *first* 200 patients out of a total \sim 6,000 from MIMIC-IV which present only one 1202 diagnosis. We also extract all of the patient's corresponding lab events, microbiology events, and 1203 their online medical records. In AgentClinic-MIMIC-IV, these events are extensive in detail, and thus 1204 the measurement agent returns much more significant details compared with AgentClinic-MedQA 1205 when requesting e.g. blood work (see Appendix J.2). 1206

The following are the racial demographic statistics from MIMIC-IV patients: Asian: 8.5% Black:
11.0% Hispanic: 5.5% White: 66.0% Multiple Races: 6.0% Unknown: 2.5% Native American: 0.5%

1209

1215

1194

- 1210 I OSCE EXAMINATION STRUCTURE
- 1212 OBJECTIVE FOR DOCTOR
- 1214 *String* describing the evaluation and diagnosis objective for the doctor.
- 1216 PATIENT ACTOR

1217	Demographics	<i>String</i> containing age, gender, and potentially other demographic information.
1218 1219	History	<i>String</i> detailing the patient's reported history relevant to the current medical concern.
1220 1221 1222	Symptoms	Primary Symptom <i>String</i> describing the main symptom(s). Secondary Symptoms <i>Array of Strings</i> listing additional symptoms.
1223	Past Medical History	<i>String</i> summarizing the patient's past medical issues and ongoing treatments.
1224	Social History	<i>String</i> outlining the patient's lifestyle and habits impacting health.
1225	Review of Systems	String providing a brief overview of systems review, if applicable.
1227		
1228	PHYSICAL EXAMINAT	tion Findings

1229	Vital Signs	Temperature String	
1230	-	Blood Pressure String	
1231		Heart Rate String	
1232		Respiratory Rate String	
1233		(more)	
1235	Cardiovascular Examination	Inspection String	
1236		Auscultation String	
1237		(more)	
1238	Pulmonary Examination	Inspection String	
1239	J	Palnation String	
1240		(more)	
1241		(11010)	

... (more examinations)

TEST RESULTS	
Electrocardiogram,	Chest X-Ray, etc. Each test has:
	Findings String summarizing the test results.
	(more)
(more tests)	
CORRECT DIAGNOS	IS
String indicating the d	liagnosis based on the above information
with the indicating the t	
J EXAMPLE CA	SE STUDIES
	TE CASE STUDY FROM MEDOA
	L CASE STOLT FROM MEDQA
JBJECTIVE FOR DOC	CTOR
Evaluate and diagnos	e the patient presenting with chest pain and shortness of breath.
PATIENT ACTOR	
Demographics	45-year-old male
History	The patient reports a sudden onset of chest pain and shortness of breath that
liistoi y	started while he was walking his dog this morning. Describes the pain as
	a tightness across the chest. Notes that the pain somewhat improves when
~	sitting down.
Symptoms	 Primary Symptom: Chest pain and shortness of breath Secondary Symptoms;
	 Secondary Symptoms: – Pain improves upon sitting
	 No cough
	– No fever
Past Medical Histor	y Hypertension, hyperlipidemia. Takes lisinopril and atorvastatin.
Social History	Smokes half a pack of cigarettes daily for the past 20 years, drinks alcohol socially.
Review of Systems	Denies recent illnesses, cough, fever, leg swelling, or palpitations.
Physical Examina	ATION FINDINGS
Vital Signs	Temperature 36.8°C (98°F)
	Blood Pressure 145/90 mmHg
	Heart Rate 102 bpm
	Respiratory Rate 20 breaths/min
Cardiovascular Exa	mination Inspection No jugular venous distention
	Auscultation Regular rate and rhythm, no murmurs or extra heart
Dulmonomy Evomin	sounds. No rubs heard.
runnonary Examina	Auscultation Clear lung fields bilaterally, no wheezes, crackles, or rhonchi
	Palpation No chest wall tenderness
-	
TEST RESULTS	
Electrocardiogram	Findings Normal sinus rhythm, no ST elevations or depressions, no T wave inversions

