

GPT-4V(ision) Unsuitable for Clinical Care and Education: A Clinician-Evaluated Assessment

Senthujan Senkaiahliyan M. Mgt, Augustin Toma MD, Jun Ma PhD, An-Wen Chan MD, Andrew Ha MD, Kevin R. An MD, Hrishikesh Suresh MD, Barry Rubin MD, and Bo Wang PhD

Abstract

OpenAI's large multimodal model, GPT-4V(ision), was recently developed for general image interpretation. However, less is known about its capabilities with medical image interpretation and diagnosis. Board-certified physicians and senior residents assessed GPT-4V's proficiency across a range of medical conditions using imaging modalities such as CT scans, MRIs, ECGs, and clinical photographs. Although GPT-4V is able to identify and explain medical images, its diagnostic accuracy and clinical decision-making abilities are poor, posing risks to patient safety. Despite the potential that large language models may have in enhancing medical education and delivery, the current limitations of GPT-4V in interpreting medical images reinforces the importance of appropriate caution when using it for clinical decision-making.



1. INTRODUCING GPT-4V(ISION)

This past year, large language models (LLMs) demonstrated impressive capabilities to perform numerous language-based tasks. They have shown capability in analyzing text, discerning patterns, and establishing connections between words [1]. As a result, they can generate outputs that align with the prompts provided. While LLMs have expressed strong performance in expert-level medical question answering, they are still unable to outperform their clinician counterparts especially in scenarios that require reasoning capabilities [2].

Generative Pre-Trained Transformer Vision (GPT-4V) is OpenAI's first large multimodal model with the ability to accept image input alongside text. [3] Multimodal learning is the ability for machine learning models to be trained on and input multiple forms of input data. They have the potential to enhance the breadth and depth of tasks that LLMs can perform across various medical disciplines. [4]

To evaluate GPT-4V's proficiency in analyzing medical images, we conducted an evaluation involving senior residents and board-certified physicians to assess its capability to accurately interpret various medical conditions and provide accurate and useful information regarding the diagnosis and management of these conditions. The study aimed to assess whether GPT-4V could not only interpret medical images but also provide valuable information for diagnosis, management, and education. Finally, we aimed to evaluate if the resulting outputs align with the safety standards for patient care.

2. DATA COLLECTION

2.1 General Conditions

In the data collection phase, a diverse set of multimodal medical images were gathered to assess the performance of GPT-4V across various medical scenarios and specialties. The breakdown of multimodal images is presented in Table 1, showcasing different modalities and their respective counts. These images were sourced from open-source libraries and repositories found on the internet.

-
- Senthujan Senkaiahliyan is with the Institute for Health Policy Management and Evaluation, Faculty of Public Health, University of Toronto and Peter Munk Cardiac Centre, University Health Network, Toronto ON, Canada.
 - Augustin Toma is with the Department of Medical Biophysics, Faculty of Medicine, University of Toronto, Toronto, ON, Canada.
 - Jun Ma is with Peter Munk Cardiac Centre, University Health Network; Department of Laboratory Medicine and Pathobiology, University of Toronto; Vector Institute, Toronto, ON Canada.
 - An-Wen Chan is with the Institute for Health Policy Management and Evaluation, Faculty of Public Health and with the Division of Dermatology, Department of Medicine, University of Toronto, Toronto, ON, Canada
 - Andrew Ha is with Peter Munk Cardiac Centre, University Health Network and the Division of Cardiology, Department of Medicine, University of Toronto, Toronto, ON, Canada
 - Kevin R. An is with the Division of Cardiac Surgery, Department of Surgery, University of Toronto, Toronto, ON, Canada
 - Hrishikesh Suresh is with the Division of Neurosurgery, Department of Surgery, University of Toronto, Toronto, ON, Canada
 - Barry Rubin is with Peter Munk Cardiac Centre, University Health Network and the Division of Vascular Surgery, Department of Surgery, University of Toronto, Toronto, ON, Canada
 - Bo Wang (Corresponding Author) is with Peter Munk Cardiac Centre, University Health Network; Department of Laboratory Medicine and Pathobiology and Department of Computer Science, University of Toronto; Vector Institute, Toronto, Canada. E-mail: bowang@vectorinstitute.ai

TABLE 1
Breakdown of Multimodal Images.

