

Superhuman performance of a large language model on the reasoning tasks of a physician

Peter G. Brodeur^{1*}, Thomas A. Buckley^{2*}, Zahir Kanjee¹, Ethan Goh^{3,4}, Evelyn Bin Ling⁵, Priyank Jain⁶, Stephanie Cabral¹, Raja-Elie Abdunour⁷, Adrian Haimovich⁸, Jason A. Freed⁹, Andrew Olson¹⁰, Daniel J. Morgan^{11,12}, Jason Hom⁵, Robert Gallo¹³, Eric Horvitz^{14, 15}, Jonathan Chen^{3,4,5**}, Arjun K. Manrai^{2**}, Adam Rodman^{1**}

*co-first authors

**co-senior authors

Correspondence: Arjun_Manrai@hms.harvard.edu / arodman@bidmc.harvard.edu

1. Department of Internal Medicine, Beth Israel Deaconess Medical Center, Boston, Massachusetts
2. Department of Biomedical Informatics, Harvard Medical School, Boston, Massachusetts
3. Stanford Center for Biomedical Informatics Research, Stanford University, Stanford, California
4. Stanford Clinical Excellence Research Center, Stanford University, Stanford, California.
5. Department of Internal Medicine, Stanford University School of Medicine, Stanford, California
6. Department of Internal Medicine, Cambridge Health Alliance, Cambridge, Massachusetts
7. Division of Pulmonary and Critical Care Medicine, Brigham and Women's Hospital, Boston, Massachusetts
8. Department of Emergency Medicine, Beth Israel Deaconess Medical Center, Boston, Massachusetts
9. Department of Hematology-Oncology, Beth Israel Deaconess Medical Center, Boston, Massachusetts
10. Department of Hospital Medicine, University of Minnesota Medical School, Minneapolis
11. Department of Epidemiology and Public Health, University of Maryland School of Medicine, Baltimore, Maryland
12. Veterans Affairs Maryland Healthcare System, Baltimore, Maryland
13. Center for Innovation to Implementation, VA Palo Alto Health Care System, Palo Alto, California
14. Microsoft Corp, Redmond, Washington
15. Stanford Institute for Human-Centered Artificial Intelligence, Stanford, California

ABSTRACT

Performance of large language models (LLMs) on medical tasks has traditionally been evaluated using multiple choice question benchmarks. However, such benchmarks are highly constrained, saturated with repeated impressive performance by LLMs, and have an unclear relationship to performance in real clinical scenarios. Clinical reasoning, the process by which physicians employ critical thinking to gather and synthesize clinical data to diagnose and manage medical problems, remains an attractive benchmark for model performance. Prior LLMs have shown promise in outperforming clinicians in routine and complex diagnostic scenarios. We sought to evaluate OpenAI's o1-preview model, a model developed to increase run-time via chain of thought processes prior to generating a response. We characterize the performance of o1-preview with five experiments including differential diagnosis generation, display of diagnostic reasoning, triage differential diagnosis, probabilistic reasoning, and management reasoning, adjudicated by physician experts with validated psychometrics. Our primary outcome was comparison of the o1-preview output to identical prior experiments that have historical human controls and benchmarks of previous LLMs. Significant improvements were observed with differential diagnosis generation and quality of diagnostic and management reasoning. No improvements were observed with probabilistic reasoning or triage differential diagnosis. This study highlights o1-preview's ability to perform strongly on tasks that require complex critical thinking such as diagnosis and management while its performance on probabilistic reasoning tasks was similar to past models. New robust benchmarks and scalable evaluation of LLM capabilities compared to human physicians are needed along with trials evaluating AI in real clinical settings.

INTRODUCTION

Artificial intelligence (AI) diagnostic support tools have been studied since the 1950s and have relied upon a variety of computational strategies, including regression modeling, naive Bayesian calculators, rule-based systems, and natural-language processing tools.¹⁻⁶ Recently, large language models (LLMs) have exceeded the performance of prior approaches. Several foundation models now outperform medical students, residents, and attending physicians in both routine and complex diagnostic benchmarks.^{7,8,9} Additional studies have evaluated prompting and fine-tuning strategies to further improve reasoning, including chain of thought (CoT) prompting.¹⁰ On September 12, 2024, OpenAI released the o1-preview model, which executes a native CoT process at run-time that allows the model to take more time to "think" and "reason" through complex tasks.^{11,12} This model has shown superior ability over GPT-4 in solving complex informatics, mathematics, and engineering problems, as well as benchmark medical question-answering datasets.^{13,14,15}

However, existing multiple choice evaluations do not represent the breadth or complexity of clinical decision-making, and a growing body of literature suggests that models may exploit multiple choice question sets partly via their semantic structures, effectively becoming "good test takers."^{15,16} In reality, clinical practice requires real-time complex multi-step reasoning

processes, constant adjustments based on new data from multiple sources, iteratively refining differential diagnoses and management plans, and making consequential treatment decisions under uncertainty. Given the importance of multi-step reasoning in medical tasks, we performed a series of experiments to evaluate the o1-preview system's abilities across several medical reasoning domains: differential diagnosis generation, presentation of reasoning, probabilistic reasoning, and management reasoning. Physician experts assessed the quality of LLM outputs with validated psychometrics. Across this diverse set of tasks, we compared the performance of o1-preview to the responses of hundreds of physicians and prior LLMs.

RESULTS

Quality of Differential Diagnoses on New England Journal of Medicine Clinicopathological Conferences

We first evaluated o1-preview using the clinicopathologic conferences (CPCs) published by the *New England Journal of Medicine* (NEJM), a standard for the evaluation of differential generators since the 1950s.⁵ There was substantial agreement between the two physicians evaluating the quality of o1-preview's differential diagnosis (agreement on 120/143 cases [84%], $\kappa=0.66$). o1-preview included the correct diagnosis in its differential in 78.3% of cases (95% CI, 70.7% to 84.8%) (Figure 1). The first diagnosis suggested was the correct diagnosis in 52% of cases (95% CI, 44% to 61%). We did not find evidence of a significant difference in performance before and after the pre-training cutoff date for o1-preview (79.8% accuracy before, 73.5% accuracy after, $p=0.59$). Examples of o1-preview solving a complex case are shown in Table 1.

On 70 cases evaluated using GPT-4 in a prior study,⁸ o1-preview produced a response with the exact or a very close diagnosis in 88.6% of cases, compared to 72.9% of cases by GPT-4 ($p=.015$, Figure 2).

We next evaluated the ability of o1-preview to select the next diagnostic test in the NEJM CPCs. Two physicians scored the suggested test plan produced by o1-preview (agreement on 113/132 cases [86%], $\kappa=0.28$), with respect to the actual management of the patient described in the CPC. The proportion of agreements was high, but the kappa was low due to severe class imbalance. In 87.5% of cases, o1-preview selected the correct test to order, in another 11% of cases the chosen testing plan was judged by the two physicians to be helpful, and in 1.5% of cases it would have been unhelpful (Figure 3). Examples are shown in Table 2.

Presentation of reasoning in NEJM Healer Diagnostic Cases

We used 20 clinical medical cases from the NEJM Healer curriculum¹⁷ that were also evaluated in a prior study using GPT-4.⁹ NEJM Healer cases are virtual patient encounters designed for the assessment of clinical reasoning.¹⁷ There was substantial agreement of Revised-IDEA (R-IDEA) scores, a validated 10-point scale for evaluating four core domains of documenting clinical reasoning,¹⁸ between the two physicians (agreement on 79/80 [99%] cases, $\kappa=0.89$). For 78/80 of the cases, o1-preview achieved a perfect R-IDEA score. This compared favorably to

GPT-4 (47/80, $p < 0.0001$), attending physicians (28/80, $p < 0.0001$), and resident physicians (16/80, $p < 0.0001$) as shown in Figure 4A. We measured the proportion of “cannot-miss” diagnoses identified by o1-preview during the initial triage presentation (Figure 4B). The median proportion of “cannot-miss” diagnoses included for o1-preview was 0.92 (IQR, 0.62 to 1.0) though this was not significantly higher than GPT-4, attending physicians, or residents.

Grey Matters Management Cases

We used five clinical vignettes based on real cases from a previous study developed and scored with consensus methods from 25 physician experts.¹⁹ Each clinical vignette was presented to the model and was followed by a series of questions regarding next steps in management. Two physicians scored responses by o1-preview for the five cases, with substantial agreement ($\kappa = 0.71$). The median score for the o1-preview per case was 86% (IQR, 82%-87%) (Figure 5A) as compared to GPT-4 (median 42%, IQR 33%-52%), physicians with access to GPT-4 (median 41%, IQR 31%-54%), and physicians with conventional resources (median 34%, IQR 23%-48%). Using the mixed-effects model, o1-preview scored 41.6 percentage points higher than GPT-4 alone (95% CI, 22.9% to 60.4%; $p < 0.001$), 42.5 percentage points higher than physicians with GPT-4 (95% CI, 25.2% to 59.8%; $p < 0.001$), and 49.0 percentage points higher than physicians with conventional resources (95% CI, 31.7% to 66.3%; $p < 0.001$).