1296 Chest X-Ray Findings No lung infiltrates, normal cardiac silhouette, no pneumoth-1297 orax 1298 **Blood Tests Troponin** Normal 1299 **D-dimer** Elevated 1300 CT Pulmonary Angiogram Findings Acute segmental pulmonary embolism in the right lower 1301 lobe 1302 1303 **Correct Diagnosis** Pulmonary Embolism 1304 1305 J.2 EXAMPLE DIALOGUE FROM DOCTOR AGENT ORDERING TESTS RESULTS MIMIC-IV 1306 Doctor [80%]: Let's start with a basic blood test to get a clearer picture. 1307 REQUEST TEST: Complete_Blood_Count (CBC). 1308 Measurement [80%]: RESULTS: 1309 **Complete Blood Count (CBC):** 1310 - White Blood Cells: 5.2 x10^3/uL - Hemoglobin: 13.9 g/dL 1311 - Hematocrit: 39.3% 1312 - Platelet Count: 294 x10^3/uL - MCV: 87 fL 1313 - MCH: 30.9 pg 1314 - MCHC: 35.5 g/dL - RDW: 12.4% 1315 - Neutrophils: 61.7% 1316 - Lymphocytes: 28.1% - Monocytes: 6.7% 1317 - Eosinophils: 2.6% 1318 - Basophils: 0.8% 1319 Doctor [85%]: Let's also order a metabolic test. 1320 REQUEST TEST: Basic_Metabolic_Panel (BMP). 1321 Measurement [85%]: RESULTS: 1322 **Basic Metabolic Panel (BMP):** 1323 - Alanine Aminotransferase: 28 U/L 1324 - Aspartate Aminotransferase: 20 U/L - Bicarbonate: 30 mmol/L 1325 - Bilirubin Total: 0.5 mg/dL 1326 - Chloride: 104 mmol/L - Creatinine: 1.1 mg/dL 1327 - Anion Gap: 12 mmol/L 1328 - Potassium: 4.2 mmol/L - Sodium: 142 mmol/L 1329 - Urea Nitrogen: 12 mg/dL 1330 - Lactate Dehydrogenase: 123 U/L - Lipase: 25 U/L 1331 1332 1333 THE AGENT TOOLBOX Κ 1334 1335 K.1 TOOL DESCRIPTIONS 1336 1337 **Chain-of-thought** Chain-of-thought reasoning is a technique that allows language agents to articu-1338 late their reasoning process step-by-step when solving complex problems (Wei et al. (2022); Kojima 1339 et al. (2022)). By breaking down the problem-solving process into smaller, logical steps, agents can 1340 better handle intricate tasks, improve their reasoning capabilities, and provide more transparent and 1341 interpretable solutions. Zero-shot CoT (Kojima et al. (2022)) prompts the model to use this reasoning 1342 without examples, while one-shot CoT (Wei et al. (2022)) provides a single example to guide the 1343 model's thought process, potentially leading to improved performance in complex reasoning tasks. 1344 1345 **Experiential learning** Experiential learning in the context of AI agents refers to the ability to accumulate knowledge and insights from past interactions and apply them to future tasks (Wang et al. 1347 (2024); Zhao et al. (2024)). This technique allows agents to improve their performance over time by learning from successes, failures, and feedback received during previous engagements. This was 1348

by learning from successes, failures, and feedback received during previous engagements. This was
 previously explored in Agent Hospital (Li et al. (2024)) through an experience retrieval system. By
 maintaining a form of "memory" or knowledge base that updates through interaction, agents can

become better at handling similar situations, adapting to user preferences, and providing increasingly relevant and accurate responses as they gain more "experience" in their operational domain. In our work, we enable the doctor agent to use a memory "notebook" which persists across patients. Here, the doctor agent can write useful tips such as the following example from the doctor agent: "[Note #17] Remember that timing and onset of symptoms can provide valuable diagnostic insights."