Modality	Count	Source
Clinical Photos*	10	Various different websites
Computed Tomography – Abdomen (CT-Abdomen)	3	https://www.eurorad.org/
Computed Tomography – Chest (CT-Chest)	5	https://www.eurorad.org/
Chest X-Ray (CXR)	10	https://www.eurorad.org/
Computed Tomography- Head (CT-Head)	5	https://www.eurorad.org/
Electrocardiogram (ECG)	10	https://ecg.bidmc.harvard.edu/maven/mavenmain.asp
Magnetic Resonance Imaging- Brain (MRI-Brain)	3	Various different websites
Abdominal X-ray (AXR)	2	https://www.eurorad.org/
Flow-Volume	3	https://medschool.co/
Musculoskeletal X-ray (MSK-Xray)*	3	Various different websites
Electroencephalogram (EEG)	4	https://www.eegatlas-online.com/index.php/en/
Esophagogastroduodenoscopy (EGD) *	2	Various different websites
Fundoscopy	3	https://stanfordmedicine25.stanford.edu/

2.2 Cardiology

The dataset used was a set of ECG waveforms sourced from the ECG Wave-Maven: A Self-Assessment Program for Students and Clinicians¹. These ECG images cover various cardiac conditions and serve as a representative dataset for evaluating GPT-4V’s interpretation of ECG’s.

2.3 Dermatology

In dermatology, clinical photos were collected from the Hellenic Dermatological Atlas², to curate a comprehensive set of dermatological conditions for assessing GPT-4V’s performance in interpretation.

3. EXPERIMENTAL SETUP

The methodology employed for this comprehensive evaluation followed a structured four-phase approach.

3.1 Dataset Curation

A diverse range of medical images and corresponding labels were selected from public datasets, encompassing various diagnostic modalities such as patient clinical photos, radiological images, ECG traces, EEG, fundoscopy, endoscopy, and colonoscopy. GPT-4V analyzed these images based on the prompts. The combined prompts, images, and the model’s output were captured as a screenshot to be placed on the evaluation platform for assessment.

3.2 Evaluation Criteria

A dual approach was adopted to assess the accuracy and reliability of GPT-4V’s interpretations. All images were evaluated by two senior surgical residents (K.R.A, H.S.) and a board-certified internal medicine physician (A.T.). ECGs and clinical photos of dermatologic conditions were additionally evaluated by a board-certified cardiac electrophysiologist (A.H.) and dermatologist (A.C.), respectively.

The following below are the questionnaires used for the evaluation

General Conditions (Diverse Modalities):

- 1) Rate the answer from 1-5.
- 2) Rate from 1-5 how comfortable you would be letting a medical student rely on this content to help learning.
- 3) Was the image interpreted correctly? (Yes/No)
- 4) Was the advice correct? (Yes/No)
- 5) Was the advice given dangerous? (Yes/No)

Cardiology (ECGs):

- 1) Rate the overall interpretation of the ECG (1-5).
- 2) Compared to a standard automated read of an ECG, would you consider this interpretation more competent? (Yes/No)
- 3) Rate from 1-5 how comfortable you would be letting a medical student rely on this content to help learning.
- 4) Would this interpretation be helpful in a medical student’s learning? (Yes/No)
- 5) General Comments:

1. <https://ecg.bidmc.harvard.edu/maven/mavenmain.asp>

2. <http://www.hellenicdermatlas.com/en/>

Dermatology (Clinical Photos):

- 1) Rate the quality of the layman's description of the rash (1-5)
- 2) Rate the quality of the medical description of the rash (1-5)
- 3) Rate the quality of the differential diagnosis (1-5)
- 4) General Comments

3.3 Evaluation Platform

Figure 1 provides a visual representation of the evaluation platform, designed as a web-based form.

