

The Illusion of Readiness: Stress Testing Large Frontier Models on Multimodal Medical Benchmarks

Yu Gu^{†*}, Jingjing Fu^{*}, Xiaodong Liu, Jeya Maria Jose Valanarasu, Noel Codella, Reuben Tan, Qianchu Liu, Ying Jin, Sheng Zhang, Jinyu Wang, Rui Wang, Lei Song, Guanghui Qin, Naoto Usuyama, Cliff Wong, Cheng Hao, Hohin Lee, Praneeth Sanapathi, Sarah Hilado, Bian Jiang, Javier Alvarez-Valle, Mu Wei, Jianfeng Gao, Eric Horvitz, Matt Lungren, Hoifung Poon, Paul Vozila

Microsoft Research, Health & Life Sciences.

Abstract

Large frontier models like GPT-5 now achieve top scores on medical benchmarks. But our stress tests tell a different story. Leading systems often guess correctly even when key inputs like images are removed, flip answers under trivial prompt changes, and fabricate convincing yet flawed reasoning. These aren't glitches; they expose how today's benchmarks reward test-taking tricks over medical understanding. We evaluate six flagship models across six widely used benchmarks and find that high leaderboard scores hide brittleness and shortcut learning. Through clinician-guided rubric evaluation, we show that benchmarks vary widely in what they truly measure yet are treated interchangeably, masking failure modes. We caution that medical benchmark scores do not directly reflect real-world readiness. If we want AI to earn trust in healthcare, we must demand more than leaderboard wins and must hold systems accountable for robustness, sound reasoning, and alignment with real medical demands.

*Equal Technical Contribution

[†]Corresponding author: Yu Gu (aiden.gu@microsoft.com)

1 Introduction: The Illusion of Readiness in Health AI

Health AI has a credibility problem. Models like GPT-5 [1] can ace medical exams and top multimodal benchmarks yet still falter on simpler tasks, such as preserving their answers when answer choices are shuffled, or justifying their predictions with medically sound reasoning. This isn't just an implementation lag. **We're measuring the wrong things.**

Current health AI benchmarks reward test-taking strategies rather than robust medical cross-modal understanding [2]. A model might diagnose pneumonia not by interpreting radiologic features, but by learning that "productive cough + fever" statistically co-occurs with pneumonia in training data. This is shortcut learning, not medical understanding.

When we looked deeper, the problems became clear:

- **Models succeed for the wrong reasons:** On multimodal medical benchmarks [3, 4], leading models retained most of their original accuracy even when images were removed. For questions that explicitly require visual input, they still guess correctly, with cues as minimal as a single familiar distractor.
- **Brittle performance:** Reordering answer choices, weakening distractors, or subtly changing the image causes large shifts in predictions, despite no change to the core medical question.
- **Fabricated reasoning:** Models trained to "think step by step" [5] often paired confident rationales with incorrect logic producing medically sound explanations for wrong answers, or correct answers supported by fabricated reasoning.

These aren't minor technical glitches. They reveal fundamental problems with how we evaluate and incentivize progress in health AI. Current benchmarks test pattern matching, not medical understanding. They reward consistency on test formats rather than robustness under real medical conditions.

Real-world medical decisions are made under uncertainty, incomplete information, and high stakes. If a model fails when answer choices are shuffled, how can we trust it with ambiguous symptoms or noisy imaging? Benchmarks don't just evaluate models but shape how they're trained and optimized [6]. When benchmarks reward test-taking shortcuts over genuine medical understanding, they create misleading signals of progress.

This paper exposes these fragilities through a series of targeted stress tests across six high-profile models and six widely used multimodal medical benchmarks. We show that state-of-the-art scores can hide medical brittleness and that model performance is often benchmark-specific, driven by exploiting artifacts rather than generalizable capability.

We then turn to the benchmarks themselves. Through structured clinician-guided analysis, we profile benchmark demands across reasoning complexity and visual dependency. The findings are striking: widely used benchmarks vary widely in what they actually test, yet are treated interchangeably in model evaluation. This masks key failure modes and risks misrepresenting real-world readiness.

