

Performance of Multimodal GPT-4V on USMLE with Image: Potential for Imaging Diagnostic Support with Explanations

Zhichao Yang, MSc^{1*}; Zonghai Yao, MSc^{1*}; Mahbuba Tasmin, BSc¹; Parth Vashisht, BSc¹;
Won Seok Jang, RN, MSc²; Feiyun Ouyang, PhD²; Beining Wang, BSc³; Dan Berlowitz, MD^{4,5};
Hong Yu, PhD^{1,2,5,6}

Author Affiliations:

¹College of Information and Computer Science, University of Massachusetts Amherst, Amherst, MA, USA

²Miner School of Computer & Information Sciences, University of Massachusetts Lowell, Lowell, MA, USA

³Shanghai Medical College, Fudan University, Shanghai, China

⁴Department of Public Health, University of Massachusetts Lowell, Lowell, MA, USA

⁵Center for Biomedical and Health Research in Data Sciences, University of Massachusetts Lowell, Lowell, MA, USA

⁶Center for Healthcare Organization and Implementation Research, VA Bedford Health Care System, Bedford, MA, USA

These authors contributed equally *: Zhichao Yang, Zonghai Yao

Corresponding Author Information:

Hong Yu, PhD

University of Massachusetts Amherst

1 University Avenue

Lowell, MA, US

Phone: [1 508 612 7292](tel:15086127292)

Email: Hong_Yu@uml.edu

Main Figures: 1; Tables: 4

Abstract

Background: Using artificial intelligence (AI) to help clinical diagnoses has been an active research topic for more than six decades. Past research, however, has not had the scale and accuracy for use in clinical decision making. The power of large language models (LLMs) may be changing this. In this study, we evaluated the performance and interpretability of Generative Pre-trained Transformer 4 Vision (GPT-4V), a multimodal LLM, on medical licensing examination questions with images.

Methods: We used three sets of multiple-choice questions with images from United States Medical Licensing Examination (USMLE), USMLE question bank for medical students (AMBOSS), and Diagnostic Radiology Qualifying Core Exam (DRQCE) to test GPT-4V's accuracy and explanation quality. We compared GPT-4V with two other large language models, GPT-4 and ChatGPT, which cannot process images. We also assessed the preference and feedback of healthcare professionals on GPT-4V's explanations.

Results: GPT-4V achieved high accuracies on USMLE (86.2%), AMBOSS (62.0%), and DRQCE (73.1%), outperforming ChatGPT and GPT-4 by relative increase of 131.8% and 64.5% on average. GPT-4V was in the 70th - 80th percentile with AMBOSS users preparing for the exam. GPT-4V also passed the full USMLE exam with an accuracy of 90.7%. GPT-4V's explanations were preferred by healthcare professionals when it answered correctly, but they revealed several issues such as image misunderstanding, text hallucination, and reasoning error when it answered incorrectly.

Conclusion: GPT-4V showed promising results for medical licensing examination questions with images, suggesting its potential for clinical decision support. However, GPT-4V needs to improve its explanation quality and reliability for clinical use.

Keywords

Multimodality, Artificial Intelligence, Large Language Model, ChatGPT, GPT-4V, USMLE, Medical License Exam, Clinical Decision Support

1-2 sentence description:

AI models offer potential for imaging diagnostic support tool, but their performance and interpretability are often unclear. Here, the authors show that GPT-4V, a large multimodal language model, can achieve high accuracy on medical licensing exams with images, but also reveal several issues in its explanation quality.

Introduction

Using computers to help make clinical diagnoses and guide treatments has been a goal of artificial intelligence since its inception.¹ The adoption of electronic health record (EHR) systems by hospitals in the US has resulted in an unprecedented amount of digital data associated with patient encounters. Computer-assisted clinical diagnostic support system (CDSS) endeavors to enhance clinicians' decisions with patient information and clinical knowledge.² There is burgeoning interest in CDSS for enhanced imaging³, often termed radiomics, in various disciplines such as breast cancer detection⁴, covid detection⁵, diagnosing congenital cataracts⁶, and hidden fracture location⁷. For a decision to be trustworthy for clinicians, CDSS should not only make the prediction but also provide accurate explanations.⁸⁻¹⁰ However, most previous imaging CDSS offers only highlight areas deemed significant by AI,¹¹⁻¹⁵ providing limited insight into the explanation of the diagnosis.¹⁶

Large language models (LLMs) could generate explanations as they are trained with reinforcement learning from human feedback to follow user requests to explain a question. Typical LLM examples include Chat Generative Pre-trained Transformer (ChatGPT), a renowned chatbot released by OpenAI in October 2022, and its successor Generative Pre-trained Transformer 4 (GPT-4) in March 2023. The influence of ChatGPT is attributed to its conversational prowess and its performance, which approaches or matches human-level competence in cognitive tasks, spanning various domains including medicine.¹⁷ ChatGPT has achieved commendable results in the United States Medical Licensing Examinations, leading to discussions about the readiness of LLM applications for integration into clinical¹⁸⁻²⁰, and educational²¹⁻²³ environments.