- **Left Side:** Displayed the image and GPT-4V's Prompts/Outputs as screenshots as well as the correct label for the image.
- **Right Side:** Comprised a section for evaluator feedback, including questions from the evaluation criteria.

The screenshot displays a web-based evaluation platform. On the left, there is a section for the radiological image and its associated text. At the top right, the page number '1/69' is visible. The right side of the platform contains a section for 'Evaluation Criteria' with five numbered questions and a 'Comment Box for Clinicians' at the bottom. A footnote at the bottom right states '*Subject to change per image'.

1/69

Evaluation Criteria*

1. Rate the answer from 1-5.
2. Rate from 1-5 how comfortable you would be letting a medical student rely on this content to help learning.
3. Was the image interpreted correctly? (Yes/No)
4. Was the advice correct? (Yes/No)
5. Was the advice given dangerous? (Yes/No)

Comment Box for Clinicians

*Subject to change per image

PROMPT: You are an expert radiology tutor. Describe to a medical student what the major abnormality on this radiological image is.

Correct Diagnostic Label

GPT 4V OUTPUT: Of course! This is an axial slice from a non-contrast CT scan of the brain. Here's a basic interpretation for a medical student:

The most notable feature in this CT image is the absence of the usual symmetry between the two cerebral hemispheres. There is a clear hypodense (darker) area on one side of the brain. This hypodense area is consistent with an area of decreased tissue density, which often represents edema (swelling) or infarct (an area of dead tissue due to lack of blood supply).

There is also evidence of midline shift. The midline structures, such as the falx cerebri and the ventricles, are pushed away from the hypodense area. This is indicative of increased pressure on one side, often due to swelling or a mass effect. Such a shift can compress vital brain structures and is a concerning finding.

Based on this image alone, the major abnormality seems consistent with a large acute or subacute infarct (stroke), possibly due to a blockage in a major cerebral artery. This has resulted in cerebral edema and a significant midline shift.

Fig. 1. Evaluation platform to collect clinician feedback on GPT4V's output.

3.4 Prompts

The prompts that were used by GPT-4V to generate outputs are listed below for each category.

General Conditions (Diverse Modalities):

- "You are an expert [insert specialty i.e. radiology, cardiology] tutor. Explain to a medical student what this patient photo indicates."
- "If indicators like arrows, asterisks, or circles are present in the image, underline and expand on their significance."

Cardiology (ECGs):

- "You are an expert cardiology tutor assisting a medical student. Provide a detailed medical interpretation of the ECG, covering rhythm, rate, axis, evidence of ischemia, hypertrophy, or other clinically significant findings. Finally, list a differential diagnosis based on the ECG findings."

Dermatology (Clinical Photos):

- "You are an expert dermatology tutor helping a medical student. Describe the rash seen in the photo in layman's terms. Next, describe it using medical terminology. Finally, list a differential diagnosis for the given image."

4. RESULTS

4.1 Performance on Multimodal Images

For multimodal images (Table 2), a total of 69 images were assessed. Several images were accompanied by multiple prompts, with each undergoing a separate assessment. The correct diagnostic label for all these images were provided to the clinician evaluator to ensure accuracy in assessment. Clinician evaluators were asked to identify if GPT-4V correctly interpreted the images and whether they felt that the interpretation given was correct and safe for patient care. The average comfort level the clinicians felt about letting medical students learn from these images was 1.8 ± 1.4 on a scale of 1-5. Out of the 69 images, only 15 were correctly interpreted with the correct advice. However, there were a concerning number of instances (30 out of 69) where dangerous advice was provided. The images spanned various modalities (Table 1), including CT scans of various body parts, ECG, MRI, CXR, and others.

TABLE 2
Multimodal Images Summary of Results.