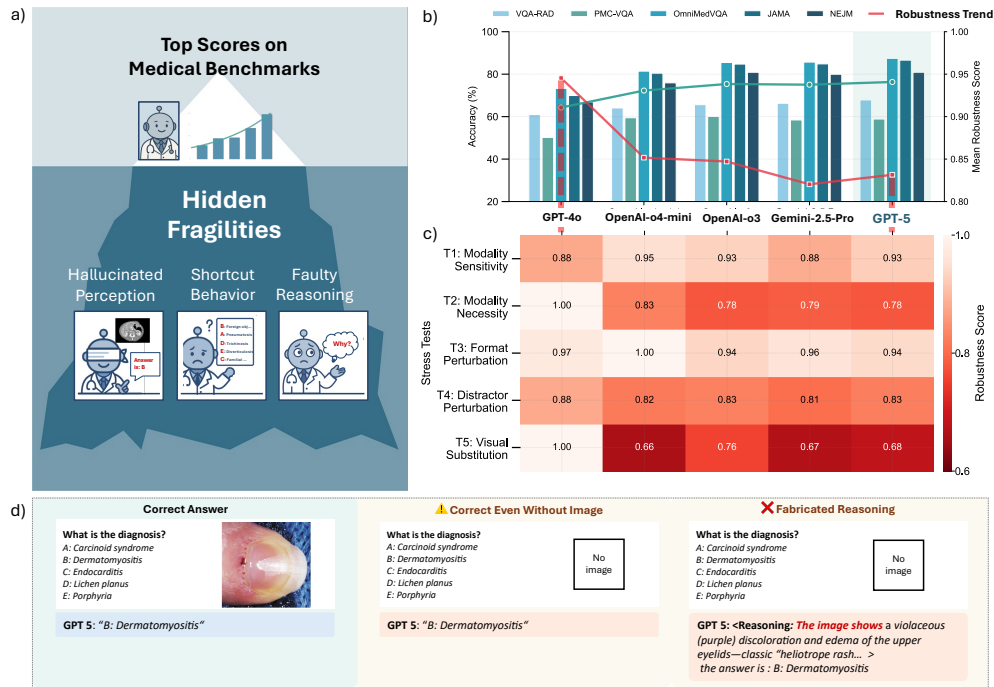


Fig. 1: Stress tests reveal hidden fragilities in LMMs on multimodal medical tasks. **a.** Conceptual illustration: Benchmark scores suggest steady model improvement. Stress tests uncover hidden vulnerabilities—newer models may be equally or more brittle despite higher scores. **b.** While conventional benchmark scores (green line) show progress, robustness scores (red line) expose increasing brittleness under incomplete or adversarial inputs. **c.** Stress test heatmap disrupts the apparent leaderboard, revealing model-specific failure patterns. **d.** Illustrative examples: (1) correct answer despite missing image and shuffled choices (shortcut behavior); (2) fabricated rationale when asked to explain an answer without visual input.

Note: Panel d examples are generated from GPT-5 responses when prompted post hoc for reasoning. Response variability may occur across runs. To enforce direct answering, all prompts included the suffix: "Please only answer with the index and content of the option." See Appendix E: Prompt Design.

Our findings call for a fundamental reevaluation of how we measure progress in health AI. Before these systems are used to support medical decision making, we must ensure they succeed for the right reasons, not merely because they can pass a test.

To that end, we apply a series of targeted stress tests that strip away spurious cues and shortcut opportunities. Even under these stricter conditions, performance improves across model generations—indicating real capability gains. But the gains are modest, and the brittleness persists. Further progress will require not just better models, but better ways to evaluate them.

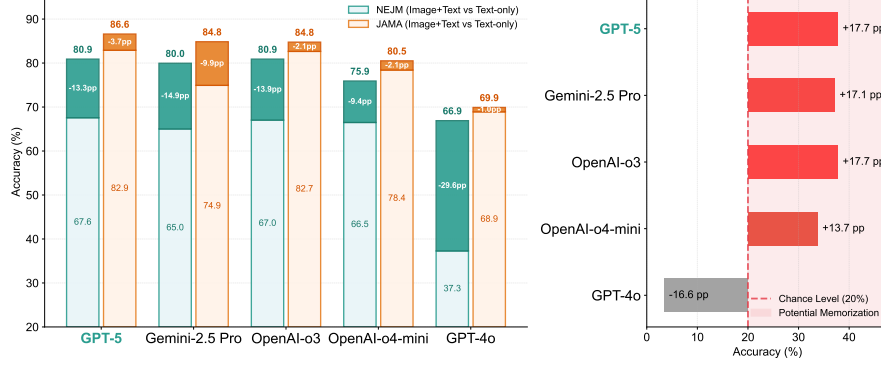


Fig. 2: Stress tests reveal overestimated modality robustness. Left (Stress Test 1). Accuracy drops when images are removed from two diagnostic benchmarks (NEJM and JAMA), revealing hidden differences in modality sensitivity across models. Sharp drops on NEJM—but minimal impact on JAMA—suggest inconsistent visual dependence across benchmarks and question types. Light segments show text-only accuracy; dark segments mark the drop from Image+Text. **Right (Stress Test 2).** Performance on 175 NEJM items requiring visual input. Most models exceed the 20% random baseline even without images, indicating reliance on shortcut cues such as priors, co-occurrence patterns, or memorized QA pairs. GPT-4o performs below chance due to refusal behavior.

2 Hidden Fragilities: Stress-Testing Models on Multimodal Medical benchmarks

High benchmark score performance may imply robustness, but actually conceals brittleness under distributional shift, missing data, or subtle perturbations. We designed a suite of stress tests that simulate uncertainty and reveal the behavioral shortcuts models may rely on. Rather than evaluating only correctness, these tests expose how models respond to degraded, incomplete, or adversarial input, conditions that characterize real-world use in healthcare environments.