One limitation of ChatGPT is that it may only read and generate text but is unable to process other data modalities, such as images. This limitation, known as the "single-modality," is a common issue among many LLMs.^{24,25} Advancements in multimodal LLM promise enhanced capabilities and integration with diverse data sources.^{26,27} OpenAI's recent introduction of GPT-4V has undeniably made strides toward bridging this divide. GPT-4V, a state-of-the-art multimodal LLM, is equipped with visual processing ability, granting it to understand and describe visual content.²⁸ By incorporating GPT-4V into current imaging CDSS, physicians can ask open-ended questions pertaining to a patient's medical evaluation - taking into account all available information including images, symptoms, lab results, allowing for an interactive experience where AI suggest both decision and explanation to support physicians.

However, the ability of GPT-4V to analyze medical images still remains unknown. GPT-4V must perform comparably to humans on assessments of medical knowledge and reasoning such that users have sufficient confidence in its responses. In this work, we aim to assess GPT-4V performance on medical licensing examination questions with images, as well as to analyze its explanation for healthcare professional interpretability.

Method

This cross-sectional study compared the performance between GPT-4V, GPT-4, and ChatGPT on medical licensing examination questions answering. This study also investigates the quality of GPT-4V explanation in answering these questions. The study protocol was deemed exempt by Institutional Review Board at the VA Bedford Healthcare System and informed consent was waived due to minimal risk to patients. This study was conducted in October 2023.

Medical Exam Data Collection

We obtained study questions from three sources. The United States Medical Licensing Examination (USMLE) consists of three steps required to obtain a medical license in the United States. The USMLE assesses a physician's ability to apply knowledge, concepts, and principles, which is critical to both health and disease management and is the foundation for safe, efficient patient care. The Step1, Step2 clinical knowledge(CK), Step3 of USMLE sample exam released from the National Board of Medical Examiners (NBME) consist of 119, 120, and 137 questions respectively. Each question contained multiple options to choose from. We then selected all questions with images, resulting in 19, 13, and 18 questions from Step1, Step2 CK, and Step3. Discipline includes but is not limited to radiology, dermatology, orthopedics, ophthalmology, and cardiology.

The sample exam only included limited questions with images. Thus, we further collected similar questions from a non-public available and registered required source: AMBOSS, a widely used question bank for medical students, which provides exam performance data given students' performance. The performance of past AMBOSS students enabled us to assess the comparative effectiveness of the model. For each question, AMBOSS associated an expert-written hint to tip the student to answer the question and a difficulty level that ranges from 1-5. Levels 1, 2, 3, 4, and 5 represent the easiest 20%, 20-50%, 50%-80%, 80%-95%, and 95%-100% of questions respectively.²⁹ We randomly selected 10 questions from each of the 5 difficulty levels. And we repeated this process for Step1, Step2 CK, and Step3. This resulted in a total number 150 of questions.

The Diagnostic Radiology Qualifying Core Exam (DRQCE), offered after 36 months of residency training, is an image-rich exam to evaluate a candidate's core fund of knowledge and clinical judgment across practice domains of diagnostic radiology. DRQCE is not publicly available and requires registration. We collected 26 questions with images from the preparation exam offered by the American Board of Radiology (ABR). Thus, we had a total of 226 questions with images from the three sources. To illustrate GPT-4V's potential as an imaging diagnostic support tool, we modified a patient case report³⁰ to resemble a typical "curbside consult" question between medical professionals.³¹

How to Answer Image Questions using GPT-4V Prompt

GPT-4V took image and text data as inputs to generate textual outputs. Given that input format (prompt) played a key role in optimizing model performance, we followed the standard prompting guidelines of the visual question-answering task. Specifically, we prompted GPT-4V by first adding the image, then appending context (i.e., patient information) and questions, and finally providing multiple-choice options, each separated by a new line. An example user prompt and GPT-4V response are shown in Figure 1. When multiple sub-images existed in the image, we uploaded multiple sub-images to GPT-4V. We did not use hint in the prompt unless otherwise specified. The response consists of the selected option as an answer, supported by a textual explanation to substantiate the selected decision. When using ChatGPT and GPT-4 models that cannot handle image data, images were omitted from the prompt. Responses were collected from the September 25, 2023 version of models. Each question was manually entered into the ChatGPT website independently (new chat window).

Evaluation Metrics

For answer accuracy, we evaluated the model's performance by comparing the model's choice with the correct choice provided by the exam board or question bank website. We defined accuracy as the ratio of the number of correct choices to the total number of questions.

We also evaluated the quality of the explanation by preference from 3 healthcare professionals (one medical doctor, one registered nurse, and one medical student). For each question from AMBOSS dataset (n=150), we asked the healthcare professionals to choose their preference between an explanation by GPT-4V, an explanation by an expert, or a tie.

Additionally, we also asked healthcare professionals to evaluate GPT-4V explanation from a sufficient and comprehensive perspective.^{32,33} They determined if the following information exists in the explanation:

1. Image interpretation: GPT-4V tried to interpret the image in the explanation, and such interpretation is sufficient to support its choice.
2. Question information: Explanations contained information related to the textual context (i.e., patient information) of the question, and such information was essential for GPT-4V's choice.
3. Comprehensive explanation: The explanation included comprehensive reasoning for all possible evidence (e.g., symptoms, lab results) that leads to the final answer.