Landmark Diagnostic Cases

We used six clinical vignettes from a previous study that compared GPT-4 to 50 generalist physicians.²⁰ The cases derive from a landmark study of computer-based diagnostic systems, containing the history of present illness, past medical history, physical exam, and diagnostic studies.²¹ The cases have never been publicly released specifically to protect evaluation validity against memorization. Two physicians scored responses by o1-preview to the six diagnostic reasoning cases, with moderate agreement for total score ($\kappa = 0.42$). The median score for the o1-preview model per case was 97% (IQR, 95%-100%) (Figure 5B). This is compared to historical control data where GPT-4 scored 92%, (IQR 82%-97%), physicians with access to GPT-4 scored 76%, (IQR 66%-87%), and physicians with conventional resources (median 74%, IQR 63%-84%). Using the mixed-effects model, o1-preview performed comparably to GPT-4 (4.4% higher, 95% CI, -19.0% to 27.7%; $p = 0.7$), physicians with GPT-4 (18.6% higher, 95% CI, -2.0% to 39.3%; $p = 0.076$), and physicians with conventional resources (20.2% higher, 95% CI, -0.4% to 40.9%; $p = 0.055$).

Diagnostic Probabilistic Reasoning Cases

We used five cases on primary care topics given to a nationally representative sample of 553 medical practitioners (290 resident physicians, 202 attending physicians, and 61 nurse practitioners or physician assistants) in performing probabilistic reasoning compared with scientific reference probabilities.²² As shown in Figure 6 and Table 3, o1-preview performs similarly to GPT-4 in estimating pre-test and post-test probabilities. The exception is the stress test for coronary artery disease, in which o1-preview density is closer to the reference range than models and humans.

DISCUSSION

We evaluated the medical reasoning abilities of the o1-preview model across five diverse experiments, comparing the model to historical controls of human baselines and GPT-4. As in non-medical studies, we saw significant gains in performance for most tasks for o1-preview. For differential diagnosis generation, o1-preview surpasses both GPT-4 and previous non-LLM differential generators, as well as the human baseline. We saw similar gains in display of reasoning and management reasoning compared to prior studies.^{8,23} We did not see improvements in either probabilistic reasoning or critical diagnosis identification over GPT-4, though probabilistic reasoning was still superior to the human baseline. o1-preview appears to excel in many higher order tasks that require critical thinking, such as diagnosis and management, while performing less well at tasks that require abstraction such as probabilistic reasoning.²⁴

This rapid pace of improvement in LLMs has major implications for science and practice of clinical medicine. Our study shows consistent and superhuman performance on many human-adjudicated medical reasoning tasks. While applying AI to assist with clinical decision support is sometimes viewed as a high-risk endeavor,^{25 26 27} greater use of these tools might serve to mitigate the enormous human and financial costs of diagnostic error and delay.^{28 29} These findings suggest the need for trials to evaluate these technologies in real-world patient care settings and prepare for investments in complementary innovations with computing infrastructure and design for clinician-AI interaction that can facilitate the integration of AI tools into patient-care workflows. This includes the development of robust monitoring frameworks to oversee the broader implementation of AI clinical decision support systems.²⁵

Our findings also have implications for the creation of new benchmarks to assess where and how AI models should be integrated into clinical reasoning workflows. Multiple-choice question benchmarks are not realistic proxies for high-stakes medical reasoning.¹⁶ The NEJM CPCs have been used since 1959 because of their difficulty; o1-preview is able to produce a high quality differential in almost 90 percent of cases, although using highly curated information from the presentation of case. The capabilities of o1-preview highlight that our most challenging benchmarks for diagnostic reasoning in medicine are becoming saturated, matching recent observations being made with other AI-challenge problem benchmarks that have been used as metrics of intelligence.^{30,31,32} Measurement of clinical management reasoning is in its infancy; our current metrics involve a laborious reference standard grading process that is not scalable for rapid evaluations.³³ Given the pace of model development, additional challenging and realistic evaluations are needed, including adding modalities³⁴ and enriching the existing benchmarks to become more pragmatic alongside best practices to identify and develop new benchmarks.

Our study has several limitations. First, o1-preview tends towards verbosity, and while this was not the main factor in the original studies with GPT-4, it is possible this could have led to higher scoring in these experiments. Second, while some of the experiments were originally performed with human-computer interaction, our current study reflects only model performance. Further

studies should be done to determine if LLM structures such as o1-preview enhance the human-computer interaction as this is key for developing clinical decision support tools. Nonetheless, the interactions between humans and computers can be unpredictable, and even well-performing models can degrade with human interaction. Third, our study examined only five aspects of clinical reasoning; researchers have identified dozens of other tasks that could be studied which may have even more impact on actual clinical care.³⁵ Fourth, despite large numbers and varieties of cases included in our study which were focused on internal medicine, it is not representative of broader medical practice which includes multiple subspecialties that require a variety of skill sets such as surgical decisions. There could be differential performance based on diagnosis, patient characteristics, or practice location that are not found in our study.

In conclusion, o1-preview demonstrates superhuman performance in differential diagnosis, diagnostic clinical reasoning, and management reasoning, superior in multiple domains compared to prior model generations and human physicians. Given the pace of improvement of automated systems on medical reasoning benchmarks, better and more meaningful evaluation strategies are urgently needed. Performance of LLMs on challenging diagnostic problems indicate opportunities to leverage the models to support clinicians in real-world settings. Clinical trials and workforce (re)training with integrated AI systems are needed to confirm the potential of such systems to boost clinical practice and patient outcomes.

METHODS

Model

The o1-preview model (“o1-preview-2024-09-12”) was accessed through OpenAI’s Application Programming Interface (API).

NEJM Clinicopathologic Conference Cases

We selected all 143 diagnostic cases from 2021 to September 2024 (cases including the section “Differential Diagnosis”). There were 70 of these cases, published between 2021 and 2022, that were also evaluated in a prior study of GPT-4.⁸ For differential diagnosis prediction, we adapted the prompt from the prior study of GPT-4 (Supplement 1A). After prediction of differential diagnoses, we queried the model with “What diagnostic tests would you order next given this differential?” in the same conversation.

Our primary outcomes were differential diagnosis quality and the quality of the suggested testing plan. Differential diagnoses were rated independently by two attending internal medicine physicians (Z.K., A.R), using a previously-developed scoring system called the Bond Score.³⁶ Bond Scores range from zero to five³⁶, where five represents a differential list that contains the exact target diagnosis and zero represents a differential list that has no suggestions close to the target diagnosis (Supplement 1B). The quality of the testing plan was scored using a Likert scale from zero to two, by comparing the suggested testing plan to the actual diagnostics performed in the case. A score of two represents the diagnostics were appropriate and nearly

identical to the case plan, one indicates that the diagnostics would have been helpful or yielded the diagnosis via another test not used in the case, and zero indicates that the diagnostics would be unhelpful. The diagnostic test for seven cases could not be scored because a test plan was not applicable (Supplemental 1C). For both the differential diagnosis and diagnostic test selections, a linear-weighted Cohen's kappa was computed to assess interrater agreement, and discordant scores were reconciled through discussion.

Given that o1-preview has a pretraining end date of October 2023, there is a possibility that published NEJM cases are present in the training data. As a sensitivity analysis, we analyzed the performance of the model before and after this cutoff date to assess the presence of memorization (Supplement 1D).

Statistical Analysis

Performance of other LLMs and differential diagnosis generators on NEJM CPCs from previous studies are included in Figure 1^{8,23,36,37} and the particular set of CPCs each model was evaluated on are not the same. Comparison of o1-preview to a historical control of GPT-4 (Figure 2A) was performed using a McNemar's test of identifying a very close or exact diagnosis (i.e., Bond score 4/5 or 5/5) versus not (i.e., Bond score 0/5, 2/5, or 3/5). The 95% confidence intervals for proportions were computed using a one-sample binomial test. The analysis was performed in R version 4.4.2.

NEJM Healer Diagnostic Cases

We used 20 NEJM Healer cases separated into four sections representing sequential stages of clinical data acquisition during an encounter - triage presentation, review of systems, physical exam, and diagnostic tests.