Parameter	Results
Total Number of Images	69
Average Comfort Level for Medical Students (Scale: 1-5)	1.8
Correctly Interpreted Images (Yes/No)	15 /69
Correct Advice Given (Yes/No)	16 /69
Was Dangerous Advice Provided? (Yes/No)	30 /69

4.2. Performance on Electrocardiograms (Cardiology)

For ECG images (Table 3), 24 images were examined. The overall interpretation of these images had an average rating of 2.25 ± 1.07 out of 5. Notably, none of these interpretations matched the competence of standard automated ECG reads as determined by the cardiac electrophysiologist. Out of the 24, only 3 responses were considered helpful for medical student learning, and in 9 cases, dangerous advice for patient care was given.

TABLE 3
ECG Summary of Results

Parameter	Results
Total Number of Images	24
Average Rating for Overall Interpretation of the ECG (Scale: 1-5)	2.25
Competence Compared to Standard Automated ECG Read (Yes/No)	0/24
Responses Considered Helpful for Medical Student Learning (Yes/No)	3/24
Was Dangerous Advice Provided? (Yes/No)	9/24

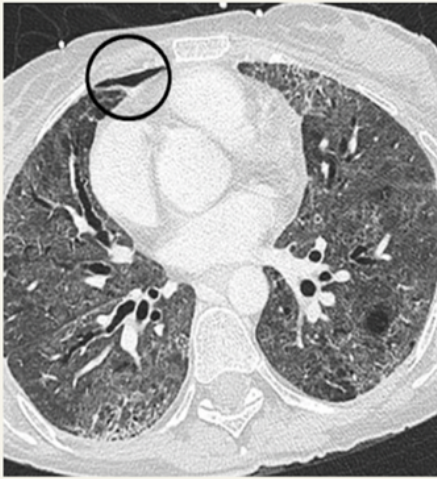
4.3 Performance on Clinical Photos (Dermatology)

For dermatology images (Table 4), out of the 49 images, the average quality of layman's description of the rash was 3 ± 1.55 out of 5. The medical descriptions and differential diagnoses of the rash averaged at 2.5 ± 1.49 and 2 ± 1.46 out of 5, respectively. The comfort level of using GPT-4V as an education tool for medical students averaged at 2 ± 1.4 out of 5. In addition, the differential diagnosis was described by the dermatologist as lacking depth and containing inaccuracies or irrelevant conditions.

TABLE 4
Clinical Photos Summary of Results

Parameter	Results
Total Number of Images	49
Average Quality of Layman's Description of the Rash (Scale: 1-5)	3 out of 5
Average Quality of Medical Description of the Rash (Scale: 1-5)	2.5 out of 5
Average Quality of Differential Diagnosis (Scale: 1-5)	2 out of 5

Figure 5 highlights direct examples of GPT-4V responses to images used in the evaluation along with clinician comments. For both cases highlighted, clinician comments indicate that GPT-4V has provided inaccurate advice that can impact patient care.



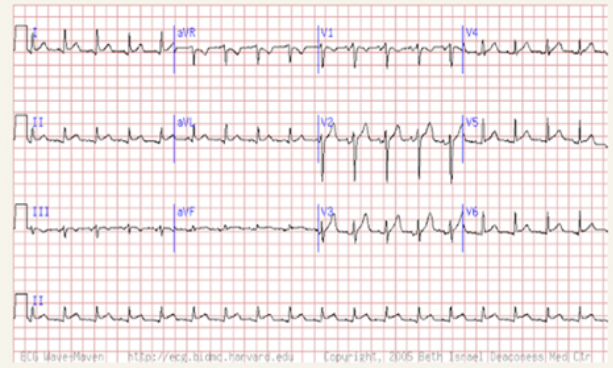
Correct Label: *Small pneumothorax*

Modality: *CT- Chest*

Prompt: *“You are a helpful radiology tutor, explain to a medical student what the abnormality in the circle is”*

Output: *Given the image and the information provided, this opacity is consistent with a pulmonary nodule. Pulmonary nodules can arise from various etiologies, both benign and granulomas*

“Missed a pneumothorax, thought it was a nodule? It does not look like a nodule at all”



Correct Label: *Acute Pericarditis*

Modality: *ECG*

Prompt: *“You are an expert cardiology tutor assisting a medical student. Provide a detailed medical interpretation of the ECG”*

Output: *Given this ECG, there aren't any overt pathologies noted. The differential would focus more on normal variants and conditions that might not manifest with clear ECG changes. Some considerations might include:*

- 1. Normal ECG in an asymptomatic patient.*
- 2. Early stages of a condition not yet manifesting on the ECG.*
- 3. Non-cardiac causes of symptoms (if the patient is symptomatic).*

“The interpretation missed the presence of diffuse ST elevation (and the most likely diagnosis of pericarditis).”