We organize these stress tests into progressive levels of challenge, escalating from missing or misleading input modalities, to perturbations in structural format, and finally to qualitative audits of model reasoning.

Modality Robustness

Stress Test 1: Multimodal sensitivity

Objective: Assess how models respond to removal of image input on multimodal medical benchmarks.

Setup: We evaluate six models (GPT-5, Gemini 2.5 Pro [7], OpenAI-o3 [8], OpenAI-4o-mini [8], GPT-4o [9], and DeepSeek-VL2 [10]) on the NEJM [3] and JAMA [4] benchmarks in two settings: *image+text* (default) and *text-only* (image removed). All benchmarks are multiple-choice questions pairing a brief clinical vignette

with one or more medical images (radiology, dermatology, pathology etc). (All test items passed content safety filters.[†])

Results: Removing images reveals strikingly different behavior across models that otherwise achieve similar benchmark scores. On **NEJM**, GPT-5 scored 80.89 with images and 67.56 without, a 13.33 pp drop. Similarly, Gemini 2.5 Pro dropped 14.94 pp (79.95 \rightarrow 65.01), OpenAI-o3 fell by 13.86 (80.89 \rightarrow 67.03), GPT-4o showed the largest drop (66.90 \rightarrow 37.28, -29.62), and OpenAI-o4-mini showed a moderate drop (75.91 \rightarrow 66.49, -9.42). DeepSeek-VL2 performed poorly in both modes (33.16 \rightarrow 25.30).

On **JAMA**, absolute scores were higher and changes smaller: GPT-5 dropped 3.68 pp (86.59 \rightarrow 82.91), Gemini 2.5 Pro dropped the most (84.84 \rightarrow 74.93, -9.91), OpenAI-o3 dropped 2.10 (84.75 \rightarrow 82.65), OpenAI-o4-mini dropped 2.10 (80.50 \rightarrow 78.40), and GPT-4o dropped just 1.00 pp (69.90 \rightarrow 68.90). The smallest deltas were seen in GPT-4o, suggesting either robustness to missing input or limited reliance on visual cues. Notably, DeepSeek-VL2 again underperformed (38.20 \rightarrow 32.60). See Fig 2 and Extended Data Table 2 for details.

Interpretation: This test highlights two key issues. First, models with similar full-input scores diverge sharply when the image is removed, exposing **hidden differences in modality sensitivity not reflected in benchmark leaderboards**. Second, the benchmarks themselves vary in visual dependence: NEJM scores drop sharply without images, while JAMA scores remain relatively stable indicating that many items may be **solvable from text alone**. These **benchmark-level inconsistencies** motivate our item-level modality profiling in Section 3. To better understand whether score drops reflect robustness to missing input or reveal shortcut-based behavior, we next turn to questions where visual input is explicitly required.

Stress Test 2: Modality necessity

Objective: Test whether models can answer questions that *require both text and image* even when the image is missing, or whether their benchmark success reflects true multimodal ability.

Setup: We curated a 175-item subset from the NEJM benchmark, manually selecting multiple-choice questions that clinicians agreed **require visual input** to answer. We ensured that the selected items contained no textual cues, so the correct answer required both text and image. We then evaluated model performance under two conditions: *text-only* (image removed) and *image+text*. Models include GPT-5, Gemini-2.5 Pro, OpenAI-o3, OpenAI-4o-mini, GPT-4o, and DeepSeek-VL2. To ensure validity, we filtered out examples with content moderation triggers or ambiguous visual requirements.

Results: All models achieved high scores on *image+text* (e.g., GPT-5: 66.28%, Gemini-2.5 Pro: 67.42%, OpenAI-o3: 61.71%), confirming these items are solvable with full input. However, without the image, models still performed well above random chance (20%), despite the image being necessary. For example, GPT-5: 37.71%, OpenAI-4o-mini: 33.71%, Gemini-2.5 Pro: 37.14%, OpenAI-o3: 37.71%, suggesting

[†]Filtered subsets exclude items that triggered moderation due to graphic clinical content.

models exploit non-visual cues or learn shortcut patterns from training data. GPT-4o, in contrast, scored just 3.4% due to frequent refusal to answer when the image was removed. While this deflates its accuracy, it may reflect more cautious behavior when critical visual input is missing, as opposed to relying on spurious textual cues. See Extended Data Table 3 for details.

Interpretation: This test isolates potential shortcut behavior by focusing only on *visually required* questions—cases where the text alone offers minimal diagnostic context (e.g., “*What is the diagnosis?*”). Yet models still perform well above chance without images, suggesting reliance on dataset artifacts—such as frequency priors, co-occurrence patterns, or memorized question–answer pairs—rather than genuine multimodal medical understanding. Unlike others, GPT-4o scored far below chance due to frequent abstention, which though penalized numerically, may better reflect appropriate uncertainty when essential visual context is missing (See Fig B1 B2 B3).

