

Vision-based ChatGPT Shows Deficits Interpreting Radiologic Images

Released: September 3, 2024

OAK BROOK, Ill. — Researchers evaluating the performance of ChatGPT-4 Vision found that the model performed well on text-based radiology exam questions but struggled to answer image-related questions accurately. The study's results were published today in *Radiology*, a journal of the Radiological Society of North America (RSNA).

Chat GPT-4 Vision is the first version of the large language model that can interpret both text and images.

[download photo](#)



Chad Klochko, M.D.

“ChatGPT-4 has shown promise for assisting radiologists in tasks such as simplifying patient-facing radiology reports and identifying the appropriate protocol for imaging exams,” said Chad Klochko, M.D., musculoskeletal radiologist and artificial intelligence (AI) researcher at Henry Ford Health in Detroit, Michigan. “With image processing capabilities, GPT-4 Vision allows for new potential applications in radiology.”

For the study, Dr. Klochko’s research team used retired questions from the American College of Radiology’s Diagnostic Radiology In-Training Examinations, a series of tests used to benchmark the progress of radiology residents. After excluding duplicates, the researchers used 377 questions across 13 domains, including 195 questions that were text-only and 182 that contained an image.

GPT-4 Vision answered 246 of the 377 questions correctly, achieving an overall score of 65.3%. The model correctly answered 81.5% (159) of the 195 text-only queries and 47.8% (87) of the 182 questions with images.

“The 81.5% accuracy for text-only questions mirrors the performance of the model’s predecessor,” he said. “This consistency on text-based questions may suggest that the model has a degree of textual understanding in radiology.”

Genitourinary radiology was the only subspecialty for which GPT-4 Vision performed better on questions with images (67%, or 10 of 15) than text-only questions (57%, or 4 of 7). The model performed better on text-only questions in all other subspecialties.

The model performed best on image-based questions in the chest and genitourinary subspecialties, correctly answering 69% and 67% of the image-containing questions, respectively. The model performed lowest on image-containing questions in the nuclear medicine domain, correctly answering only 2 of 10 questions.

The study also evaluated the impact of various prompts on the performance of GPT-4 Vision.

- Original: You are taking a radiology board exam. Images of the questions will be uploaded. Choose the correct answer for each question.
- Basic: Choose the single best answer in the following retired radiology board exam question.
- Short instruction: This is a retired radiology board exam question to gauge your medical knowledge. Choose the single best answer letter and do not provide any reasoning for your answer.
- Long instruction: You are a board-certified diagnostic radiologist taking an examination. Evaluate each question carefully and if the question additionally contains an image, please evaluate the image carefully in order to answer the question. Your response must include a single best answer choice. Failure to provide an answer choice will count as incorrect.
- Chain of thought: You are taking a retired board exam for research purposes. Given the provided image, think step by step for the provided question.

Although the model correctly answered 183 of 265 questions with a basic prompt, it declined to answer 120 questions, most of which contained an image.

“The phenomenon of declining to answer questions was something we hadn’t seen in our initial exploration of the model,” Dr. Klochko said.

The short instruction prompt yielded the lowest accuracy (62.6%).

On text-based questions, chain-of-thought prompting outperformed long instruction by 6.1%, basic by 6.8%, and original prompting style by 8.9%. There was no evidence to suggest performance differences between any two prompts on image-based questions.

“Our study showed evidence of hallucinatory responses when interpreting image findings,” Dr. Klochko said. “We noted an alarming tendency for the model to provide correct diagnoses based on incorrect image interpretations, which could have significant clinical implications.”

Dr. Klochko said his study’s findings underscore the need for more specialized and rigorous evaluation methods to assess large language model

performance in radiology tasks.

“Given the current challenges in accurately interpreting key radiologic images and the tendency for hallucinatory responses, the applicability of GPT-4 Vision in information-critical fields such as radiology is limited in its current state,” he said.

“Performance of GPT-4 with Vision on Text- and Image-based ACR Diagnostic Radiology In-Training Examination Questions.” Collaborating with Dr. Klochko were Nolan Hayden, M.D., Spencer Gilbert, B.S., Laila M. Poisson, Ph.D., and Brent Griffith, M.D.