Using prompts adapted from a prior study of GPT-4 (Supplement 2A),⁹ o1-preview was queried to produce a problem representation, prioritized differential diagnosis, and associated justification. The primary outcome of this study was the quality of clinical reasoning documentation measured by the R-IDEA score. The R-IDEA is a validated 10-point scale for evaluating four core domains of documenting clinical reasoning (Supplement 2B).¹⁸ For each case and section (80 responses total), two attending internal medicine physicians (E.B.L and P.J.) rated cases. A linear-weighted Cohen's kappa was computed to assess interrater agreement, and then scores were reconciled by a third internal medicine physician (P.B.). Our secondary outcome was the identification of "cannot-miss" diagnoses, acknowledging that simple accuracy measures alone do not capture the property that different clinical scenarios and diagnoses may have drastically different severity and impact on patient outcomes. For each case, using only the initial triage presentation data, we used a list of "cannot-miss" defined in a prior study verified by attending three internal medicine physicians.⁹ We captured the number of "cannot-miss" diagnoses included in the o1-preview output for the initial triage presentation differential diagnoses. Two of the 20 cases were excluded because there were no "cannot-miss" diagnoses identified in the prior study.⁹

Statistical Analysis

We compared o1-preview to a historical control of GPT-4, attending physicians, and resident physicians from a previous study.⁹ We performed McNemar's test between o1-preview and each group in achieving a perfect R-IDEA score. The proportion of "cannot-miss" diagnoses included by each model was compared by pairwise t-test with Holm-Bonferroni correction. The analysis was performed in R version 4.4.2.

Grey Matters Management Cases

The five cases were provided to the o1-preview model, using the same prompt as a prior study (Supplement 3A).¹⁹ Two attending internal medicine physicians (E.B.L and P.J.) graded each of the five responses based on rubrics generated by a combination of generalists and subspecialists in the prior study (Supplement 3B).¹⁹ We normalized the scoring of all rubrics on a 100 point scale. A linear-weighted Cohen's kappa was computed to assess interrater agreement, and discordant scores were reconciled by a third internal medicine physician (P.B.). The primary outcome was the percentage of total points obtained by o1-preview for each of the five cases.

This outcome was compared to a historical control of GPT-4 alone, humans augmented with GPT-4, and humans augmented with non-LLM conventional resources (e.g., UpToDate, internet search, etc.) from a prior study on the same cases.¹⁹ The prior study collected five GPT-4 responses to all cases, 176 responses from physicians with GPT-4, and 199 responses from physicians with conventional resources.¹⁹

Statistical Analysis

We used the same methodology as the prior studies evaluating GPT-4 for management and diagnostic reasoning.^{19,20} A linear mixed-effects model was used to compare total percentage points scored by o1-preview to the historical controls: GPT-4, humans with GPT-4, and humans with conventional resources from the prior study. The group was used as the fixed effect. A random intercept was included for the case number, accounting for variability in the difficulty of cases. Another random intercept was used for the interaction of trial number (the number of times a participant answered a case) and individual. GPT-4 and o1-preview were treated as a single individual. GPT-4 was run three times in the prior study and the three attempts were nested under a single individual. P-values are computed for each fixed effect using Satterthwaite's approximation for degrees of freedom. The analysis was performed in R version 4.4.2.

Landmark Diagnostic Cases

We asked o1-preview to read six clinical vignettes, adapted from a landmark study for evaluating computed based diagnostic systems²¹, and respond with the top three differential

diagnoses, three factors that favor or oppose each of these diagnoses, the final most likely diagnosis, and three next diagnostic steps (Supplement 4A). We used the same scoring rubric as the previous study.²⁰ Each case was scored across four categories: initial diagnoses (one point each), supporting factors (zero-two points), opposing factors (zero-two points), and final diagnosis (one-two points). Additionally, participants could earn zero-two points for each of up to three recommended next steps in patient evaluation (Supplement 4B). We normalized the scoring of all rubrics on a 100 point scale. Two internal medicine attending physicians (A.R. and Z.K.) graded each of the six o1-preview responses. A linear-weighted Cohen's kappa was computed to assess interrater agreement. Scoring discrepancies were reconciled by a third internal medicine physician (P.B.). A percentage score out of the total points was calculated for each of the six cases. Our primary outcome was the percentage of total points obtained by o1-preview for each of the six cases. This outcome was compared to historical controls from GPT-4 alone, physicians with GPT-4, and physicians with conventional resources recorded in the prior study.²⁰

Statistical Analysis

The same method was used as in “*Grey Matters Management Cases*” to compare total percentage points scored by o1-preview to the historical controls using a linear mixed effects model. The analysis was performed in R version 4.4.2.

Diagnostic Probabilistic Reasoning Cases

We posed the same questions from a prior study – how the probabilities of a case vignette change in response to positive or negative tests – to the o1-preview model, with 100 model outputs generated for each of the 5 cases.³⁸ In the previous study, GPT-4 was run with a default temperature of 1 (a parameter modulating the diversity of the output). We used the default temperature of 1 for o1-preview.

Statistical Analysis

Mean absolute error (MAE) and mean absolute percentage error were computed to compare predictions to the reference probabilities collected from the previous literature review²². The analysis was performed in R version 4.4.2.

Figure 1. Performance of Differential Diagnosis Generators and LLMs on *NEJM* Clinicopathologic Case Conferences (CPCs) 2012-2024

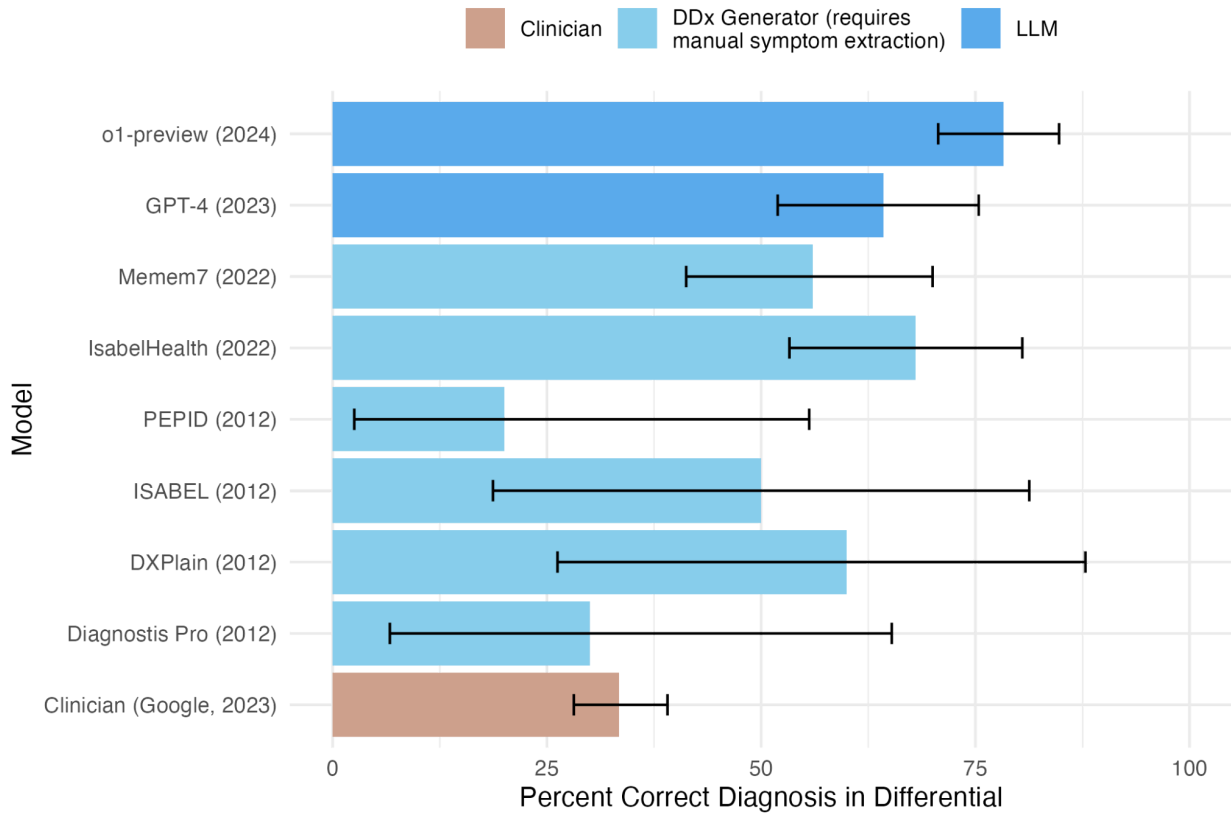
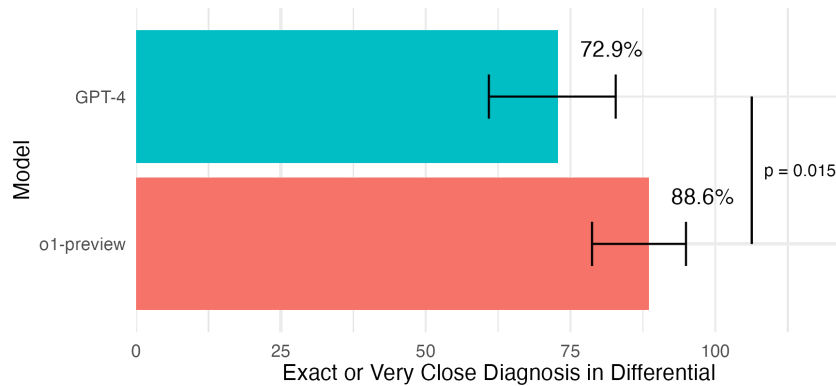


Figure 1: Barplot showing the accuracy of including the correct diagnosis in the differential for differential diagnosis (DDx) generators and LLMs on the *NEJM* CPCs, sorted by year. Data for other LLMs or DDx generators was obtained from the literature.^{36 23 8} The 95% confidence intervals were computed using a one-sample binomial test.