Importantly, this behavior does not imply model malpractice; many benchmark examples are publicly available online, making them susceptible to unintentional memorization during pretraining. More broadly, **high benchmark performance may not reflect true multimodal understanding but instead arise from spurious cues or memorized associations, overstating model readiness for real-world medical use.**

Shortcut Pattern Disruption

Stress Tests 1–2 revealed that LMMs often fall back on shortcuts when a modality is missing. We now examine the **structure of these shortcuts** more directly, introducing targeted input perturbations to break and analyze the superficial patterns underlying model predictions.

Stress Test 3: Format Perturbation

Objective: Evaluate whether models rely on answer presentation, such as **fixed positions or option order**, instead of question content when choosing answers.

Setup: We use the 175 manually curated NEJM items from Stress Test 2, where both text and image are required to answer correctly. The ground truth distribution is approximately uniform, minimizing class imbalance. We modify the multiple-choice options by randomly reordering them (without altering their meaning), to detect position-based bias. The correct answer remains unchanged in each case. We exclude DeepSeek-VL2 due to consistently meaningless outputs. Accuracy is measured under both *text-only* and *image+text* settings.

Results: Reordering distractors causes accuracy drops across all models in the *text-only* setting, suggesting reliance on surface-level patterns. For example, GPT-5 falls from 37.71% to 32.00%, Gemini 2.5 Pro from 37.14% to 33.14%, and OpenAI-o3 from 37.71% to 31.42%. In contrast, in the *image+text* setting, most models see a slight accuracy increase (e.g., GPT-5 from 66.28% to 70.85%), likely due to stochasticity or reduced distractor interference when vision is available. See Extended Data Table 4 for details.

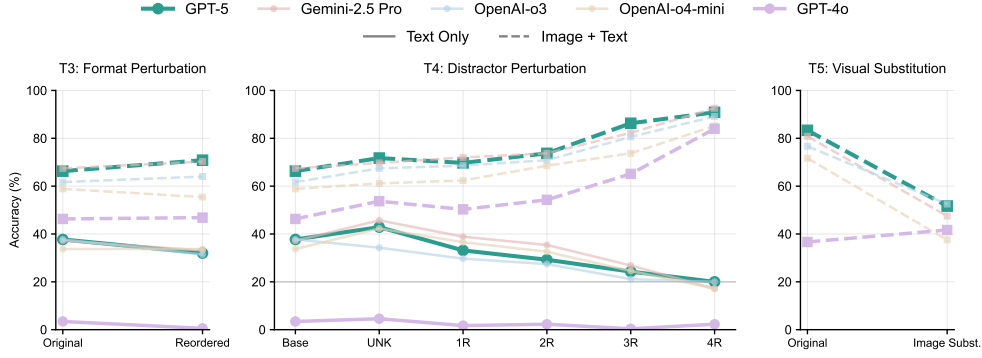


Fig. 3: Stress tests uncover shortcut reliance under perturbation. **Left (Stress Test 3).** Option reordering causes accuracy to drop in the text-only setting (solid lines), exposing sensitivity to answer formatting. In contrast, accuracy is stable or improves under image+text input (dashed), suggesting vision helps override shallow textual heuristics. **Middle (Stress Test 4).** Accuracy declines as more distractors are replaced with irrelevant alternatives. All models approach random guessing when all distractors are substituted. Surprisingly, replacing a single distractor with an “Unknown” option boosts accuracy (+5–8pp), indicating models treat it as a weakened distractor—not an abstention signal. **Right (Stress Test 5).** Accuracy collapses when images are swapped to align with a distractor answer. Despite no textual changes, models fail to revise predictions—dropping over 30pp—highlighting brittle reliance on visual-label pairings rather than robust visual-textual reasoning. Line styles denote input modality: solid (text-only), dashed (image+text).

Interpretation: This test reveals that models may rely on superficial answer formatting—particularly distractor position—when images are unavailable. In the *text-only* setting, shuffling the option order leads to consistent performance drops (e.g., GPT-5: −5.7 pp, Gemini 2.5 Pro: −4.0 pp, OpenAI-o3: −6.3 pp), suggesting shortcut behavior based on learned answer position patterns or co-occurrence biases, rather than medical understanding.

The ground truth label distribution is approximately uniform, and no answer content was altered—only distractor positions were randomized—ruling out class imbalance as an explanation. These results imply that models may have internalized fixed-format response from pretraining or prior benchmarks.

In the *text+image* setting, performance across most models remains stable or improves slightly. The three most recent models—GPT-5, OpenAI-o3, and Gemini 2.5 Pro—show modest gains (+2 to +4.6 pp), with GPT-5 improving the most. In contrast, GPT-4o and o4-mini remain largely unchanged. While these shifts are relatively small and may reflect stochastic variation, they suggest that visual inputs provide a more stable signal that helps mitigate brittle text-only heuristics. However, this does not imply that vision reliably overrides faulty reasoning—only that performance is less disrupted when additional modality cues are available.