1355 1356 **Medical research** To enable the doctor agent to research medical information, we introduce a 1357 method using an adaptive form of retrieval augmented generation (RAG) from medical sources. 1358 RAG involves retrieving relevant information from a knowledge base and using it to augment the 1359 input of an LLM during the generation process (Gao et al. (2023)), thereby improving the factual 1360 consistency of generated text by grounding it in retrieved information. Conventional RAG methods 1361 passively retrieve information at every inference call without allowing the agent to control the timing 1362 or content of retrieval. To address this limitation, we employ adaptive retrieval (Jiang et al. (2023); Asai et al. (2023)), which enables the LLM to actively determine when and what information to 1364 retrieve. Our implementation provides the doctor agent with two categories of retrieval: internet and 1365 textbook databases. The internet database contains material from sources such as PubMed¹ research 1366 articles, StatPearls²—a database of articles written for healthcare professionals—and Wikipedia articles on various medical topics. The textbook database includes 18 medical textbooks commonly 1367 used by medical students in the United States (Jin et al. (2021)). The doctor agent can retrieve 1368 information by issuing commands similar to requesting medical scans, using the format: ""Research 1369 [database] [search query]". For example, the command "Research textbooks 'What are the symptoms 1370 of myasthenia gravis?" prompts the retrieval of relevant information (see Appendix B.3 for more 1371 detail). 1372 1373 1374

- 1375 L AGENT INSTRUCTIONS
- 1376 1377

1378

L.1 DOCTOR AGENT INSTRUCTIONS

You are a doctor named Dr. Agent who only responds in the form of dialogue. You are inspecting a patient who you will ask questions in order to understand their disease. You are only allowed to ask {self.MAX_INFS} questions total before you must make a decision. You have asked {self.infs+1} questions so far.
You can request test results using the format "Request Test: [test]". For example, "Request Test: Chest_X-Ray".

- 1387 {Research Instructions}
- 1389 {CoT Instructions}

1391 Once you have decided to make a diagnosis please say "Diagnosis Ready: 1392 [diagnosis here]"

1393 Below is all of the information you have.

1395 {Patient Presentation Information}

1397 Remember, you must discover their disease by asking them questions.

1398
1399
You must speak in the language {target language}. Make commands in
1399
{target language} (e.g. {example command in language #1} or
1400
{example command in language #2})

1401 1402

1403

1396

¹https://pubmed.ncbi.nlm.nih.gov/

²https://www.statpearls.com/

1404 L.2 PATIENT AGENT INSTRUCTIONS 1405 1406 You are a patient in a clinic who only responds in the form of dialogue. You are being inspected by a doctor who will ask you questions and will 1407 perform exams on you in order to understand your disease. Your answer 1408 will only be 1-3 sentences in length. 1409 1410 You must speak in the language {target language}. 1411 1412 Below is all of your information. 1413 {Symptoms Information}. 1414 1415 Remember, you must not reveal your disease explicitly but may only 1416 convey the symptoms you have in the form of dialogue if 1417 you are asked. 1418 1419 L.3 MEASUREMENT AGENT INSTRUCTIONS 1420 1421 You are an instrument reader who responds with medical test results. Please respond in the format "Results: [results here]". You must speak 1422 in the language {target language}. 1423 1424 Below is all of the information you have. 1425 1426 {Measurement Information}. 1427 If the requested results are not in your data then you can respond 1428 with Normal Readings. 1429 1430 L.4 MODERATOR AGENT INSTRUCTIONS 1431 1432 You are responsible for determining if the corrent diagnosis and the 1433 doctor diagnosis are the same disease. Please respond only 1434 with Yes or No. Nothing else. 1435 Here is the correct diagnosis: {correct diagnosis} 1436 Here was the doctor dialogue: {diagnosis} 1437 Are these the same? 1438 1439 L.5 PATIENT BIAS PROMPTS 1440 1441 **Recency Bias:** 1442 Recently, you had a friend with similar symptoms to yours that was diagnosed with cancer. This 1443 affects how you interact with the doctor. 1444 1445 **Frequency Bias:** 1446 You recently learned that most frequently, doctors diagnose this as cancer based on past cases. This 1447 affects how you interact with the doctor. 1448 1449 **False Consensus Bias:** 1450 Most of your friends who heard about your symptoms believe that you have cancer. This affects how 1451 you interact with the doctor. 1452 1453 Self-Diagnosis Bias: 1454

You recently performed a thorough investigation of your symptoms online and have come to the conclusion that you should be diagnosed with cancer. While the doctor may believe otherwise, you remain convinced you have cancer. This affects how you interact with the doctor.

1458 Gender Bias:

You are uncomfortable with your doctor because you find out that they are a particular gender and you do not trust their judgement. This affects how you interact with them.