Finally, for each question answered incorrectly, we asked healthcare professionals to check if the explanation contained any of the following errors:

1. Image misunderstanding: if the sentence in the explanation showed an incorrect interpretation of the image. Example: GPT-4V said that a bone in the image was for the hand, but it was the foot.
2. Text hallucination: if the sentence in the explanation contained something that is incorrect. (Other than the image) Example: Claiming Saxenda was just insulin.
3. Reasoning error: if the sentence did not properly convert the information (either image or text) to an answer. Example: GPT-4V identified a patient trip occurring in the last 3 months and despite the fact that chagas disease usually develops 10~20 years after infection, it still diagnosed the patient as having chagas disease.
4. Non-medical error: For any non-medical error, use this class. GPT is known to struggle with tasks requiring precise spatial localization, such as identifying chess positions on the board. It is known to struggle with calculations, such as $1 + 1 = ?$.

Statistical Analysis

GPT-4V's accuracies on the AMBOSS dataset were compared between different difficulties using unpaired chi-square tests with a significance level of 0.05. All analysis was conducted in Python software (version 3.10.11).

Results

Overall Answer Accuracy

For questions with image, GPT-4V achieved an accuracy of 84.2%, 85.7%, 88.9% in Step1, Step2CK, and Step3 of USMLE questions accordingly, outperforming ChatGPT and GPT-4 by 42.1% and 21.1% in Step1, 35.7% and 21.4% in Step2CK, 38.9% and 22.2% in Step3 (Table 1). Similarly, GPT-4V achieved an accuracy of 73.1%, outperforming ChatGPT (19.2%) and GPT-4 (26.9%) in DRQCE.

For all questions in the USMLE sample exam (including ones without image), GPT-4V achieved an accuracy of 88.2%, 90.8%, 92.7% in Step1, Step2CK, and Step3 of USMLE questions accordingly (Table 1), passing the standard for the USMLE (about 60%). However, it achieved limited accuracy in several medical disciplines such as anatomy (25.0%), emergency medicine (25.0%), and pathology (50.0%). A grasp of images is essential to correctly answer the majority of questions in these disciplines.

Accuracy Decreases When Difficulty Increases

When asking GPT-4V questions without the hint, it achieved an accuracy of 60%, 64%, and 66% for AMBOSS Step1, Step2CK, and Step3. GPT-4V was in the 72nd, 76th, and 80th percentile with AMBOSS users who were preparing for Step1, Step2CK, and Step3 respectively. Table 2 shows a decreasing trend in performance as question difficulty increased in the AMBOSS dataset ($P < 0.05$). However, the decreasing trend was not observed when the GPT-4V was

questioned with the hint. Out of 55 wrong answers without the hint, 17 were corrected by hints.

An example and detailed analysis are provided in the supplementary material figure 1.

Quality of Explanation

We first evaluated the user's preference among GPT-4V generated explanations and expert generated explanations. When GPT-4V answered incorrectly, it was no surprise that healthcare professionals overwhelmingly preferred expert explanations as shown in Table 3. When GPT-4V answered correctly, healthcare professionals favoring experts only exceeded favoring for GPT-4V by 4 votes, out of a total of 95 votes.

We further evaluated the quality of the GPT-4V generated explanation by verifying if explanation includes image and question text interpretation in Table 4. When examining the 95 correct answers, 84.2% (n=80) of the responses contained an interpretation of the image, while 96.8% (n=92) aptly captured the information presented in the question. On the other hand, for the 55 incorrect answers, 92.8% (n=51) interpreted the image, and 89.1% (n=49) depicted the question's details. In terms of explanation comprehensiveness, GPT-4V offered a comprehensive one in 79.0% (n=75) of correct responses. In contrast, only 7.2% (n=4) of the wrong responses had a comprehensive explanation that led to the final choice.

We also evaluated the explanations that lead to GPT-4V in answering incorrectly across 4 metrics as outlined above: image misunderstanding, text hallucination, reasoning error, and non-medical error. Among questions with wrong answers (n=55), we found that 76.3% (n=42) of questions included misunderstanding of the image, 45.5% (n=25) of questions included logic error, 18.2% (n=10) of questions included text hallucination, and no questions included non-medical errors.

Case Study on Consult Conversation

The consultation conversation, regarding a 45-Year-Old woman with hypertension, fatigue, and altered mental status, is provided in supplementary material figure 2. We found that the interactive design of GPT-4V allowed physicians to seek additional information by posing follow-up questions. Specifically, GPT-4V initially provided an irrelevant response when asked to interpret the CT scan. However, it was able to adjust its response and accurately identify the potential medical condition depicted in the image after receiving a physician's visual hint - an arrow pointed to a part of the CT scan where physicians desired GPT-4V to analysis.

Through comparing GPT-4V response with the case report, we also found that GPT-4V generally offered responses that were clear and coherent. When asked about differential diagnosis, GPT-4V explained why 3 diseases should be listed (Primary Aldosteronism, Hypertension, and Cushing's Syndrome) along with its explanation which were deemed relevant by a medical doctor. Following a query about the subsequent steps to ascertain the origin of the anomaly, GPT-4V recommended a PET-CT scan. Utilizing the patient's PET-CT scan, it was able to locate a tumor in the mediastinum, lending credence to the suspicion of Cushing's Syndrome. Finally, GPT-4V asked for further studies, such as a biopsy of the mass, to confirm the diagnosis.