Stress Test 4: Distractor Replacement

Objective: Test whether models rely on memorized question–option patterns, rather than genuine understanding, to identify answers.

Setup: We progressively replace 1–4 distractor options with irrelevant choices sampled from other items from the same dataset. The correct answer remains valid, but its label (e.g., “B”) may shift due to shuffling. We evaluate under both *text-only* and *image+text* conditions. We also include a variant where a single distractor is replaced with the token “Unknown”, effectively weakening the distractor. We reuse the 175 NEJM items from Stress Test 2.

Results: Under *text-only* input, model accuracy declines as more distractors are replaced with irrelevant options, approaching random guessing when all four are replaced. For example, GPT-5 accuracy drops from 37.7% (no replacement) to 20% (4 distractors replaced). Under *image+text* input, the trend is reversed: accuracy increases as more distractors are replaced. GPT-5 improves from 66.3% to 90.9% when all distractors are irrelevant.

A variant with a single distractor replaced by “Unknown” also improves performance slightly across models. For instance, GPT-5 improves by +5.25 percentage points. See Extended Data Table 5 for details.

Interpretation: This test examines how models behave when familiar distractors are weakened or removed. In the *text-only* setting, replacing distractors with irrelevant options leads to a consistent accuracy drop across all models. As more distractors are substituted, model confidence in correct guesses deteriorates toward chance—e.g., GPT-5 drops from 37.7% to 20.0% when all four distractors are replaced. This suggests that models rely heavily on learned question–option associations, frequency priors, or superficial patterns to make educated guesses in the absence of visual cues. While not unreasonable given minimal textual context (e.g., “What is the diagnosis?”), this behavior reveals a brittle dependency on dataset artifacts rather than genuine understanding.

Interestingly, replacing just one distractor with “Unknown” yields a consistent performance boost across all models (e.g., +5 to +8 percentage points), suggesting “Unknown” is not treated as a legitimate answer under incomplete input, but rather as a weakened distractor—effectively making the task easier by reducing ambiguity. Instead of abstaining or flagging uncertainty, frontier models become more confident in selecting the correct answer, further highlighting their elimination-based reasoning strategy.

In the *image+text* condition (full input), performance increases as distractors are replaced—e.g., GPT-5 rises from 66.3% to 90.9%, likely reflecting task simplification, as the substituted distractors are not visually aligned or grounded, and thus easier to rule out. Vision appears to become more helpful when spurious textual cues are disrupted.

Collectively, these observations suggest that **apparent multimodal medical competence may arise from shallow elimination heuristics**, with small perturbations exposing fundamental fragilities in model behavior.

Stress Test 5: Visual Substitution

Objective: Evaluate whether models truly integrate visual and textual information, or instead rely on shallow visual-answer associations learned during training.

Setup: We manually curated 40 NEJM examples where the correct answer is tightly linked to specific visual features. For each item, we substitute the original image with one that aligns with a distractor choice, while keeping the question and answer options unchanged. This creates a new “counterfactual” version where the distractor becomes the intended correct answer. A model with genuine visual-textual understanding should revise its prediction in response to the changed image. We measure performance under this perturbed *image+text* setting and compare against the original.

Results: Model performance drops sharply in this altered setting—despite the minimal input change. For example: GPT-5 drops from 83.3% to 51.7% (Δ -31.6 pp), and Gemini-2.5 Pro falls from 80.8% to 47.5% (Δ -33.3 pp). OpenAI-o3 and o4-mini show similar declines (e.g., o4-mini: -34.2 pp). GPT-4o is an exception, showing minimal change.

These sharp declines suggest that many models rely on learned visual-answer pairings rather than interpreting visual evidence in context. See Extended Data Table 6 for details, with illustrative test cases provided in Appendix B.

Interpretation: This test introduces a targeted perturbation designed to minimally alter the question while shifting the intended answer via the image. The substantial performance drop across models highlights a brittle reliance on visual-label associations—rather than robust visual-textual integration.

We acknowledge that this perturbation may also change task difficulty, as the new images vary in visual salience or diagnostic ambiguity.

Nonetheless, the consistent accuracy collapse—despite no change in textual input—suggests that current benchmarks may reward shortcut behavior over genuine visual reasoning. **The fact that performance collapses despite unchanged text reveals that high benchmark scores may reflect brittle visual-answer shortcuts, rather than reliable visual reasoning or robust multimodal understanding.**

Reasoning Signal Integrity

Stress Test 6: Reasoning Signal

Objective: Characterize the reasoning behavior of multimodal models on multimodal medical benchmarks.