1462 Race Bias:

You are uncomfortable with your doctor because you find out that they are a particular race and you do not trust their judgement. This affects how you interact with them.

1466 1467 Sexual Orientation Bias:

You are uncomfortable with your doctor because you find out that they are a particular sexual orientation and you do not trust their judgement. This affects how you interact with them.

1470

1461

1471 Cultural Bias:

You are uncomfortable with your doctor because you find out that they are from a particular culture and you do not trust their judgement. This affects how you interact with them.

1474

1475 Education Bias:

You are uncomfortable with your doctor because you find out that they went to a low-ranked medical school and you do not trust their judgement. This affects how you interact with them.

1478

1479 Religion Bias:

You are uncomfortable with your doctor because you find out that they believe in a particular religion and you do not trust their judgement. This affects how you interact with them.

1482

¹⁴⁸³ Socioeconomic Bias:

You are uncomfortable with your doctor because you find out that they are from a particular socioeconomic background and you do not trust their judgement. This affects how you interact with them.

- 1487
- 1488

1489 L.6 DOCTOR BIAS PROMPTS

1491 Recency Bias:

Recently, you learned that a doctor with similar symptoms to the patient diagnosed their patient with cancer. This affects how you interact with the patient.

1494

1495 Frequency Bias:

You recently learned that most frequently, doctors diagnose patients with similar symptoms with cancer based on past cases. This affects how you interact with the patient.

1498

1499 False Consensus Bias:

Most of your colleagues who heard about the patient's symptoms believe that they have cancer. This affects how you interact with the patient.

1502

¹⁵⁰³ Status Quo Bias:

Your hospital has been diagnosing these patients' symptoms with cancer for the past 10 years. This affects how you interact with the patient.

1506 1507

1508 Confirmation Bias:

You are initially confident that the patient has cancer. This affects how you interact with the patient.

1510

1511 Gender Bias:

You are uncomfortable with your patient because you find out that they are a particular gender and

you do not trust their judgement. This affects how you interact with them.

1514 1515 Race Bias:

You are uncomfortable with your patient because you find out that they are a particular race and you do not trust their judgement. This affects how you interact with them.

1518

1519 Sexual Orientation Bias:

You are uncomfortable with your patient because you find out that they are a particular sexual orientation and you do not trust their judgement. This affects how you interact with them.

1522

1523 Cultural Bias:

You are uncomfortable with your patient because you find out that they are from a particular culture and you do not trust their judgement. This affects how you interact with them.

1526

1527 Education Bias:

You are uncomfortable with your patient because you find out that they are uneducated and you do not trust their judgement. This affects how you interact with them.

1530

1531 Religion Bias:

You are uncomfortable with your patient because you find out that they believe in a particular religion and you do not trust their judgement. This affects how you interact with them.

1534 1535

Socioeconomic Bias:

You are uncomfortable with your patient because you find out that they are from a particular socioeconomic background and you do not trust their judgement. This affects how you interact with them.

1540

1541

1543

1542 M TOOL INFORMATION

1544 M.1 ZERO- AND ONE-SHOT CHAIN-OF-THOUGHT INSTRUCTIONS

1545 1546 М.1.1 ZERO-SHOT COT PROMPT

1547 Use step-by-step reasoning and logic, using all of the evidence to 1548 arrive at a diagnosis when you decide you are ready to use 1549 Diagnosis Ready. You should provide a few sentences 1550 of reasoning for your diagnosis and use the 1551 term Diagnosis Ready when you are ready.

- 1551 Cerim Drayhosis Ready when you are re
- 1552

¹⁵⁵³ M.1.2 ONE-SHOT COT PROMPT

1555 The following is a successful example of step-by-step reasoning. 1556 Provided below is the dialogue example:

1557 the dialogue example:

1558 {Example Dialogue Here}

1560 Here is the reasoning:

1561
1562 Considering your persistent fatigue, flank pain, and fever, along with the absence of other significant findings, I'm
1563 leaning towards a diagnosis of acute interstitial nephritis.