Fig. 2. Evaluation of GPT-4V's Interpretations on Medical Images with Expert Feedback

5. DISCUSSION AND LIMITATIONS

While GPT-4V demonstrates moderate proficiency in processing diverse medical imaging modalities and identifying specific features, it is important to note that the model occasionally falls short in recognizing overt findings. In addition, it's important to consider that the public-facing version of GPT-4V, as part of alignment efforts to not explicitly provide directives, may have impacted its performance on certain medical tasks.

Nevertheless, this evaluation of GPT-4V is not without its limitations. Firstly, our utilization of public-facing images, which might have potentially been part of the model's training datasets, should, in theory, have augmented its performance. Yet, GPT-4V's performance, especially with these images, were poor. This raises concerns about the depth and diversity of its training dataset. Secondly, as we provided GPT-4V with standalone images devoid of a broader clinical context, we

expected clinicians to consider this aspect when evaluating the model’s efficacy. It should be emphasized that diagnoses are not formed solely on a single picture and, in the absence of patient history, GPT-4V’s output should be evaluated with this consideration in mind.

The most glaring concern lies in the model’s accuracy, particularly with ECG interpretations. Instances where GPT-4V misinterprets severe conditions as benign poses significant risk for patient care. Without insight on the training datasets, a comprehensive evaluation will need to be conducted to uncover any harms in misrepresentation or potential bias. From our evaluation of GPT-4V’s performance, it’s evident that proprietary LLMs should strongly consider aligning with open-source principles. This is particularly crucial as many healthcare institutions are exploring collaborations with them for deployment in clinical and operational environments [5]. The Department of Health and Human Services within the United States is spearheading initiatives in this area, emphasizing the necessity for diverse and representative training data to ensure the ethical application of AI [6].

While LLMs have showcased the capability to tailor their responses based on user input and changing contexts, it’s noteworthy that our assessment was conducted during GPT-4V’s initial selective release. Since then, it appears that guardrails have been implemented to ensure that responses related to medical images remain generalized and descriptive rather than prescriptive.

Newer LLMs are being designed to address specific challenges within the medical field. An exemplar of this is Clinical Camel, a model that has been fine-tuned with medical datasets to enhance its performance significantly when addressing clinical inquiries, surpassing the capabilities of its pre-trained model [7]. With these developments, there’s an untapped potential for these models to become multimodal, offering a chance to develop comprehensive tools that support healthcare professionals provided they undergo thorough evaluation and validation in real-world clinical settings.

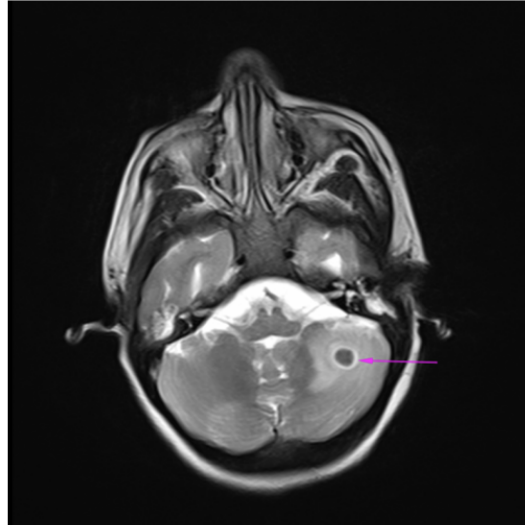
Considering the enthusiasm around Large Language Models (LLMs) and the suggestion that they will revolutionize the medical sphere, in our view GPT-4V’s current performance fails to offer merit to those claims. Our human evaluation substantiates healthcare regulatory bodies and OpenAI’s own advice on not using it as a substitute for clinician-based decision making [3]. While GPT-4V’s functionality as a multimodal foundation model—capable of processing both text and image inputs—is noteworthy, in its current form, significant concerns remain regarding its diagnostic accuracy and ability to interpret various medical image modalities.