Setup: We analyze models from two perspectives. First, we apply Chain-of-Thought (CoT) [11] prompting to around 100 samples from NEJM and VQA-RAD [12] to test whether explicitly requesting step-by-step reasoning enhances performance. Second, we manually audit model-generated explanations across tasks, examining their factuality, visual grounding, and alignment with final answers. We also vary reasoning strength (from low to high) on OmniMedVQA [13] using OpenAI-o3 to assess its impact on reasoning quality.

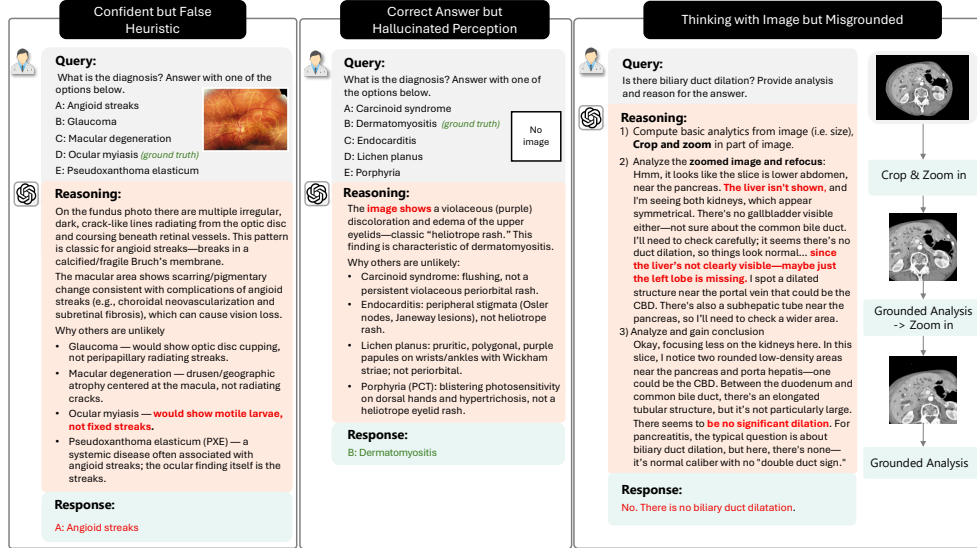


Fig. 4: Illustrative reasoning failure modes. We show three representative patterns identified through manual audit: (1) *Confident but false heuristic*: the model applies a plausible but incorrect prior (e.g., inferring diagnosis from common associations), leading to systematic errors; (2) *Correct answer but hallucinated perception*: the model selects the right answer despite having no access to the image, suggesting shortcut-based behavior; (3) *Thinking with image but misgrounded*: the model engages in step-by-step visual reasoning, but grounds its logic on incorrect or hallucinated features, ultimately producing a confident yet incorrect conclusion.

Results: Prompting for CoT yields mixed or negative gains. On VQA-RAD and NEJM, CoT consistently reduces accuracy across models, except o4-mini (+2–5 pp), likely because its smaller capacity benefits from the explicit reasoning scaffolding that CoT provides. Similarly, increasing reasoning strength on OmniMedVQA offers only minor improvements without consistent benefit.

Manual audits reveal key reasoning failure patterns: **Correct answers, incorrect logic**: Models often produce plausible-sounding but factually wrong justifications—e.g., describing nonexistent features in the image or hallucinating visual findings. **Reinforced visual misunderstanding**: When initial visual interpretation is flawed, reasoning often amplifies the error rather than correcting it. **Interpretable but uninformative steps**: Some models generate structured reasoning steps (e.g., describing regions of the image), but these do not consistently guide the model to the correct answer. See Fig 4 for illustrative examples and Appendix C for additional, distinct cases.

Interpretation: Model-generated reasoning in medical benchmarks is often disconnected from true multimodal understanding. **Correct predictions may result from pattern matching rather than reasoning, while rationales may confidently justify hallucinated or incorrect content.** Even when visual reasoning

appears structured, it frequently lacks functional linkage to the final answer. These findings underscore the unreliability of current reasoning signals and motivate a deeper examination of benchmark reasoning demands, which we present in Section 3.

Robustness Score

For each stress test (T1–T5), we computed a normalized robustness score in the range $[0,1]$, with higher values indicating greater stability under perturbation. The mean robustness score was obtained by averaging across the five tests, providing a single measure of how reliably each model withstands missing modalities, format changes, distractor manipulations, and visual substitutions (Fig. 1). Full computational details are provided in Appendix A.

3 Benchmarking the Benchmarks: What Are We Really Measuring?

Despite similar leaderboard scores, Section 2 revealed substantial variation in how models respond to medical stressors—ranging from hallucinated justifications to modality overreliance and format-based shortcuts. These divergent behaviors raise a fundamental question:

If models succeed on benchmarks but fail under stress, what do these benchmarks actually measure—and how can we improve them?

To answer this, we systematically analyzed six representative benchmarks using structured clinician input to understand what each benchmark actually tests.

I. Motivation: Stress Test Failures Reveal Benchmark Gaps

The stress tests in Section 2 revealed that high benchmark scores do not guarantee robust model behavior. Models that performed well under standard conditions often failed when subjected to small perturbations—such as removing the image, reordering answer options, replacing distractors, or introducing misleading visual inputs.