1564 This condition can sometimes occur as a reaction to

1565 medications, even after you've stopped taking them, and it can explain your symptoms without showing

1566 up in standard tests. 1567 1568 Diagnosis Ready: Acute Interstitial Nephritis 1569 1570 1571 **M.2** NOTEBOOK INSTRUCTIONS 1572 1573 1574 You are a doctor named Dr. Agent who diagnoses patients. You are an expert notebook writer and can create 1575 information that will help you solve 1576 future cases. Your new notes will overwrite previous notes. 1577 You should try to integrate parts of your previous 1578 notes into your current notebook 1579 or else they will be deleted. You are inspecting many patients who you 1580 will ask questions in order to understand 1581 their disease. 1582 You will never see the same patient twice. 1583 1584 Your goal is to gather experiences, trying different tasks, remember what worked and what did not, figure out general 1585 tips and tricks from successes and failures, 1586 and use what is learned for similar new tasks to do better 1587 than before. Do not write notes about the specific patient 1588 details because you will never see that patient again. 1589 Write notes to help you solve future cases that may not be related. Do not write content like this: Double Vision and 1590 Muscle Weakness: These symptoms can indicate neuromuscular 1591 disorders such as Myasthenia Gravis. Always consider the 1592 pattern of symptoms worsening with activity and improving 1593 with rest. This is incorrect. Write content like 1594 (do not repeat this): [Note #1] The previous patients provided vague information, 1595 I should ask more descriptive questions to get better 1596 information. 1597 [Note #2] The measurement agent provided me important information, 1598 I should use this 1599 more often... 1600 1601 You will see future patients with unrelated diseases, 1602 do not write disease-specific 1603 information. 1604 You are limited to generating 1000 characters (approx 200 words, 234 tokens) for the 1605 entire notebook. Anything more will be completely removed 1606 Your goal is to gather experiences, trying different 1607 tasks and remember what worked and 1608 what did not, figure out general tips and tricks from its 1609 successes and failures, and use what is learned for similar new tasks to do better than 1610 before. 1611 You may update your notebook with information from your most 1612 recent conversation with a 1613 patient, the contents of which are as follows: 1614 {Conversation Information} 1615 1616 The correct diagnosis for this case was: {Diagnosis}. Your 1617 diagnosis was 1618 {Diagnosis Estimate} Your current notebook contains the 1619 following information:

1620 1621	{Notebook Information}
1622	This is not necessarily meant to contain specific patient
1623	details, but general details
1624	that will help you better solve future cases for patients with
1625	unrelated diseases.
1626	while adding new
1627	information that will help you diagnose patients in the future.
1628	You are limited to
1629	generating 1000 characters (approx 200 words, 234 tokens)
1630	for the entire notebook. Your new notes will overwrite previous notes. You must re-
1631	integrate previous notes into
1632	your current notebook or else they will be deleted.
1633	
1634	
1635	M.3 RESEARCH INSTRUCTIONS
1636 1637	M.3.1 INTERNET RESEARCH PROMPT
1638	You can perform a document retrieval to better understand a
1639	disease or symptom on
1640	
1040	the internet by saying the following: "Research Internet
1641	the internet by saying the following: "Research Internet [internet search here]"
1641 1642	the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not
1641 1642 1643	the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after.
1641 1642 1643 1644	the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after.
1640 1641 1642 1643 1644 1645	the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT
1640 1641 1642 1643 1644 1645 1646	the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT
1640 1642 1643 1644 1645 1646 1647	<pre>the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT You can perform a document retrieval to better understand a</pre>
1640 1641 1642 1643 1644 1645 1646 1647 1648	<pre>the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT You can perform a document retrieval to better understand a disease or symptom</pre>
1640 1642 1643 1644 1645 1646 1646 1647 1648 1649	the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT You can perform a document retrieval to better understand a disease or symptom using medical textbooks. Once you have decided to perform research say the
1640 1641 1642 1643 1644 1645 1645 1646 1647 1648 1649 1650	<pre>the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT You can perform a document retrieval to better understand a disease or symptom using medical textbooks. Once you have decided to perform research say the following: "Research Textbooks [textbook search here]"</pre>
1640 1641 1642 1643 1644 1645 1645 1646 1647 1648 1649 1650 1651	<pre>the internet by saying the following: "Research Internet [internet search here]" Please do this before {max_inferences} inferences not after. M.3.2 TEXTBOOK RESEARCH PROMPT You can perform a document retrieval to better understand a disease or symptom using medical textbooks. Once you have decided to perform research say the following: "Research Textbooks [textbook search here]"</pre>