Discussion

A prevailing research direction involves thoroughly investigating ChatGPT's ability to handle a diverse set of medical examination questions.^{21,34-36} GPT-4 was recently introduced as the upgraded model for ChatGPT. Studies showed that GPT-4 outperforms ChatGPT in various

medical tasks.^{31,37} The collective insights illustrated the power of ChatGPT and GPT-4 in medical exam answering and the potential for medical decision support. However, previous evaluations were only limited to questions without images.

In this study on the evaluation of medical exam questions with images, we found that GPT-4V selects more correct choices compared to ChatGPT and GPT-4 as shown in Table 1. Hence, when evaluating all questions in the USMLE sample exam, GPT-4V achieved an accuracy of 90.7% outperforming ChatGPT (58.5%) and GPT-4 (83.8%). The passing standard for the USMLE was typically set at 60%, indicating that the GPT-4V performed at a level similar to a medical graduate in the final year of study. The accuracy of GPT-4V highlights its grasp over biomedical and clinical sciences, essential for medical practice, but also showcases its ability in patient management and problem-solving skills,³⁸ both of which indicate the potential for clinical routines, such as summarizing radiology reports³⁹ and differential diagnosis^{40,41}. We further explored GPT-4V's potential in CDSS through a consultation conversation. By posing follow-up queries and providing visual hints from physicians, GPT-4V showed promise in CDSS applications, including interpreting CT scan interpretation, differential diagnoses, and follow-up exams recommendation. However, the quality of GPT-4V's explanation to support its clinical decision making remained an open question, bridging the focus of our study to the next segment of our analysis.

In terms of explanation quality, we found that more than 80% of responses from GPT-4V provided an interpretation of the image and question of its answer selection, regardless of correctness. This suggested that GPT-4V consistently took into account both the image and question elements while generating responses. Figure 1 illustrates an example of high-quality explanation that utilizes both text and image in answering a hard question. More than 70% of students answered incorrectly on the first try, because both bacterial pneumonia and pulmonary

embolism may involve symptoms such as cough. To differentiate them, GPT-4V correctly interpreted the X-ray with a radiologic sign of Hampton hump, which further increased the suspicion of pulmonary infarction rather than pneumonia.⁴² To show the need for X-ray as mentioned in the explanation, we removed the image from the input, and GPT-4V switched the answer to bacterial pneumonia while also acknowledging the possibility of pulmonary infarction. This change in response demonstrated the high quality of the GPT-4V explanation, as its explanation about X-ray was not fictional and it truly needed the X-ray to answer this question. The high quality of GPT-4V explanations was also supported by experts' preference voting. When comparing explanations generated by experts with ones generated by GPT-4V, experts' preference for the expert over the GPT-4V was minimal (n=4) when GPT-4V correctly answered the question (n=95). Previous studies have shown limited utilization of CDSS as most of them offered limited decision explanation and thus gained limited trust among physicians (unlike their colleagues).⁴³⁻⁴⁶ Thus, GPT-4V could enhance the effectiveness and trustworthiness of CDSS by providing high-quality, expert-preferred explanations, encouraging broader adoption and more confident utilization among physicians.

We also found that the accuracy of GPT-4V was related to the comprehensiveness of the explanation. GPT-4V offered a comprehensive explanation in 79.0% of correct responses. In contrast, only 7.2% of the wrong answers had a comprehensive explanation. Thus, the absence of key information may be the cause of inaccurate answers. These observations suggested that enhancing the performance of models can be achieved by training GPT-4V with more intricate, clinical-specific data with insights from experienced physicians such as UpToDate,⁴⁷ and medical research literature such as PubMed.²⁶

Image misunderstanding was the primary reason why GPT-4V answered incorrectly. Out of 55 wrong responses, 42 (76.3%) were due to misunderstanding of the image. In comparison, only

10 (18.2%) of the mistakes were attributed to text misinterpretation. GPT-4V's proficiency in processing images was considerably lagging behind its text-handling capability. To circumvent its image interpretation issue, we tried to additionally prompt GPT-4V with a short hint that described the image. We found that 40.5% (17 out of 42) responses switched to the correct answer. Corrections from the hint indicated that GPT-4V could be easily persuaded. Within a conversational interface, medical professionals can readily guide and refine GPT-4V's initial outputs. This adaptability could be advantageous for physicians, as it allows for real-time adjustments and ensures that the generated information aligns more closely with the clinical context or the specific details of a patient's case.⁴⁷ With customized hint from physicians, GPT-4V enhances the utility and reliability as an auxiliary tool.

Another drawback of GPT-4V involved its tendency to produce factually inaccurate responses, a problem often referred to as the hallucination effect, which is prevalent among many large language models such as GPT-4V.⁴⁸ We found that more than 18% of GPT-4V explanations contain hallucinations. Thus when designing clinical support tools for high-risk situations such as patient diagnosing, it is crucial to integrate GPT-4V and a probabilistic model with confidence interval, indicating the reliability of the response.⁴⁹ This would enhance the reliability of the CDSS response when additional physician review is warranted.¹⁶

Limitations

This study has several limitations. First, our findings are constrained in their applicability due to the modest sample size. We gathered 226 questions that included images, which might not comprehensively represent all medical disciplines. Second, while GPT-4V has demonstrated proficiency in medical license examination, its real-world applicability, especially in dynamic,

user-interactive scenarios, remains untested. Therefore, while the results are promising, extrapolating the efficacy of GPT-4V to broader clinical applications requires appropriate benchmarks and further research.²⁰