Radiology is edited by Linda Moy, M.D., New York University, New York, N.Y., and owned and published by the Radiological Society of North America, Inc. (<https://pubs.rsna.org/journal/radiology>)

RSNA is an association of radiologists, radiation oncologists, medical physicists and related scientists promoting excellence in patient care and health care delivery through education, research and technologic innovation. The Society is based in Oak Brook, Illinois. ([RSNA.org](https://www.rsna.org))

For patient-friendly information on medical imaging, visit [RadiologyInfo.org](https://radiologyinfo.org).

Images (JPG, TIF):

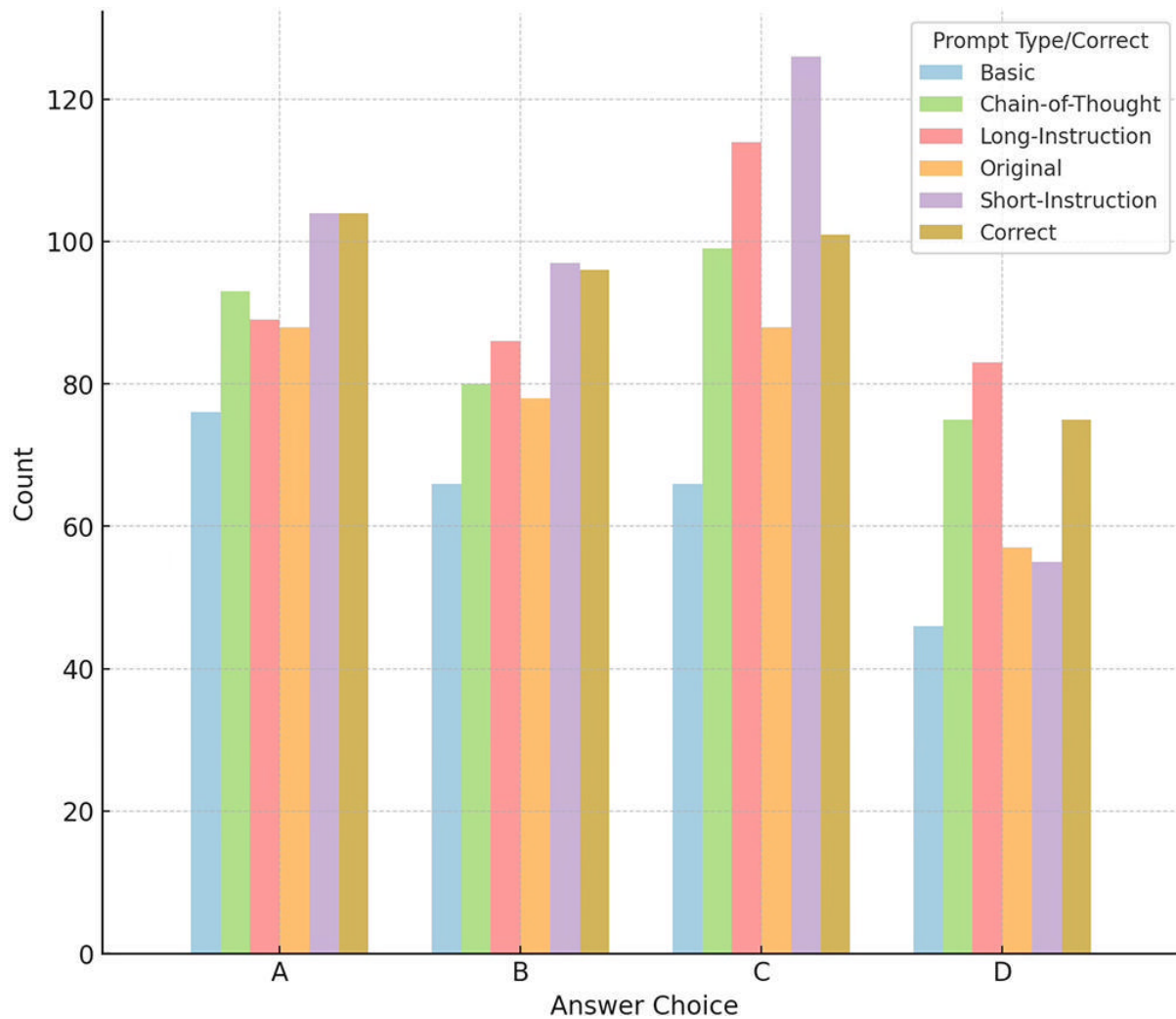
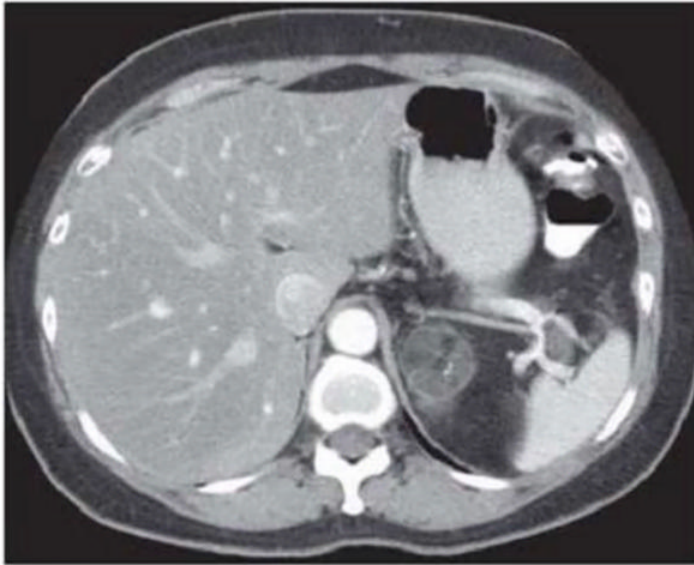