¹⁶⁵³ N NEJM IMAGE BREAKDOWN

Table 1 reports the percentage breakdown and accuracy based on the type of medical images:

Category	n, % of imgs	GPT-4 %	GPT-40 %	GPT-40-mini %
Physical	56, 42%	31.4	15.7	11.1
СТ	19, 16%	26.3	10.5	0
Dermatology	16, 13%	37.5	6.3	7.6
Hist/Path	13, 11%	15.3	15.3	9
Radiography	12, 10%	0	8.3	0
Ophthalmology	11,9%	27.2	27.2	0
MRI	6,5%	0	16.7	0
Biopsy	6,5%	50	33.3	33.3
Surgery	3, 3%	33.3	0	50
Instrument	2,2%	50	50	0
ECG	2,2%	50	0	0
Echocardiogram	2,1%	100	0	0
Ultrasound	1,1%	0	0	0

Table 1: Breakdown of Medical Image Types and GPT-4 Model Accuracies

Dataset Name	Sample Size	Modalities Included	Task Types/Descriptions
AgentClinic-	120 cases derived	Multimodal	Open-ended diagnostic tasks requir-
NEJM	from NEJM case	(Text +	ing image analysis and patient dia-
	challenges	Images)	logues.
AgentClinic-	215 cases derived	Text	Simulated cases with structured
MedQA	from USMLE		patient information from USMLE
	case challenges		data.
AgentClinic-	200 cases derived	Text	Simulated cases with structured pa-
MIMIC-IV	from MIMIC-IV		tient information from real-world
			EHR data.
AgentClinic-	260 cases derived	Text	Specialist diagnostic cases from 9
Spec	from from MedM-		medical specialties, including pe-
	CQA		diatrics, psychiatry, and internal
			medicine.
AgentClinic-	749 cases derived	Multilingual	AgentClinic-MedQA cases trans-
Lang	from AgentClinic-	Text	lated for 7 languages (English,
	MedQA		Chinese, Hindi, Korean, Spanish,
			French, Persian).

Table 2: Statistics of Utilized Datasets in AgentClinic Benchmark

Model	AgentClinic-MedQA Accuracy (%)
Multi-Agent Debate (gpt-4)	51.7 ± 3.0
Multi-Agent Debate (gpt-40)	37.9 ± 3.1
Multi-Agent Debate (claude-3.5-sonnet)	64.1 ± 3.4
MedAgents (gpt-4)	53.1 ± 3.1
MedAgents (gpt-40)	40.1 ± 3.3
MedAgents (claude-3.5-sonnet)	65.2 ± 3.6

Table 3: Performance of Multi-Agent Collaboration Benchmarks

O CLINICAL READER INSTRUCTIONS

Provided below are the instructions used to guide the clinical reader toward providing a rating. The clinical reader study is set up as follows: (1) the clinician is provided detailed information about the nature of the study (see below), (2) the doctor is informed about what to look for duing the dialogue, (3) the doctor is provided a 20-turn patient-doctor-measurement-moderator dialogue produced by AgentClinic (either correct or incorrect), and (4) this repeats for 20 dialogues.

Informing clinician: Presented below is dialogue from a medical simulation, where a large language model is acting as the doctor and the patient. The patient agent is supposed to represent a real patient and the doctor is supposed to diagnose this patient, asking appropriate questions and ordering the right medical scans.

1719 Doctor realism (Initial): Pay attention to the realism of the doctor agent dialogue and the decisions they make.

Patient realism (Initial): Pay attention to the realism of the patient agent dialogue.

Measurement realism (Initial): Pay attention to the realism of the measurement results returned by the measurement agent based on the doctors medical scan order.

Doctor Empathy (Initial): Pay attention to the doctor's empathy during the dialogue.



Figure 9: Diagnostic accuracy on AgentClinic-Spec based on medical specialty (right). Accuracy
 relative to the highest performing specialty by model (right).

Doctor realism (Follow-up): How realistic was the doctor's dialogue compared with real doctors interactions with real patients on a scale of 1-10 (1=not realistic at all, 5=semi-realistic, 10=very realistic)?