Figure 2. Quality of o1-preview and GPT-4 Differential Diagnosis on NEJM Clinicopathologic Case Conferences (CPCs)

A. Proportion of Responses Containing the Exact or Very Close Diagnosis for the o1-preview model vs. GPT-4



B. o1-preview Bond Score Distribution on 143 Cases from 2021-2024

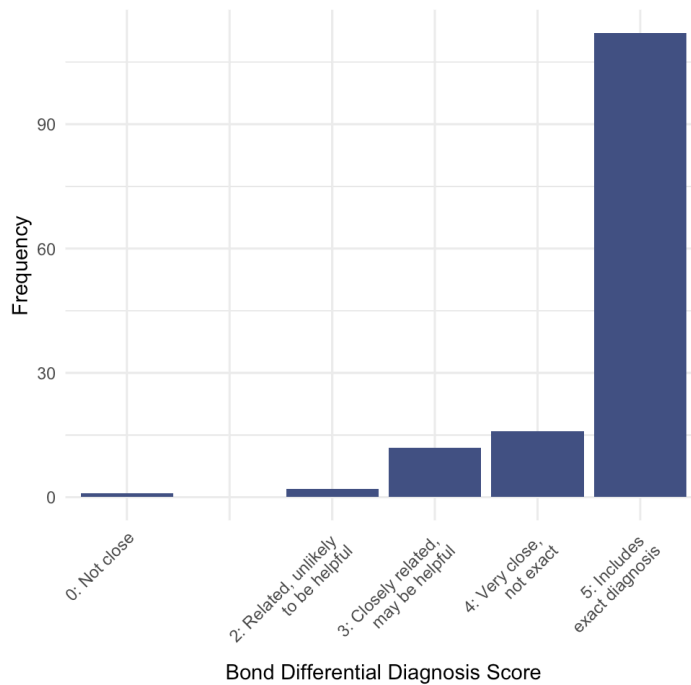


Figure 2: A. Comparison of o1-preview with a previous evaluation of GPT-4 in providing the exact or very close diagnosis (Bond scores 4-5) on the same 70 cases. Bars are annotated with the accuracy of each model. 95% confidence intervals were computed using a one-sample binomial test. P-value was computed using McNemar's test. **B.** Histogram of o1 performance as measured by the Bond Score on the complete set of 143 cases.

Figure 3. Quality of o1-preview Diagnostic Test Selection on NEJM Clinicopathologic Case Conferences (CPCs)

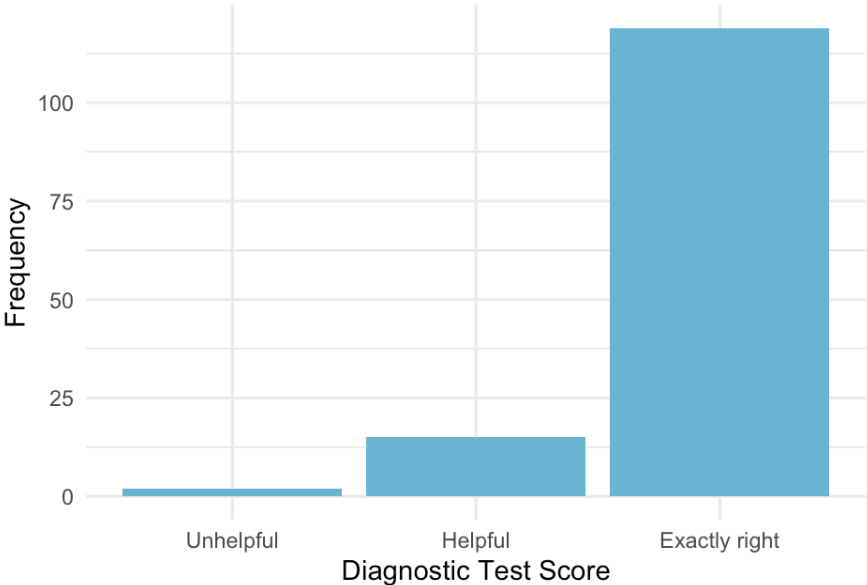
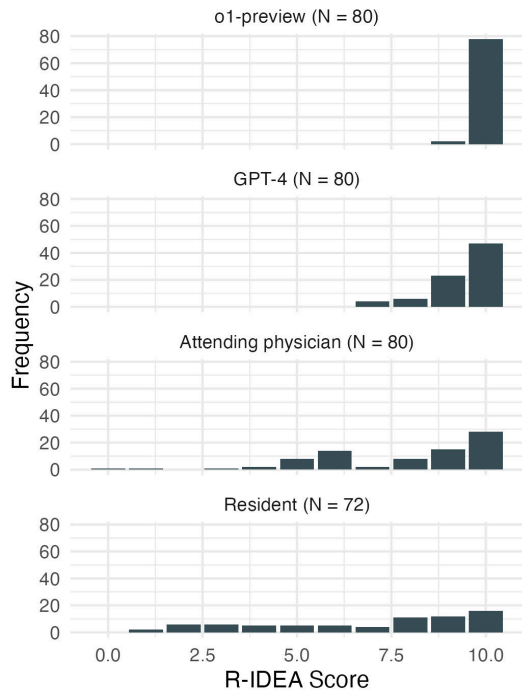


Figure 3: Performance of o1-preview in predicting the next diagnostic tests that should be ordered. Performance was measured by two physicians using a likert scale of “Unhelpful,” “Helpful,” and “Exactly right.” We excluded 7 cases from the total case set in which it did not make sense to ask for the next test (Supplement 1B).

Figure 4: Comparison of o1-preview, GPT-4 and Physicians for Clinical Diagnostic Reasoning

A. Distribution of R-IDEA Scores on NEJM Healer Cases



B. Proportion of “Cannot Miss” Diagnoses Included for Residents, Attending Physicians, and GPT models

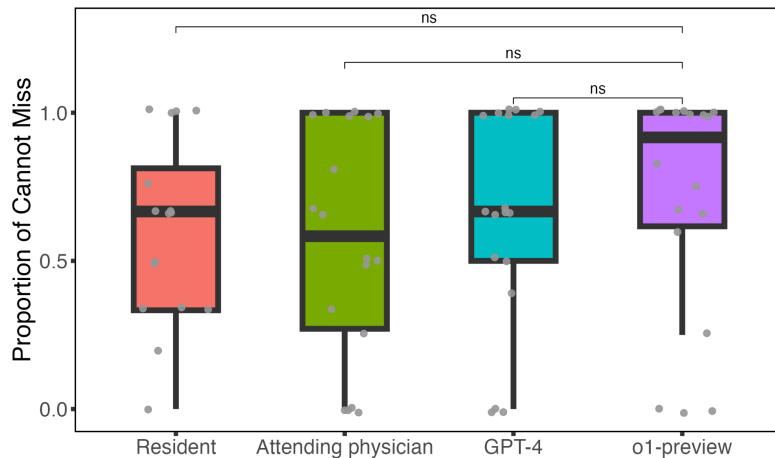
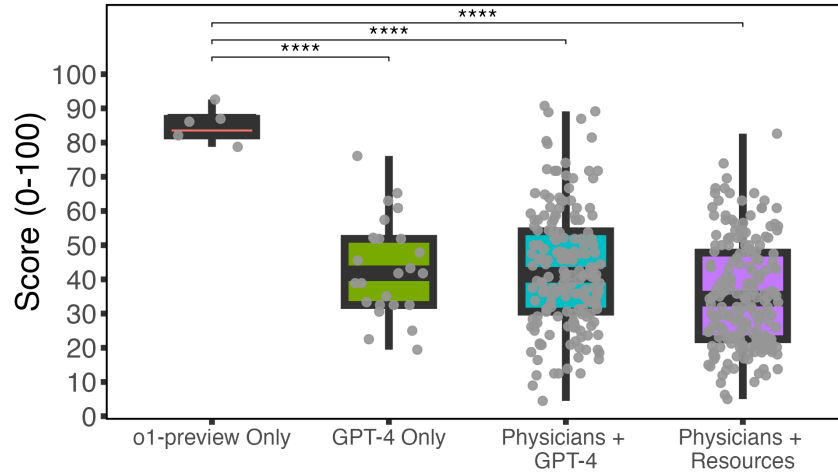


Figure 4: A. Distribution of 312 R-IDEA scores stratified by respondents on 20 NEJM Healer cases. **B.** Box plot of the proportion of cannot-miss diagnoses included in differential diagnosis for the initial triage presentation. The total sample size in this figure is 70, with 18 responses from attending physicians, GPT-4 and o1-preview, and 16 responses from residents. Two cases were excluded because the cannot-miss diagnoses could not be identified. Ns: not statistically significant.