Figure 1. Bar graph shows the distribution of answers provided by GPT-4 with vision (GPT-4V; OpenAI), categorized by prompt phrase, as well as the distribution of correct answers.
[High-res \(TIF\) version](#)

5. Based upon this abdominal CT image from a 39-year-old woman, what is the MOST likely diagnosis?



- A. Adenoma
- B. Lymphangioma
- C. Metastasis
- D. Myelolipoma



The image shows an abdominal CT with a well-defined, heterogeneous lesion in the liver that has areas of fat attenuation and denser regions, which is suggestive of soft tissue or hematopoietic elements.

Given the appearance and characteristics of the lesion, the MOST likely diagnosis is:

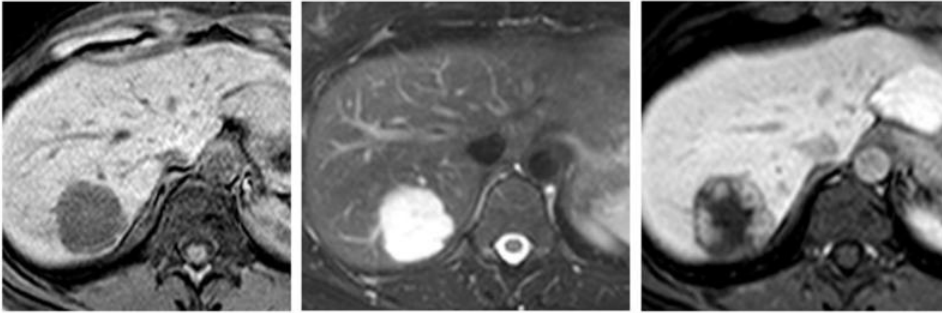
D. Myelolipoma

Myelolipomas are rare benign tumors that primarily arise in the adrenal glands but can occasionally be found elsewhere, such as in the liver. They are composed of mature adipose tissue and hematopoietic elements.

Figure 2. Screenshot of an example question from the American College of Radiology Diagnostic Radiology In-Training Examinations and response of GPT-4 with vision (GPT-4V; OpenAI). In this example, GPT-4V does not address the obvious left adrenal gland lesion and describes a nonexistent hepatic lesion, yet still arrives at the correct answer.

[High-res \(TIF\) version](#)

1. This 44-year-old woman is being evaluated for a focal liver lesion detected on an abdominal sonogram. Based on the provided pre- and post- contrast MRI images, what is the MOST likely diagnosis?