Conclusion

While GPT-4V showcased remarkable accuracy across a spectrum of medical disciplines and varying difficulty levels in this study, it is paramount that further refinement be undertaken, particularly in enhancing its explanatory capabilities prior to any clinical assimilation. Medical students and professionals must be acutely aware of its limitations and consistently cross-verify with authoritative sources. Notably, state-of-the-art LLMs like GPT-4V, even with their sophisticated capabilities, are merely on the brink of supplanting physicians. Their performance in specialized examinations, though noteworthy, still exhibits imperfections, which can lead to consequential inaccuracies and uncertainties. Coupled with the well-known ethical concerns, the credibility and readiness of LLMs for clinical settings remain under scrutiny. However, the preliminary results are promising, suggesting that the present technology is poised to influence clinical practices. As research and development persist, we anticipate a more extensive and profound integration of AI in the medical domain.

Reference

1. Shortliffe EH, Cimino JJ. *Biomedical Informatics: Computer Applications in Health Care and Biomedicine*. Springer; 2014.
2. Sutton RT, Pincock D, Baumgart DC, Sadowski DC, Fedorak RN, Kroeker KI. An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ Digital Medicine*. 2020;3.
3. Rajpurkar P, Lungren MP. The Current and Future State of AI Interpretation of Medical Images. *The New England journal of medicine*. 2023;388 21:1981-1990.

4. Aggarwal R, Sounderajah V, Martin G, et al. Diagnostic accuracy of deep learning in medical imaging: a systematic review and meta-analysis. *NPJ Digital Medicine*. 2021;4. <https://api.semanticscholar.org/CorpusID:233139020>
5. Wang L, Lin ZQ, Wong A. COVID-Net: a tailored deep convolutional neural network design for detection of COVID-19 cases from chest X-ray images. *Scientific Reports*. 2020;10. <https://api.semanticscholar.org/CorpusID:215768886>
6. Long E, Lin H, Liu Z, et al. An artificial intelligence platform for the multihospital collaborative management of congenital cataracts. *Nature Biomedical Engineering*. 2017;1. <https://api.semanticscholar.org/CorpusID:113460889>
7. Rayan JC, Reddy N, Kan JH, Zhang W, Annapragada AV. Binomial Classification of Pediatric Elbow Fractures Using a Deep Learning Multiview Approach Emulating Radiologist Decision Making. *Radiology Artificial intelligence*. 2019;1 1:e180015.
8. Bussone A, Stumpf S, O'Sullivan D. The Role of Explanations on Trust and Reliance in Clinical Decision Support Systems. *2015 International Conference on Healthcare Informatics*. Published online 2015:160-169.
9. Panigutti C, Beretta A, Giannotti F, Pedreschi D. Understanding the impact of explanations on advice-taking: a user study for AI-based clinical Decision Support Systems. *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. Published online 2022. <https://api.semanticscholar.org/CorpusID:248419322>
10. Gaube S, Suresh H, Raue M, et al. Non-task expert physicians benefit from correct explainable AI advice when reviewing X-rays. *Scientific reports*. 2023;13(1):1383.
11. Singh A, Mohammed AR, Zelek JS, Lakshminarayanan V. Interpretation of deep learning using attributions: application to ophthalmic diagnosis. In: *Optical Engineering + Applications*. ; 2020. <https://api.semanticscholar.org/CorpusID:221616930>
12. Eitel F, Ritter K. Testing the Robustness of Attribution Methods for Convolutional Neural Networks in MRI-Based Alzheimer's Disease Classification. In: Suzuki K, Reyes M, Syeda-Mahmood T, et al., eds. *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support*. Springer International Publishing; 2019:3-11.
13. Papanastasiopoulos Z, Samala RK, Chan HP, et al. Explainable AI for medical imaging: deep-learning CNN ensemble for classification of estrogen receptor status from breast MRI. In: *Medical Imaging*. ; 2020. <https://api.semanticscholar.org/CorpusID:216291456>
14. Shamout FE, Shen Y, Wu N, et al. An artificial intelligence system for predicting the deterioration of COVID-19 patients in the emergency department. *NPJ Digital Medicine*. 2021;4. <https://api.semanticscholar.org/CorpusID:220968946>
15. Pei Y, Wang G, Cao H, et al. A deep-learning pipeline to diagnose pediatric intussusception and assess severity during ultrasound scanning: a multicenter retrospective-prospective study. *NPJ Digital Medicine*. 2023;6. <https://api.semanticscholar.org/CorpusID:263264783>