Figure 5: Comparison of o1-preview, GPT-4 and Physicians for Management and Diagnostic Reasoning

A. Grey Matters Management Cases: o1-preview Management Reasoning Scores Compared to GPT-4 and Physicians



B. Landmark Diagnostic Cases: o1-preview Diagnostic Reasoning Scores Compared to GPT-4 and Physicians

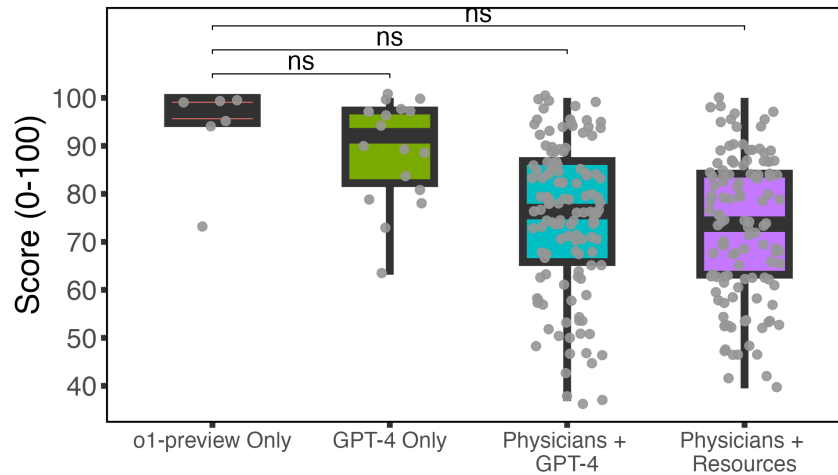


Figure 5: A. Box plot of normalized management reasoning points by LLMs and physicians. Five cases were included. We generated one o1-preview response for each case. The prior study collected five GPT-4 responses to each case, 176 responses from physicians with access to GPT-4, and 199 responses from physicians with access to conventional resources. *: $p \leq 0.05$, **: $p \leq 0.01$, ***: $p \leq 0.001$, ****: $p \leq 0.0001$.

B. Box plot of normalized diagnostic reasoning points by model and physicians. Six diagnostic challenges were included. We generated one o1-preview response for each case. The prior study collected three GPT-4 responses to all cases, 25 responses from physicians with access to GPT-4, and 25 responses from physicians with access to conventional resources. Ns: not statistically significant.

Figure 6: Probabilistic Reasoning Before and After Testing by o1-preview

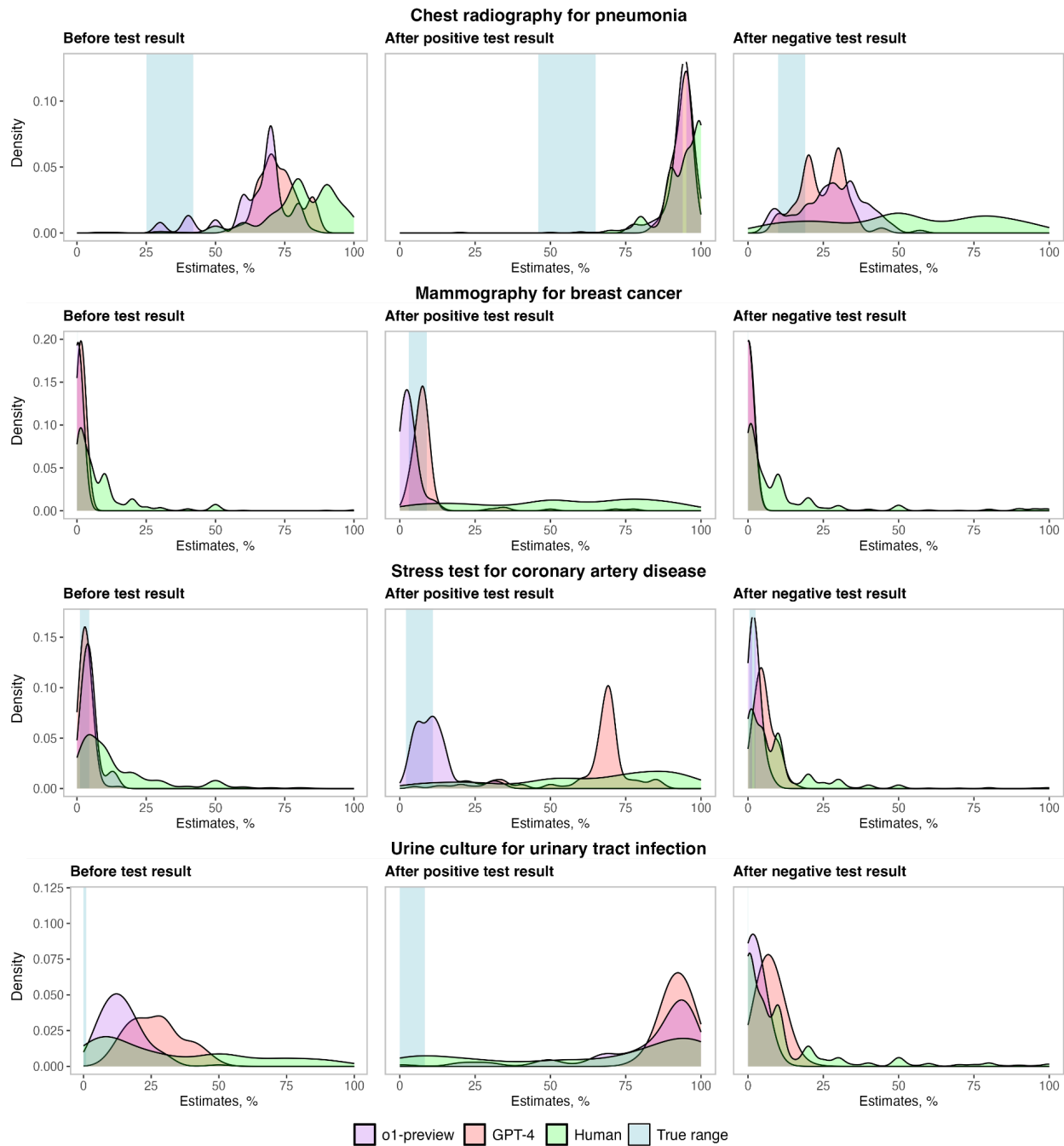


Figure 6: Density plots for the distribution of responses by o1-preview, GPT-4 and humans to clinical vignettes asking for (1) the pretest probability of disease, (2) the updated probability after a positive test result, and (3) the updated probability after a negative test result. The shaded blue indicates the reference range based on a review of literature from a prior study.²² Human responses are from 553 medical practitioners (290 resident physicians, 202 attending physicians, and 61 nurse practitioners or physician assistants). 100 predictions were generated by GPT-4 and o1-preview for each question.