REFERENCES

- [1] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting, “Large language models in medicine,” *Nature Medicine*, vol. 29, no. 8, pp. 1930–1940, 2023.
- [2] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfohl *et al.*, “Large language models encode clinical knowledge,” *Nature*, vol. 620, no. 7972, pp. 172–180, 2023.
- [3] OpenAI, “Gpt-4v(ision) system card,” 2023.
- [4] J. N. Acosta, G. J. Falcone, P. Rajpurkar, and E. J. Topol, “Multimodal biomedical ai,” *Nature Medicine*, vol. 28, no. 9, pp. 1773–1784, 2022.
- [5] A. J. Nashwan, A. A. Abujaber, and A. Abujaber, “Harnessing the power of large language models (llms) for electronic health records (ehrs) optimization,” *Cureus*, vol. 15, no. 7, 2023.
- [6] B. Meskó and E. J. Topol, “The imperative for regulatory oversight of large language models (or generative ai) in healthcare,” *npj Digital Medicine*, vol. 6, no. 1, p. 120, 2023.
- [7] A. Toma, P. R. Lawler, J. Ba, R. G. Krishnan, B. B. Rubin, and B. Wang, “Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding,” *arXiv preprint arXiv:2305.12031*, 2023.

SUPPLEMENTARY NOTES

Below are additional case studies from the evaluation highlighting examples of GPT-4V's output and comments from the evaluators.



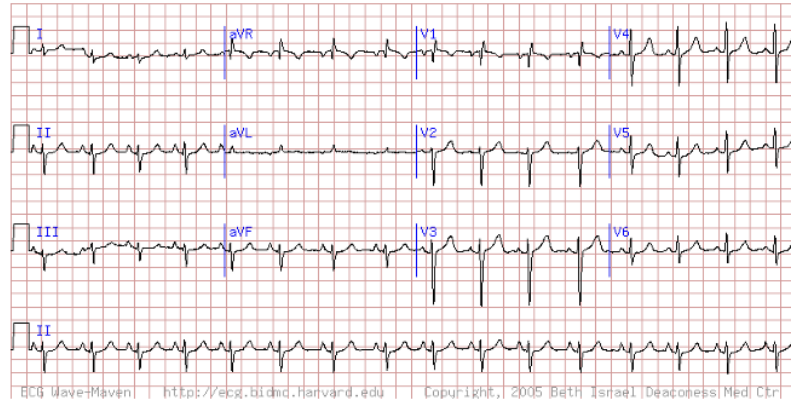
Correct Label:	Well-defined hypointense lesion with hyperintense rim with surrounding moderate perilesional oedema
Prompt	You are an expert radiology tutor. Describe to a medical student what the major abnormality on this radiological image is
GPT- 4V Output	<p>The image provided appears to be a T2-weighted axial MRI of the brain. The most notable abnormality is the large, well-defined, hyperintense (bright on T2-weighted images) lesion in the posterior fossa, which seems to be compressing the fourth ventricle. The lesion is located centrally within the cerebellum. Given its location, size, and signal characteristics, one would consider the possibility of a cerebellar tumor. A common tumor in this location, especially in pediatric patients, is a medulloblastoma, but the exact diagnosis would depend on further imaging studies, clinical context, and ultimately, histopathological examination.</p> <p>The ventricular system (seen as the dark, fluid-filled spaces in the brain) might also be enlarged (hydrocephalus) secondary to obstruction of the fourth ventricle by the mass. This can lead to an increase in intracranial pressure.</p> <p>It's essential for the medical students to understand that while imaging provides valuable insights, the full clinical picture, including patient symptoms, age, and other factors, should be <u>taken into account</u> when making a diagnosis or management decision.</p>
Clinician Comment	<i>It's just right about tumour, all else is off</i>