16. Shen Y, Heacock L, Elias J, et al. ChatGPT and Other Large Language Models Are Double-edged Swords. *Radiology*. Published online 2023:230163.
17. OpenAI. GPT-4 Technical Report. *ArXiv*. 2023;abs/2303.08774. <https://api.semanticscholar.org/CorpusID:257532815>
18. Goodman RS, Patrinely JR, Stone J Cosby A, et al. Accuracy and Reliability of Chatbot Responses to Physician Questions. *JAMA Network Open*. 2023;6(10):e2336483-e2336483. doi:10.1001/jamanetworkopen.2023.36483
19. Decker H, Trang K, Ramirez J, et al. Large Language Model–Based Chatbot vs Surgeon-Generated Informed Consent Documentation for Common Procedures. *JAMA Network Open*. 2023;6. <https://api.semanticscholar.org/CorpusID:263774434>
20. Ayers JW, Poliak A, Dredze M, et al. Comparing Physician and Artificial Intelligence Chatbot Responses to Patient Questions Posted to a Public Social Media Forum. *JAMA internal medicine*. Published online 2023. <https://api.semanticscholar.org/CorpusID:258375371>
21. Kung TH, Cheatham M, Medenilla A, et al. Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digital Health*. 2022;2.
22. Thirunavukarasu AJ, Hassan R, Mahmood S, et al. Trialling a Large Language Model (ChatGPT) in General Practice With the Applied Knowledge Test: Observational Study Demonstrating Opportunities and Limitations in Primary Care. *JMIR Medical Education*. 2023;9. <https://api.semanticscholar.org/CorpusID:258259005>
23. Cooper AZ, Rodman A. AI and Medical Education - A 21st-Century Pandora's Box. *The New England journal of medicine*. Published online 2023. <https://api.semanticscholar.org/CorpusID:260322445>
24. Khader F, Müller-Franzes G, Wang T, et al. Multimodal Deep Learning for Integrating Chest Radiographs and Clinical Parameters: A Case for Transformers. *Radiology*. 2023;309 1:e230806.
25. Topol EJ. As artificial intelligence goes multimodal, medical applications multiply. *Science*. 2023;381 6663:adk6139.
26. Zhang S, Xu Y, Usuyama N, et al. Large-Scale Domain-Specific Pretraining for Biomedical Vision-Language Processing. *ArXiv*. 2023;abs/2303.00915. <https://api.semanticscholar.org/CorpusID:257280046>
27. Tu T, Azizi S, Driess D, et al. Towards Generalist Biomedical AI. *ArXiv*. 2023;abs/2307.14334. <https://api.semanticscholar.org/CorpusID:260164663>
28. Yang Z, Li L, Lin K, et al. The Dawn of LMMs: Preliminary Explorations with GPT-4V(ision). *ArXiv*. 2023;abs/2309.17421. <https://api.semanticscholar.org/CorpusID:263310951>
29. AMBOSS. AMBOSS Question difficulty. Published 10/15/2023. <https://support.amboss.com/hc/en-us/articles/360035679652-Question-difficulty>

30. Pallais JC, Fenves AZ, Lu MT, Glomski K. Case 18-2018: A 45-Year-Old Woman with Hypertension, Fatigue, and Altered Mental Status. *The New England journal of medicine*. 2018;378 24:2322-2333.
31. Lee P, Bubeck S, Petro J. Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine. *The New England journal of medicine*. 2023;388 25:2399.
32. Yu M, Chang S, Zhang Y, Jaakkola T. Rethinking Cooperative Rationalization: Introspective Extraction and Complement Control. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics; 2019:4094-4103. doi:10.18653/v1/D19-1420
33. Zaidan O, Eisner J, Piatko C. Using “Annotator Rationales” to Improve Machine Learning for Text Categorization. In: *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*. Association for Computational Linguistics; 2007:260-267. <https://aclanthology.org/N07-1033>
34. Bhayana R, Krishna S, Bleakney RR. Performance of ChatGPT on a Radiology Board-style Examination: Insights into Current Strengths and Limitations. *Radiology*. Published online 2023:230582.
35. Gilson A, Safranek CW, Huang T, et al. How Does ChatGPT Perform on the United States Medical Licensing Examination? The Implications of Large Language Models for Medical Education and Knowledge Assessment. *JMIR Medical Education*. 2023;9. <https://api.semanticscholar.org/CorpusID:256663603>
36. Brin D, Sorin V, Vaid A, et al. Comparing ChatGPT and GPT-4 performance in USMLE soft skill assessments. *Scientific Reports*. 2023;13(1):16492. doi:10.1038/s41598-023-43436-9
37. Nori H, King N, McKinney SM, Carignan D, Horvitz E. Capabilities of GPT-4 on Medical Challenge Problems. *ArXiv*. 2023;abs/2303.13375. <https://api.semanticscholar.org/CorpusID:257687695>
38. The Federation of State Medical Boards (FSMB) and the National Board of Medical Examiners® (NBME®). Step 3 - United States Medical Licensing Examination. Published online 2023. <https://www.usmle.org/step-exams/step-3>
39. Elkassem AMA, Smith AD. Potential Use Cases for ChatGPT in Radiology Reporting. *AJR American journal of roentgenology*. Published online 2023. <https://api.semanticscholar.org/CorpusID:258003533>
40. Hirose T, Harada Y, Yokose M, Sakamoto T, Kawamura R, Shimizu T. Diagnostic Accuracy of Differential-Diagnosis Lists Generated by Generative Pretrained Transformer 3 Chatbot for Clinical Vignettes with Common Chief Complaints: A Pilot Study. *International Journal of Environmental Research and Public Health*. 2023;20. <https://api.semanticscholar.org/CorpusID:256936867>