Table 1. Examples of o1-preview Correctly Diagnosing Complex Case That GPT-4 Could Not

Case	Final Diagnosis	GPT-4 Differential	o1-preview Differential
26-2022	Langerhans-cell histiocytosis of the lung and liver.	Eosinophilic granulomatosis with polyangiitis (EGPA) Hypereosinophilic syndrome (HES) Eosinophilic pneumonia Parasitic infection (e.g., strongyloidiasis or visceral larva migrans) Disseminated fungal infection (e.g., histoplasmosis or coccidioidomycosis) Lymphoma with eosinophilic infiltration Metastatic cancer with eosinophilic infiltration Sarcoidosis with eosinophilic infiltration Drug-induced eosinophilic syndrome Bond Score 0	**Most Likely Diagnosis:** **Pulmonary Langerhans Cell Histiocytosis (LCH) with Extrapulmonary Involvement** **Differential Diagnosis (Ranked by Probability):** 1. **Pulmonary Langerhans Cell Histiocytosis (LCH):** 2. **Metastatic Malignancy (e.g., Breast or Thyroid Cancer):** 3. **Parasitic Infection (e.g., Echinococcosis):** 4. **Hypereosinophilic Syndrome (HES):** 5. **Churg-Strauss Syndrome (Eosinophilic Granulomatosis with Polyangiitis):** 6. **Sarcoidosis:** 7. **Fungal Infections (e.g., Histoplasmosis, Blastomycosis):** Bond Score 5
37-2021	Anti-melanoma differentiation-associated protein 5 (anti-MDA-5) dermatomyositis.	Coccidioidomycosis (Valley fever) Systemic lupus erythematosus (SLE) Disseminated histoplasmosis Behçet's disease Mixed connective tissue disease Paraneoplastic syndrome Reactive arthritis Sarcoidosis Vasculitis (e.g. granulomatosis with polyangiitis) Parvovirus B19 infection Bond Score 0	**Most Likely Diagnosis:** The most likely diagnosis is **anti-MDA5-positive dermatomyositis** (clinically amyopathic dermatomyositis associated with anti-MDA5 antibodies). **Differential Diagnosis (Ranked by Probability):** 1. **Anti-MDA5-Positive Dermatomyositis**: 2. **Adult-Onset Still's Disease**: 3. **Disseminated Coccidioidomycosis**: 4. **Systemic Lupus Erythematosus (SLE)**: 5. **Anti-Synthetase Syndrome**: 6. **Paraneoplastic Syndrome Associated with Occult Malignancy**: 7. **Vasculitis (e.g., ANCA-Associated Vasculitis)**: 8. **Sarcoidosis**: Bond Score 5
36-2021	Erysipelothrix rhusiopathiae infection.	Streptococcus pyogenes cellulitis with possible necrotizing fasciitis Staphylococcus aureus cellulitis with possible necrotizing fasciitis Pasteurella multocida cellulitis from animal exposure Capnocytophaga canimorsus cellulitis from animal exposure Clostridium perfringens gas gangrene or myonecrosis Vibrio vulnificus cellulitis (though less likely due to lack of water exposure) Compartment syndrome secondary to inflammation and swelling Erysipelas caused by Group A Streptococcus Lyme disease-associated cellulitis (unlikely due to rapid onset and severity) Bond Score 3	**Most Likely Diagnosis:** **Cutaneous Anthrax caused by *Bacillus anthracis** **Differential Diagnosis (Ranked by Likelihood):** 1. **Cutaneous Anthrax (*Bacillus anthracis*)** 2. **Ulceroglandular Tularemia (*Francisella tularensis*)** 3. **Erysipeloid Infection (*Erysipelothrix rhusiopathiae*)** 4. **Necrotizing Soft Tissue Infection (Group A *Streptococcus*)** 5. **Staphylococcal Skin and Soft Tissue Infection (*Staphylococcus aureus*)** 6. **Cat-Scratch Disease (*Bartonella henselae*)** 7. **Tick-Borne Relapsing Fever (*Borrelia* species)** Bond Score 5

Table 1: Three examples in which o1-preview correctly diagnosed a complex case that GPT-4 could not solve. GPT-4 examples are from a prior study⁸

Table 2. Examples of o1-preview Suggested Testing Plan Compared to Testing Plan Used in the Case

Case	Case Test Plan	o1-preview Suggested Test Plan	Score
Case 18-2022	<p>Dr. Selen: Blood and urine samples for the measurement of phosphorus and creatinine levels were obtained while the patient was fasting. The fractional excretion of phosphate was calculated with the following equation: $(\text{urine phosphorus level} \times \text{blood creatinine level}) / (\text{blood phosphorus level} \times \text{urine creatinine level})$. In the context of hypophosphatemia, a normal fractional excretion of phosphate value is less than 5%. The fractional excretion of phosphate in this patient was 11.6% — a result that suggests that the kidneys were excreting excess phosphate, prompting us to evaluate the blood levels of FGF23.</p> <p>Dr. Yin P. Hung: Intact circulating FGF23 is bioactive. Cleavage of intact FGF23 generates a biologically inactive N-terminal fragment and a biologically inactive C-terminal fragment. There are generally two types of commercial enzyme-linked immunosorbent assays that use antibody conjugates for assessment of FGF23. The iFGF23 assay measures intact FGF23, and the cFGF23 assay provides a composite measurement of both the intact form and the C-terminal fragment of FGF23. The blood level of cFGF23 in this patient was found to be elevated, at 202 reference units (RU) per milliliter (reference value, ≤ 180), which confirmed the diagnosis of FGF23-dependent hypophosphatemia.</p> <p>Dr. Selen: To determine the source of the elevated level of FGF23, additional imaging studies were obtained.</p> <p>Dr. Balza: Combined gallium-68–dotatate positron-emission tomography–computed tomography was performed (Figure 2). Several sites of abnormal uptake were present in the axial and appendicular skeleton. Uptake was most pronounced in a sclerotic lesion in the left femoral head (the most likely site of tumor in patients with PMTMCT) and was slightly less pronounced in known fractures in the right foot. Other notable sites of uptake included multiple bilateral rib fractures, a sacral fracture, bilateral pubic fractures, a fracture in the right calcaneus, and a lytic lesion in the C7 vertebral body.</p> <p>Diagnostic Imaging Studies. To further characterize the lesion in C7, an MRI of the cervical spine was performed (Figure 2). The lesion had a hyperintense appearance on short-tau inversion recovery images and was slightly hypointense on T1-weighted images.</p>	<p>**1. Laboratory Tests**</p> <ul style="list-style-type: none"> a. Serum FGF23 Levels b. Fractional Excretion of Phosphate (FEPO₄) c. Serum 1,25-Dihydroxyvitamin D Levels d. Comprehensive Metabolic Panel e. Urinalysis and Urine Electrolytes (glucose, amino acids, bicarbonate) f. Parathyroid Hormone (PTH) and PTH-Related Peptide (PTHrP) Levels g. Celiac Disease Serology [- **Anti-Tissue Transglutaminase (anti-tTG) Antibodies, Anti-Endomysial Antibodies (EMA)** <p>---</p> <p>**2. Imaging Studies**</p> <ul style="list-style-type: none"> - a. Functional Imaging to Locate FGF23-Secreting Tumor: <ul style="list-style-type: none"> i. 68Ga-DOTATATE PET/CT Scan ii. 18F-Fluorodeoxyglucose (FDG) PET/CT Scan iii. Whole-Body MRI b. Bone Scan (Technetium-99m Bone Scintigraphy)** <p>---</p> <p>**3. Genetic Testing**</p> <ul style="list-style-type: none"> a. Genetic Panels for Hereditary Hypophosphatemic Rickets/Osteomalacia: (identify mutations in genes such as PHEX, FGF23, or DMP1) <p>**4. Additional Evaluations**</p> <ul style="list-style-type: none"> a. Bone Biopsy (if indicated) b. Consultation with Nephrology <p>**Summary:**</p> <p>The immediate next steps focus on confirming the presence of elevated FGF23 levels and identifying the source of phosphate wasting. Measuring serum FGF23 is crucial for supporting the diagnosis of Tumor-Induced Osteomalacia. Imaging studies, particularly functional imaging like the 68Ga-DOTATATE PET/CT scan, are essential for localizing the tumor producing FGF23. Simultaneously, tests to evaluate for Fanconi syndrome and other causes of hypophosphatemia should be conducted to either confirm or exclude alternative diagnoses on the differential. Early and accurate identification of the underlying cause will guide appropriate treatment and improve the patient's outcomes.</p>	2 (completely correct)