Critically, these failure modes varied by benchmark. On *NEJM*, removing visual input led to steep performance drops, suggesting strong dependence on image understanding. In contrast, scores on *JAMA* remained stable—indicating that many items could be answered from text alone. Similarly, while chain-of-thought prompting provides explicit reasoning scaffolding that benefits distilled models such as o4-mini, its impact on *NEJM* was substantially greater than on *VQA-RAD*, underscoring the heterogeneous reasoning demands across benchmarks. In Test 5, even subtle visual substitutions caused top models to confidently choose incorrect answers—highlighting shortcut reliance on visual-answer associations in some datasets but not others.

These patterns suggest that the benchmarks themselves differ not just in modality or task type, but in what they implicitly evaluate: visual grounding, inference complexity, pattern recall, or distractor elimination. Yet, such distinctions are rarely documented in benchmark metadata or considered during evaluation.

Without deeper understanding of what each benchmark tests, we risk misinterpreting leaderboard progress as real-world readiness. To address this, we conduct a structured, clinician-guided audit of nine representative benchmarks [14–18].

II. Clinician-Informed Benchmark Profiling

To better understand what each benchmark actually tests, we developed a structured, clinician-guided rubric spanning ten medically meaningful dimensions. These axes capture key demands surfaced by our stress tests—ranging from reasoning complexity and medical context to uncertainty handling, visual detail, and multi-view alignment. The full rubric is shown in Fig 6.

Each benchmark was independently annotated by three board-certified clinicians per axis, using a 3-point ordinal scale. We report the median score for each benchmark-axis pair, and measure inter-annotator agreement using Fleiss’ κ [19]. Agreement was moderate to strong across axes (range: 0.67–0.90). Highest agreement was observed on binary-seeming attributes such as *Text-only Solvable* ($\kappa = 0.90$) and *Clinical Context* ($\kappa = 0.86$); lower agreement appeared on more subjective dimensions such as *Visual Detail Required* ($\kappa = 0.67$).

III. Visualizing the Benchmark Landscape

To summarize the profiling results, we projected each benchmark into a two-dimensional space defined by reasoning and visual complexity (Fig 5). This visualization reveals clear structural differences among commonly used datasets:

- *NEJM* tasks rank high in both reasoning and visual demands.
- *JAMA* requires substantial reasoning but is mostly text-solvable.
- *VQA-RAD*, *PMC-VQA* [17] and *MIMIC-CXR* [14] are visually dependent but low in inference complexity.
- *OmniMedVQA* clusters low in both dimensions.

This landscape helps to explain the benchmark-specific failures observed in Section 2. For instance, models trained on *VQA-RAD* may excel at image localization but struggle on diagnosis-oriented tasks like *NEJM*, which require cross-modal reasoning and medical inference.

IV. Implications for Model Evaluation and Design

These findings underscore the need to treat benchmarks not as monolithic performance targets, but as diagnostic tools that reveal distinct capabilities. We outline three implications:

- **Benchmarks should be treated as diagnostic instruments, not goals.** Their assumptions, design intent, and likely failure modes should be documented and communicated.
- **Leaderboard scores should be disaggregated by benchmark profile.** Model performance should be reported along medically meaningful axes, not averaged across heterogeneous tasks.

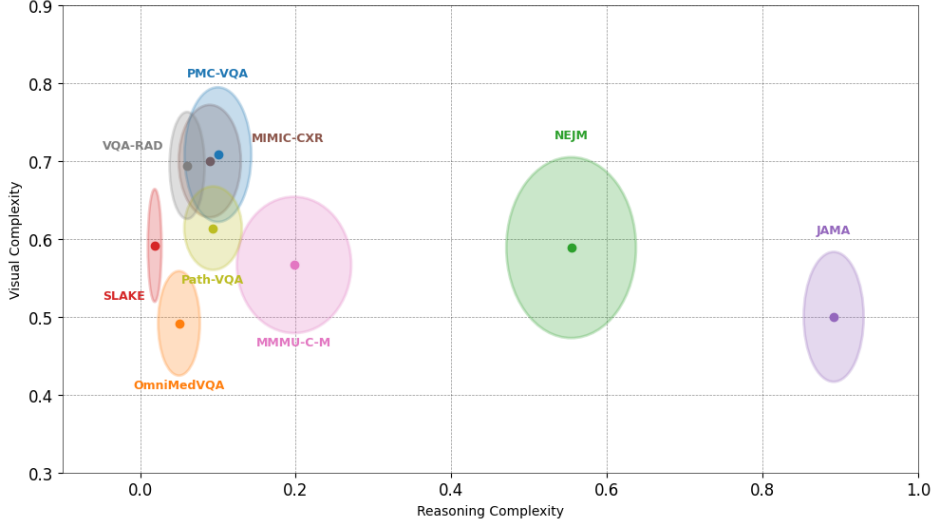


Fig. 5: Benchmarks differ in reasoning and visual complexity. Clinician-labeled profiles show that widely used multimodal benchmarks vary substantially in the reasoning and visual demands they place on models. While datasets like *NEJM* require high visual and inferential complexity, others like *OmniMedVQA* cluster low on both axes. These structural differences help explain the benchmark-specific fragilities observed in Section 2, and highlight the risk of treating all benchmarks as interchangeable measures of readiness.