41. Shea YF, Lee CMY, Ip WCT, Luk DWA, Wong SSW. Use of GPT-4 to Analyze Medical Records of Patients With Extensive Investigations and Delayed Diagnosis. *JAMA Network Open*. 2023;6. <https://api.semanticscholar.org/CorpusID:260885460>
42. Patel UB, Ward TJ, Kadoch MA, Cham MD. Radiographic features of pulmonary embolism: Hampton's hump. *Postgraduate Medical Journal*. 2014;90:420-421.
43. Liberati EG, Ruggiero F, Galuppo L, et al. What hinders the uptake of computerized decision support systems in hospitals? A qualitative study and framework for implementation. *Implementation Science*: IS. 2017;12. <https://api.semanticscholar.org/CorpusID:9726465>
44. Strohm L, Hehakaya C, Ranschaert ER, Boon WPC, Moors EHM. Implementation of artificial intelligence (AI) applications in radiology: hindering and facilitating factors. *European Radiology*. 2020;30:5525-5532.
45. Cauwenberge DV, Biesen W van, Decruyenaere JM, Leune T, Sterckx S. "Many roads lead to Rome and the Artificial Intelligence only shows me one road": an interview study on physician attitudes regarding the implementation of computerised clinical decision support systems. *BMC Medical Ethics*. 2022;23. <https://api.semanticscholar.org/CorpusID:248547001>
46. Jones C, Thornton J, Wyatt JC. Artificial intelligence and clinical decision support: clinicians' perspectives on trust, trustworthiness, and liability. *Medical law review*. Published online 2023. <https://api.semanticscholar.org/CorpusID:258844404>
47. Lourenco AP, Slanetz PJ, Baird GL. Rise of ChatGPT: It May Be Time to Reassess How We Teach and Test Radiology Residents. *Radiology*. Published online 2023:231053.
48. Ji Z, Lee N, Frieske R, et al. Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys*. 2022;55:1-38.
49. Jiang S, Xu YY, Lu X. ChatGPT in Radiology: Evaluating Proficiencies, Addressing Shortcomings, and Proposing Integrative Approaches for the Future. *Radiology*. 2023;308 1:e231335.

Tables

Table 1. Performance of ChatGPT, GPT-4, and GPT-4V on USMLE sample exam from NBME.

Performance on USMLE Sample Exam - Step1	questions with Image (n=19)	all questions (n=119)
ChatGPT	42.1%	55.1%
GPT-4	63.2%	81.5%
GPT-4V	84.2%	88.2%
Passing	-	~60%
Performance on USMLE Sample Exam - Step2CK	questions with Image (n=14)	all questions (n=120)
ChatGPT	50.0%	59.1%
GPT-4	64.3%	80.8%
GPT-4V	85.7%	90.8%
Passing	-	~60%
Performance on USMLE Sample Exam - Step3	questions with Image (n=18)	all questions (n=137)
ChatGPT	50.0%	60.9%
GPT-4	66.7%	88.3%
GPT-4V	88.9%	92.7%
Passing	-	~60%

Table 2. Performance of GPT-4V on USMLE AMBOSS. For each step, overall: n=50; difficulty 1: n=10; difficulty 2: n=10; difficulty 3: n=10; difficulty 4: n=10; difficulty 5: n=10.

GPT-4V accuracy on AMBOSS-Step1						
	Overall	Question difficulty level				
		1	2	3	4	5
Without Hint	60%	70%	70%	30%	70%	60%
Expert Hint	84%	80%	80%	80%	90%	90%
GPT-4V accuracy on AMBOSS-Step2CK						
	Overall	Question difficulty level				
		1	2	3	4	5
Without Hint	64%	80%	70%	70%	50%	50%
Expert Hint	86%	100%	90%	100%	70%	70%
GPT-4V accuracy on AMBOSS-Step3						
	Overall	Question difficulty level				
		1	2	3	4	5
Without Hint	66%	80%	90%	60%	50%	50%
Expert Hint	88%	90%	90%	90%	90%	80%

Table 3. Healthcare professionals preferred explanation for 150 AMBOSS questions.

	Correct (n=95)			Incorrect (n=55)		
	prefer expert	balance	prefer GPT-4V	prefer expert	balance	prefer GPT-4V
Step1	4	23	3	16	4	0
Step2CK	10	15	7	18	0	0
Step3	5	23	5	13	4	0