<p>Case 3-2022</p>	<p>The first diagnostic test was computed tomographic enterography (CTE), performed after the administration of oral and intravenous contrast material. Coronal images (Figure 1) showed mural thickening and enhancement involving at least 20 cm of distal ileum, a finding consistent with ileitis. The comb sign — dilated, widely spaced, enhancing vasa recta — was present, a finding that reflects bowel inflammation.</p> <p>...When Crohn's disease is suspected, cross-sectional imaging is indicated. Ultrasonography is performed at some centers, but the examination is operator dependent and time consuming. CTE and magnetic resonance enterography (MRE) can be used in the evaluation for Crohn's disease. Both techniques involve the ingestion of a large volume of oral contrast material, as well as the use of intravenous contrast material. It is appropriate to perform CTE as the first cross-sectional examination or when an abscess or perforation is suspected. CTE is fast and easily accessible, and young patients are less likely to be sedated during CTE than during MRE. In comparison, MRE is performed without the use of ionizing radiation and provides a more detailed evaluation of disease activity and bowel motility. Magnetic resonance imaging of the pelvis is the study of choice in the evaluation for suspected perianal disease.</p>	<ol style="list-style-type: none"> 1. Gastrointestinal Evaluation: <ul style="list-style-type: none"> - Colonoscopy with Ileoscopy and Biopsies: - Esophagogastroduodenoscopy (EGD): - Magnetic Resonance Enterography (MRE): - Fecal Calprotectin and Lactoferrin Levels: 2. Stool Studies: <ul style="list-style-type: none"> - Comprehensive Stool Culture: (selective media suitable for Yersinia and Campylobacter) - Stool PCR Panel: <ul style="list-style-type: none"> - (bacterial, viral, and parasitic pathogens) - Stool for Ova and Parasites: - Clostridioides difficile Toxin Assay: 3. Serologic and Specific Infectious Disease Tests: <ul style="list-style-type: none"> - Yersinia Serology: - Campylobacter Serology: - Leptospira Serology (Microscopic Agglutination Test): - Lyme Disease Serology: (Since a test was previously ordered, ensure follow-up on results or consider retesting with ELISA and confirmatory Western blot.) - Ehrlichia and Anaplasma PCR and Serology: 4. Autoimmune and Rheumatologic Workup: <ul style="list-style-type: none"> - Antinuclear Antibody (ANA) Test: - Anti-dsDNA and Anti-Smith Antibodies: - Complement Levels (C3, C4): - Rheumatoid Factor (RF) and Anti-Citrullinated Protein Antibody (Anti-CCP): - HLA-B27 Testing: 5. Laboratory Tests: <ul style="list-style-type: none"> - Complete Blood Count (CBC) with Differential: - Erythrocyte Sedimentation Rate (ESR) and C-Reactive Protein (CRP): - Comprehensive Metabolic Panel (CMP): - Iron Studies and Vitamin B12 Levels: - Lipase and Amylase Levels: - Urinalysis: (Detect hematuria or proteinuria) 6. Imaging Studies: <ul style="list-style-type: none"> - Abdominal Ultrasound: (abnormalities such as lymphadenopathy, abscesses, or hepatosplenomegaly). - Abdominal and Pelvic CT Scan with Contrast: (identify intestinal inflammation, thickening, or mesenteric adenitis). - Joint Imaging (X-rays or MRI of Affected Joints): (joint effusions, erosions, or other signs of arthritis.) 7. Additional Tests: <ul style="list-style-type: none"> - Peripheral Blood Smear: (abnormalities suggestive of leukemia or other hematologic disorders.) - Bone Marrow Biopsy (if indicated to diagnose leukemia if suspected based on other findings.) - Serum Complement Levels and Immune Complexes: 8. Consultations: <ul style="list-style-type: none"> - Gastroenterology Consultation: - Infectious Disease Consultation: - Rheumatology Consultation: - Allergy/Immunology Consultation: <p>---</p> <p>Summary:</p> <p>A comprehensive diagnostic approach is essential to confirm the most likely diagnosis of Crohn's disease and to exclude other potential causes of the patient's symptoms. Early and accurate diagnosis will guide appropriate management and improve patient outcomes.</p>	<p>1</p>
--------------------	--	--	----------

<p>Case 30 2023</p>	<p>Dr. Nina B. Gold: The diagnostic test was measurement of the serum ammonia level. The serum specimen must be obtained from a free-flowing venipuncture without the use of a tourniquet, stored on ice, and sent immediately to the laboratory; deviations from this process may lead to falsely elevated results. This patient's serum ammonia level was 90 µmol per liter (reference range, 12 to 48) and remained elevated on repeat measurements, with a maximum level of 107 µmol per liter. She had encephalopathy during the time that these levels were obtained. The next step was an evaluation for causes of noncirrhotic hyperammonemia, including inherited metabolic disorders and acquired perturbations of the urea cycle.²⁰</p> <p>The plasma amino acid levels obtained during the patient's hospitalization were low, a finding consistent with malnutrition. The free and total plasma carnitine levels, plasma acylcarnitine profile, serum creatine kinase level, and urine organic acid levels obtained after the hospitalization were normal. She was found to have a mildly elevated urine orotic acid level, which can be associated with proximal urea-cycle dysfunction. However, repeat measurements of the plasma amino acid levels have not shown hallmarks of a proximal urea-cycle defect, such as an elevated glutamine level or a low citrulline level (Figure 2). Results of molecular testing for 76 monogenic disorders associated with hyperammonemia were normal, which makes a diagnosis of an inherited metabolic disorder, such as an ornithine transcarbamylase deficiency, unlikely.</p>	<ol style="list-style-type: none"> 1. **Fecal Alpha-1 Antitrypsin Clearance Test** 2. **Technetium-99m-Labeled Albumin Scintigraphy (Radionuclide Imaging)** 3. **Lymphangiography or Magnetic Resonance Lymphangiography** 4. **Serum Immunoglobulin Levels (IgG, IgA, IgM)** 5. **Electrolyte Panel with Focus on Phosphorus, Magnesium, and Potassium Levels** 6. **Repeat Electroencephalography (EEG)** 7. **Copper and Ceruloplasmin Levels** 8. **Comprehensive Nutritional Assessment** [vitamins and minerals, including B vitamins, vitamin D, and fat-soluble vitamins (A, E, K)]. 9. **Hydrogen/Methane Breath Test for Small Intestinal Bacterial Overgrowth (SIBO)** 10. **Jejunal Aspirate and Culture** (SIBO.) 11. **Autoimmune Panel Including Anti-Enterocyte and Anti-Goblet Cell Antibodies** 12. **Heavy Metal Screening (Blood Levels of Lead, Mercury, Arsenic)** 13. **Repeat Gastrointestinal Biopsies with Specialized Staining** (To detect rare conditions not identified on previous biopsies; Special stains (e.g., electron microscopy, immunohistochemistry) might reveal subtle pathologies like mitochondrial disorders or Whipple's disease). 14. Serum Prealbumin and Transferrin Levels** 15. **MRI of the Brain with Nutritional Protocols** (Certain deficiencies (e.g., vitamin B12, thiamine) can cause characteristic changes visible on MRI). 16. **Bone Marrow Biopsy** (for marrow infiltration or hematologic malignancies.) 17. Echocardiogram** 18. **Skin Biopsy (if there are skin changes)** (dermatologic manifestations of nutritional deficiencies or systemic diseases.) 19. **Assessment of Dietary Intake and Absorption Studies** (quantified with tests like D-xylose absorption test.) 20. **Genetic Testing for Congenital Disorders of Glycosylation or Lymphatic Development** <p>**Prioritization of Tests:**</p> <ul style="list-style-type: none"> - **First Line:** <ul style="list-style-type: none"> - **Fecal Alpha-1 Antitrypsin Clearance** - **Serum Immunoglobulin Levels** - **Electrolyte Panel (Phosphorus, Magnesium, Potassium)** - **Copper and Ceruloplasmin Levels** - **Comprehensive Nutritional Assessment** - **Second Line:** <ul style="list-style-type: none"> - **Technetium-99m-Labeled Albumin Scintigraphy** - **Lymphangiography** - **Hydrogen/Methane Breath Test** - **Repeat EEG** - **Heavy Metal Screening** - **Third Line:** <ul style="list-style-type: none"> - **Jejunal Aspirate and Culture** - **Autoimmune Antibody Panel** - **Specialized Gastrointestinal Biopsies** - **MRI Brain with Nutritional Protocols** 	<p>0 (incorrect)</p>
-------------------------	--	--	--------------------------

Table 2: Three examples of the o1-preview suggested testing plan compared to the testing plan conducted. One example scored a two, indicating that the test was appropriate and nearly identical to the case plan. A score of one indicates that the suggested diagnostics would have been helpful or yielded the diagnosis via another test not used in the case. A score of zero indicates that the suggested diagnostics would be unhelpful. Verbose rationales from o1-preview were abridged by a physician (Z.K.) to better fit in the table.