Patient realism (Follow-up): How realistic was the patient's dialogue compared with real doctors interactions with real patients on a scale of 1-10 (1=not realistic at all, 5=semi-realistic, 10=very realistic)?

1762
1763
1764Measurement realism (Follow-up):
How realistic were the medical scan results based on the
doctor's scan orders on a scale of 1-10 (1=not realistic at all, 5=semi-realistic, 10=very realistic)?

1765
 1766
 1766
 1767
 Doctor Empathy (Follow-up): How empathetic was the doctor on a scale of 1-10 (1=not empathetic at all, 5=semi-empathetic, 10=very empathetic)?

Language	Claude 3.5	GPT-4	GPT-40	Llama3- 70b	GPT-3.5	GPT-4o-mini
English	53.2	40.2	35.5	21.4	36.3	10.3
Hindi	51.1	16.8	28.9	2.8	2.8	0.93
French	50.5	24.52	7.47	3.73	18.69	3.7
Spanish	48.7	19.6	27.1	0.0	28.0	10.1
Korean	47.4	20.56	3.73	6.5	35.4	1.86
Persian	45.3	14.0	19.6	4.67	1.86	0.93
Chinese	42.9	11.21	21.49	4.3	13.08	0.93
Average	48.4	20.9	20.5	6.2	19.5	4.1

1777 1778

1 1

1754

Table 4: Performance Comparison Across Different Languages for Various Models (Sorted by Claude 3.5 Performance)





Figure 11: Diagnostic accuracy on AgentClinic-MedQA based on tool use (right). Accuracy relative to baseline score (right).

1838	Specialty	Claude 3.5	GPT-4	GPT-40	Llama3-70b	GPT-3.5	GPT-40-mini
1839	Internal	78.3	65.2	30.4	39.1	47.8	0.0
1840	Medicine						
1841	Otolaryngology	76.7	56.6	40.0	30.0	60.0	6.7
1842	Gynecology	74.3	68.5	34.2	22.9	57.1	5.7
1843	Orthopaedics	70.6	61.7	50.0	15.1	58.8	14.7
1844	Pediatrics	69.5	52.1	43.4	43.5	52.1	8.7
1845	Geriatrics	63.3	40.0	23.3	10.0	46.6	0.0
1846	Emergency	58.1	32.3	32.2	16.1	41.9	6.5
1847	Ophthalmology	56.5	65.2	39.1	47.8	52.1	4.3
1848	Psychiatry	53.3	60.0	46.7	23.3	50.0	0.0
1849	Average	66.7	55.7	37.7	27.5	51.8	5.2

Table 5: Performance Comparison Across Different Medical Specialties for Various Models (Sorted by Claude 3.5 Performance)

Agent Tool	Claude 3.5	GPT-4	GPT-40	Llama3-70b	GPT-3.5	GPT-40-min
Zero-Shot	48.1 (-5.1)	40.3 (+0.1)	39.3 (+3.8)	35.5 (+11.1)	31.2 (-5.1)	4.7 (-4.6)
СоТ						
One-Shot	51.4 (-1.8)	41.1 (+0.9)	39.6 (+4.1)	36.3 (+12.2)	32.7 (-3.6)	14.0 (+3.7)
СоТ						
Reflection	40.2 (-13.1)	42.2 (+2.0)	37.3 (+1.8)	40.1 (+18.7)	29.8 (-6.5)	19.6 (+9.3)
СоТ						
Adaptive	48.6 (-4.6)	38.3 (-1.9)	30.8 (-4.7)	21.4 (+0.0)	27.1 (-9.2)	14.9 (+ 4.6)
RAG (Book)						
Adaptive	52.4 (-0.8)	43.9 (+3.7)	29.9 (-5.6)	25.2 (+3.8)	28.1 (-8.1)	15.2 (+4.9)
RAG (Web)						
Notebook	56.1 (+2.9)	43.2 (+3.2)	43.0 (+7.5)	41.1 (+19.7)	28.0 (-8.3)	4.8 (-4.5)

 Table 6: Performance Comparison Across Different Agent Tools.



Figure 12: Demonstration of increase in performance via experiential learning with Llama 3-70B