- A. Focal nodular hyperplasia
- B. Hemangioma
- C. Hepatocellular carcinoma
- D. Hepatic adenoma

From the provided pre- and post-contrast MRI images, we can observe a liver lesion that demonstrates iso- to hypointense signal on the T1-weighted image (left) and hyperintense signal on the T2-weighted image (middle). The lesion also shows enhancement in the arterial phase after contrast administration (right).

Given the imaging characteristics:

- A. Focal nodular hyperplasia (FNH) typically has a central scar and shows iso- to hypointense signal on T1 and hyperintense on T2. It enhances in the arterial phase. However, the central scar is not clearly evident in this case.
- B. Hemangioma would typically show a nodular peripheral enhancement pattern in the arterial phase, filling in centripetally. The T2 hyperintense signal is consistent with hemangioma.
- C. Hepatocellular carcinoma (HCC) is more commonly associated with cirrhosis and often shows arterial phase enhancement with washout in the portal venous or delayed phases, which is not clearly demonstrated here.
- D. Hepatic adenoma is typically T1-hypointense and T2-hyperintense, but they don't usually show the typical enhancement pattern as seen in this lesion.

Based on the provided images, the MOST likely diagnosis is:

B. Hemangioma.

Figure 3. Screenshot of an example question from the American College of Radiology Diagnostic Radiology In-Training Examinations and response of GPT-4 with vision (GPT-4V; OpenAI). In this example, GPT-4V provides reasoning for its correct answer. [High-res \(TIF\) version](#)

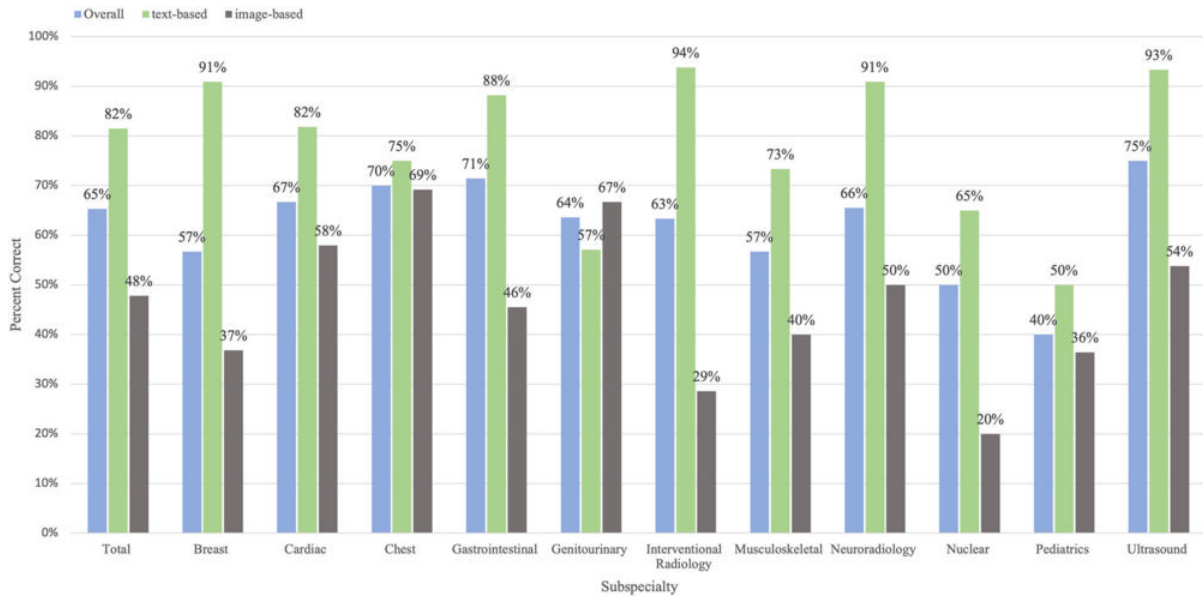


Figure 4. Bar graph shows the performance of GPT-4 with vision (GPT-4V; OpenAI) by subspecialty and question type. The only subspecialty where GPT-4V performed better in questions with images compared with text-only questions was in the genitourinary category. In all other subspecialties, ChatGPT performed better on text-only questions. The general and physics categories were excluded from the graph because these categories did not include any image-based questions.

[High-res \(TIF\) version](#)

Resources:

[Editorial](#)

[Study abstract](#)