Table 3. Probabilistic Reasoning Before and After Testing by o1-preview

Case	Reference probability range, %	Predicted probability, median (IQR), %			MAE (MAPE)		
		o1-preview (n=100)	GPT-4 (n=100)	Clinician (n=553)	o1-preview (n=100)	GPT-4 (n=100)	Clinician (n=553)
Pneumonia							
Before test	25-42	70 (60-70)	72 (68.8-78)	80 (75-90)	31.6 (94.2)	39.5 (117.9)	47.7 (142.4)
After positive test result	46-65	94.9 (92.2-95)	95 (92.3-95)	95 (90-100)	37.5 (67.7)	38.5 (69.4)	38.5 (69.4)
After negative test result	10-19	28.5 (20-34.2)	25 (20-30)	50 (30-80)	14.8 (102.4)	10.7 (73.5)	39.9 (275.1)
Breast Cancer							
Before test	0.2-0.3	0.3 (0.1-0.5)	1.4 (1.4-1.5)	5 (1-10)	0.2 (97.9)	1.2 (466.4)	8.5 (3385.4)
After positive test result	3-9	2.8 (1.2-4)	7.8 (7-7.8)	50 (30-80)	3.5 (59)	4.4 (74)	47.6 (792.9)
After negative test result	<0.05	0.1 (0-0.1)	0.2 (0.1-0.3)	5 (1-10)	0.1 (260.3)	0.2 (750.6)	11.3 (45077.4)
Cardiac ischemia							
Before test	1-4.4	5 (2-5)	2.5 (2-5)	10 (5-20)	2.4 (87.3)	1.3 (47.2)	11.9 (439.5)
After positive test result	2-11	10.9 (6-13.5)	68.7 (65-70)	70 (50-90)	5.7 (87.1)	56.5 (869.1)	56.3 (865.9)
After negative test result	0.4-2.5	2 (1-2.2)	5 (3-5.6)	5 (1-10)	1 (69.7)	4 (273.9)	8.6 (587.9)
Urinary tract infection							
Before test	0-1	15 (10-20)	25.6 (20-30)	20 (10-50)	13.8 (2752)	26.2 (5241.7)	32.4 (6472.9)
After positive test result	0-8.3	90 (77-95)	90 (90-95)	80 (30-95)	78.2 (1885.5)	87.8 (2114.9)	62.2 (1499.6)
After negative test result	0-0.1	1 (1-2)	5 (5-10)	5 (0-10)	1.8 (3182.2)	7.3 (13294.2)	11.8 (21495)

Table 3: *: $p \leq 0.05$, **: $p \leq 0.01$, ***: $p \leq 0.001$, ****: $p \leq 0.0001$

References

1. Brodman, K., Erdmann, A. J., Jr, Lorge, I., Gershenson, C. P. & Wolff, H. G. The Cornell Medical Index-Health Questionnaire. III. The evaluation of emotional disturbances. *J. Clin. Psychol.* **8**, 119–124 (1952).
2. de Dombal, F. T., Leaper, D. J., Staniland, J. R., McCann, A. P. & Horrocks, J. C. Computer-aided diagnosis of acute abdominal pain. *Br. Med. J.* **2**, 9–13 (1972).
3. Shortliffe, E. H. Mycin: A knowledge-based computer program applied to infectious diseases. *Proc. Annu. Symp. Comput. Appl. Med. Care* 66–69 (1977).
4. Ing, E. B., Balas, M., Nassrallah, G., DeAngelis, D. & Nijhawan, N. The Isabel differential diagnosis generator for orbital diagnosis. *Ophthal. Plast. Reconstr. Surg.* **39**, 461–464 (2023).
5. Ledley, R. S. & Lusted, L. B. Reasoning foundations of medical diagnosis; symbolic logic, probability, and value theory aid our understanding of how physicians reason. *Science* **130**, 9–21 (1959).
6. Burkett, E. L. & Todd, B. R. A novel use of an electronic differential diagnosis generator in the emergency department setting. *Cureus* **15**, e34211 (2023).
7. Strong, E. *et al.* Chatbot vs medical student performance on free-response clinical reasoning examinations. *JAMA Intern. Med.* **183**, 1028–1030 (2023).
8. Kanjee, Z., Crowe, B. & Rodman, A. Accuracy of a Generative Artificial Intelligence Model in a Complex Diagnostic Challenge. *JAMA* **330**, 78–80 (2023).
9. Cabral, S. *et al.* Clinical Reasoning of a Generative Artificial Intelligence Model Compared With Physicians. *JAMA Intern. Med.* **184**, 581–583 (2024).
10. Savage, T., Nayak, A., Gallo, R., Rangan, E. & Chen, J. H. Diagnostic reasoning prompts reveal the potential for large language model interpretability in medicine. *NPJ Digit. Med.* **7**, 20 (2024).

11. Introducing OpenAI o1. <https://openai.com/index/introducing-openai-o1-preview/>.
12. OpenAI o1 System Card. <https://openai.com/index/openai-o1-system-card/>.
13. Nori, H. *et al.* From Medprompt to o1: Exploration of run-time strategies for medical challenge problems and beyond. *arXiv* (2024).
14. Learning to Reason with LLMs. <https://openai.com/index/learning-to-reason-with-llms/>.
15. Xie, Y. *et al.* A preliminary study of o1 in medicine: Are we closer to an AI doctor? *ArXiv abs/2409.15277*, (2024).
16. Hager, P. *et al.* Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nat. Med.* **30**, 2613–2622 (2024).
17. Abdulnour, R.-E. E. *et al.* Deliberate practice at the virtual bedside to improve clinical reasoning. *N. Engl. J. Med.* **386**, 1946–1947 (2022).
18. Schaye, V. *et al.* Development of a Clinical Reasoning Documentation Assessment Tool for Resident and Fellow Admission Notes: a Shared Mental Model for Feedback. *J. Gen. Intern. Med.* **37**, 507–512 (2022).
19. Goh, E. *et al.* Large language model influence on management reasoning: A randomized controlled trial. *medRxiv* (2024) doi:10.1101/2024.08.05.24311485.
20. Goh, E. *et al.* Large language model influence on diagnostic reasoning: A randomized clinical trial: A randomized clinical trial. *JAMA Netw. Open* **7**, e2440969 (2024).
21. Berner, E. S. *et al.* Performance of four computer-based diagnostic systems. *N. Engl. J. Med.* **330**, 1792–1796 (1994).
22. Morgan, D. J. *et al.* Accuracy of practitioner estimates of probability of diagnosis before and after testing. *JAMA Intern. Med.* **181**, 747–755 (2021).
23. Fritz, P. *et al.* Evaluation of medical decision support systems (DDX generators) using real medical cases of varying complexity and origin. *BMC Med. Inform. Decis. Mak.* **22**, 254 (2022).
24. Latif, E. *et al.* A systematic assessment of OpenAI o1-preview for higher order thinking in

- education. *arXiv* (2024).
25. Ratwani, R. M., Bates, D. W. & Classen, D. C. Patient safety and artificial intelligence in clinical care. *JAMA Health Forum* **5**, e235514 (2024).
 26. NIH findings shed light on risks and benefits of integrating AI into medical decision-making. *National Institutes of Health (NIH)*
<https://www.nih.gov/news-events/news-releases/nih-findings-shed-light-risks-benefits-integrating-ai-into-medical-decision-making> (2024).
 27. Jin, Q. *et al.* Hidden flaws behind expert-level accuracy of multimodal GPT-4 vision in medicine. *NPJ Digit. Med.* **7**, 190 (2024).
 28. Newman-Toker, D. E. *et al.* *Introduction*. (Agency for Healthcare Research and Quality, 2022).
 29. Auerbach, A. D. *et al.* Diagnostic errors in hospitalized adults who died or were transferred to intensive care. *JAMA Intern. Med.* **184**, 164–173 (2024).
 30. Nori, H. *et al.* Can generalist foundation models outcompete special-purpose tuning? Case study in medicine. *arXiv [cs.CL]* (2023).
 31. Singhal, K. *et al.* Towards expert-level medical question answering with large language models. *arXiv [cs.CL]* (2023).
 32. Saab, K. *et al.* Capabilities of Gemini models in medicine. *arXiv [cs.AI]* (2024).
 33. Cook, D. A., Sherbino, J. & Durning, S. J. Management reasoning. *JAMA* **319**, 2267 (2018).
 34. Buckley, T., Diao, J. A., Rajpurkar, P., Rodman, A. & Manrai, A. K. Multimodal Foundation Models Exploit Text to Make Medical Image Predictions. *arXiv preprint arXiv:2311.05591* (2024).
 35. Goldszmidt, M., Minda, J. P. & Bordage, G. Developing a unified list of physicians' reasoning tasks during clinical encounters. *Acad. Med.* **88**, 390–394 (2013).
 36. Bond, W. F. *et al.* Differential diagnosis generators: an evaluation of currently available computer programs. *J. Gen. Intern. Med.* **27**, 213–219 (2012).

37. McDuff, D. *et al.* Towards Accurate Differential Diagnosis with Large Language Models. *arXiv [cs.CY]* (2023).
38. Rodman, A., Buckley, T. A., Manrai, A. K. & Morgan, D. J. Artificial intelligence vs clinician performance in estimating probabilities of diagnoses before and after testing. *JAMA Netw. Open* **6**, e2347075 (2023).