Fig. 3. Case Study 1- MRI



Correct Label: Mediastinal free air	
Prompt	You are an expert radiology tutor. Describe to a medical student what the asterisks on this radiological image are
GPT- 4V Output	<p>Certainly! This is a cross-sectional CT (computed tomography) image of the thorax, viewed in the axial plane.</p> <p>The asterisks are highlighting:</p> <ol style="list-style-type: none"> 1. The trachea (in the middle at the top): It is the main airway that connects the larynx to the bronchi, allowing air to enter the lungs. On the CT, the trachea appears as a round to oval structure with air inside, making it appear dark/black, and surrounded by a thin white rim which represents the tracheal wall. 2. The main bronchi (on either side just below the trachea): These are the primary divisions of the trachea that lead into the right and left lungs. Like the trachea, the bronchi are also filled with air and appear dark/black with a surrounding thin white rim. 3. The heart (in the center at the bottom): The structure below, appearing more centrally located and with a different texture than the lungs, represents the heart. The asterisk in this context might be highlighting the upper part of the heart or the great vessels. <p>In a clinical context, radiologists often use annotations or markers like these asterisks to draw attention to structures or potential areas of concern. In this educational setting, they are likely being used to highlight and identify key anatomical structures for the medical student.</p>
Clinician Comment	You can't miss free air in the mediastinum like that

Fig. 4. Case Study 2- CT



Prompt	You are an expert cardiology tutor assisting a medical student. Provide a detailed medical interpretation of the ECG. Your interpretation should include an assessment of the rhythm, rate, axis, and any evidence of ischemia, hypertrophy, or other clinically significant findings. Finally, list a differential diagnosis based on the ECG findings.
GPT- 4V Output	<p>1. Rate: To determine the rate, we can use the 300-rule. Count the number of large squares between two consecutive R waves and divide 300 by that number. Based on the image, the rate is around 75-80 beats per minute.</p> <p>2. Rhythm: The P waves appear regular and have a consistent morphology. The PR interval appears consistent as well. QRS complexes follow each P wave without any dropped beats, which suggests a normal sinus rhythm.</p> <p>3. Axis: The QRS complex is predominantly positive in lead I and negative in aVF. This suggests a leftward axis (between -30 and -90 degrees).</p> <p>4. Ischemia: There is ST segment depression seen in leads V4-V6. This can be suggestive of subendocardial ischemia.</p> <p>5. Hypertrophy:</p> <ul style="list-style-type: none"> • Left Ventricular Hypertrophy (LVH) criteria: Increased R wave amplitude in the left-sided chest leads (V5-V6) and deep S waves in the right-sided chest leads (V1-V3). This ECG does show some increase in R wave amplitude in V5 and V6 and deep S waves in V1-V3, which can be suggestive of LVH. • There's no clear evidence of Right Ventricular Hypertrophy (RVH) on this ECG. <p>6. Other clinically significant findings:</p> <ul style="list-style-type: none"> • No clear evidence of bundle branch block or pre-excitation. • The T waves appear upright and normal in most leads, except for the T wave inversions in V4-V6 which could be related to the ischemia or strain pattern. <p>Differential Diagnosis Based on the ECG Findings:</p> <ol style="list-style-type: none"> 1. Subendocardial ischemia. 2. Left Ventricular Hypertrophy, possibly due to conditions like hypertension. 3. Cardiac strain pattern. 4. Non-specific T wave changes which can be seen in a variety of conditions including electrolyte abnormalities.
Clinician Comment	The actual rate is about 96 bpm, not 75-80 bpm. Estimation of the axis is not correct as the actual axis is more in keeping with right (or borderline) axis deviation, and not left axis deviation. The interpretation called out ischemia which appears to not be the case.

Fig. 5. Case Study 3- ECG