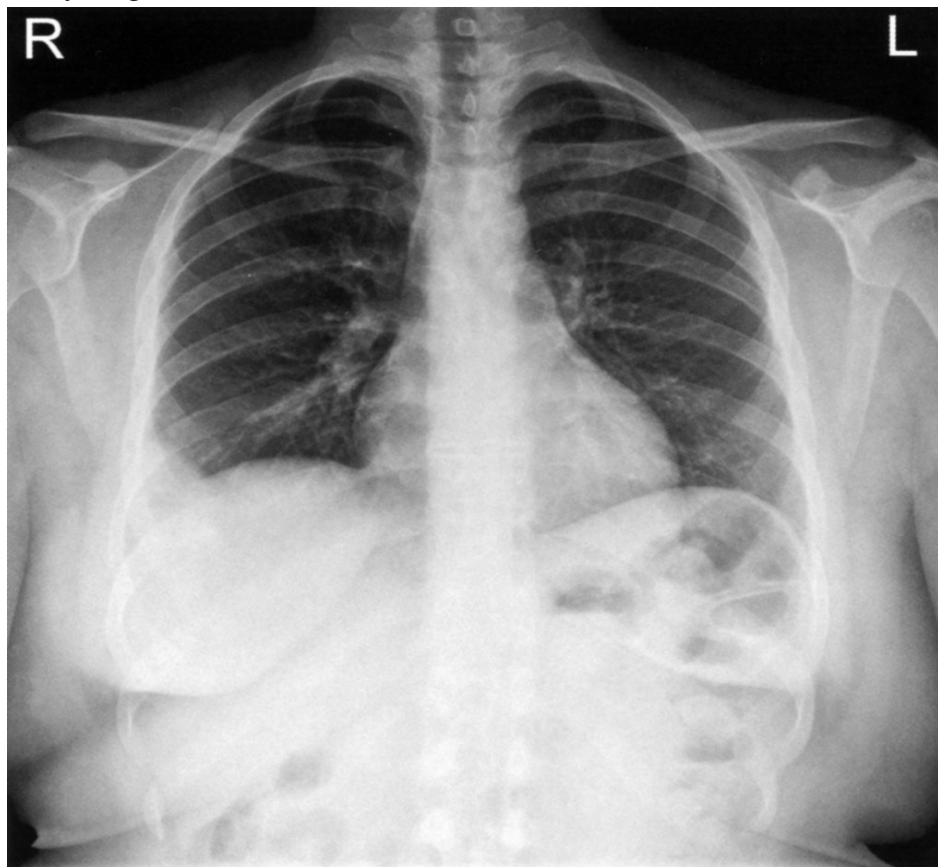
Table 4. Qualitative analysis of GPT-4V response quality for AMBOSS questions. Given each exam question, healthcare professionals evaluate GPT-4V explanations for the following criteria. Image Interpretation: GPT-4V sufficiently supported its choice through its interpretation of the image; Question Information: The explanation provided essential patient information from the textual question context; Comprehensive Explanation: The explanation included comprehensive reasoning for all possible that leads to the final answer.

	Correct (n=95)	Incorrect (n=55)
Image Interpretation	84.2%	92.8%
Question Information	96.8%	89.1%
Comprehensive Explanation	79.0%	7.2%

Figures

Figure 1. A case study of GPT-4V in answering USMLE questions with chest x-ray. Input prompt and output response are separated by "GPT Response Below". Option D is the correct answer and is selected by GPT-4V without hint.

Question: A 55-year-old woman comes to the physician because of a 4-day history of chest pain and coughing up blood. The chest discomfort is characterized as piercing and aggravated by coughing. She encountered symptoms of a sore throat and nasal congestion 11 days prior. She was diagnosed with multiple sclerosis 15 years ago and utilizes a wheelchair for movement. A 40-year history of smoking a pack of cigarettes daily is noted, but she abstains from alcohol. Her current prescriptions include ocrelizumab and dantrolene. Temperature is 37.9°C (100.2°F), pulse is 105/min, and blood pressure is 110/60 mm Hg. Lung auscultation reveals occasional inspiratory crackles in the right lower lobe, while the cardiac exam shows no abnormalities. Neurologic examination shows stiffness and decreased sensation of the lower extremities; there is diffuse hyperreflexia. An x-ray of the chest is provided. Which of the following is the most likely diagnosis?



Options:

A: Sarcoidosis

- B: Pericarditis
- C: Bacterial pneumonia
- D: Pulmonary embolism
- E: Pneumothorax

Let's think step by step. First, analyze each option one by one, and then give your final prediction and reasons to be sure we have the right answer.

GPT Response Below:

Let's go through the information provided and evaluate each option based on the clinical picture and the chest x-ray:

****Clinical Information**:**

- 55-year-old woman
- Smoked for 40 years (20 pack-years)
- Recent upper respiratory tract symptoms
- 3-day history of chest pain and cough with blood-tinged sputum
- Chest pain is sharp, stabbing, and worsened with coughing
- History of multiple sclerosis and uses a wheelchair
- On ocrelizumab and dantrolene
- Vital signs show mild fever and elevated pulse
- Auscultation shows crackles in the right lower lung

****Chest X-ray**:**

Hampton hump is a radiologic sign which consists of a shallow wedge-shaped opacity in the periphery of the lung.

****Options**:**

A: ****Sarcoidosis**** - Sarcoidosis often presents with bilateral hilar lymphadenopathy on chest x-ray. The given image doesn't seem to show this. Additionally, the acute symptoms don't align well with a diagnosis of sarcoidosis.

B: ****Pericarditis**** - This would present with chest pain that's positional and relieved by sitting forward. The x-ray provided doesn't suggest pericardial effusion. Unlikely.

C: ****Bacterial pneumonia**** - Given her symptoms and the crackles in the right lower lung, this is a possibility. However, the chest x-ray would typically show a lobar consolidation or infiltrate which isn't clearly seen here.

D: **Pulmonary embolism** - Given her history of immobility (uses a wheelchair) and the acute onset of sharp, stabbing chest pain exacerbated by coughing, this is a strong contender. The presence of a Hampton Hump should further increase suspicion of a PE.

E: **Pneumothorax** - This would present with sudden onset chest pain and shortness of breath. On x-ray, you'd expect to see a clear line indicating the collapsed lung and no lung markings beyond that line. This is not evident in the provided image.

Conclusion:

Based on the provided information and the chest x-ray, the most likely diagnosis is:

D: Pulmonary embolism.

The reason for this choice is the combination of her risk factors (immobility, smoking), her clinical presentation (sharp, stabbing chest pain exacerbated by coughing, blood-tinged sputum), and the fact that the x-ray doesn't clearly support the other diagnoses.

Conflicts of Interest

The authors declare no conflict of interests.