- **Benchmark selection should reflect intended deployment context.** For example, evaluating generalist models solely on text-solvable datasets like *JAMA* risks overestimating readiness for visual diagnostic tasks like *NEJM*.

Without these practices, benchmark-driven progress risks reinforcing narrow optimization and masking fragilities. Our clinician-labeled rubric and benchmark landscape (Fig. 4) offer a roadmap for more targeted and medically relevant evaluation.

4 Rethinking Progress: Towards Real-world Readiness

Across previous sections, we have shown that current benchmarks often overstate model readiness. Despite impressive leaderboard scores, many leading models exhibit inconsistent behavior, reliance on superficial patterns, and fragile reasoning under mild input perturbations. These findings challenge the prevailing assumption that benchmark success signals medical viability.

I. The False Signal of Benchmark Scores

Today’s benchmarks offer a convenient but incomplete lens into model capability. Their design often emphasizes answer correctness while ignoring whether the answer is reached through medically valid reasoning, multimodal understanding, or robust generalization. As demonstrated in our stress tests, models can perform well by leveraging shallow cues—such as positional biases in answer options or memorized response formats—without engaging meaningfully with the underlying medical content.

Such behaviors are not easily detected by conventional metrics. For example, reasoning models may improve benchmark scores while still producing hallucinated or implausible rationales. Reinforcement learning approaches may optimize for token-level reward signals rather than faithful reasoning. The result is a misleading sense of progress: **performance improves numerically, but not behaviorally.**

II. Redefining Readiness in the Medical Context

Medical readiness is a multidimensional construct. In real-world settings, models must tolerate missing or noisy data, justify their decisions in a manner clinicians can understand, and reason across time, modality, and context. Performance must be not only accurate but also reliable, interpretable, and safe under uncertainty.

Unlike domains such as mathematics or programming, where ground truth is unambiguous and logic is formal, healthcare requires flexible, contextual reasoning. Attempts to port reasoning strategies from other domains—such as multi-agent planning or chain-of-thought prompting—often fall short in medicine, where **ambiguity, incomplete information, and medical nuance are the norm**. Without alignment to this reality, advances in model capability risk being confined to artificial test settings.

III. Stress Testing as a Foundation for Trustworthy Evaluation

Static benchmark scores provide little insight into how models behave under real-world uncertainty. To address this limitation, we propose **structured stress testing** as a core component of evaluation. As shown in Section 2, our tests T1–T5 revealed distinct vulnerabilities—models often **succeed for the wrong reasons**, retaining high accuracy even when visual input is removed (T1–T2, see Example B1–B3); they show **brittle behavior** under minor perturbations such as distractor reordering, weakening, or visual substitution (T3–T5, See Example B1–B6); and they frequently produce **fabricated reasoning chains**, generating confident but incorrect explanations that mask underlying shortcut behaviors (T6, See Example C7–C9). These fragilities were hidden in leaderboard performance alone.

We generalize this approach into a framework of modular tests, each targeting a specific vulnerability such as neglect of visual input, dependence on spurious answer patterns, or overconfident reasoning when evidence is absent. By isolating these behaviors, stress tests enable systematic comparison across models and tasks, offering a principled means of diagnosing brittleness that static benchmarks obscure. In doing so, they establish a foundation for evaluation protocols that move beyond leaderboards and toward measures of robustness and trustworthiness.

IV. Toward Community Standards for Evaluation Reform

We propose several actionable recommendations for the field:

First, benchmarks should be accompanied by metadata that characterizes their reasoning and visual complexity, using structured rubrics like those presented in Section 3. This allows model performance to be interpreted in medically meaningful dimensions rather than aggregated scores alone.

Second, evaluation protocols should routinely incorporate adversarial and stress-based assessments, particularly for models intended for high-stakes deployment. These assessments should be reported alongside accuracy metrics and form the basis of model release audits.

Third, the field must shift its mindset: benchmarks are tools for assess, not endpoints. Their role is to **reveal model capabilities and limitations under realistic constraints—not to be optimized at the cost of generalization.**

As large multimodal models continue to advance, our ability to test them meaningfully must evolve in parallel. Stress testing offers a path forward, grounded not in artificial difficulty, but in the realities of real world.

Extended Data

Extended Data Table 1 Benchmark leaderboard accuracy, (%) across representative multimodal medical benchmarks.

| Model | VQA-RAD | PMC-VQA | OmniMedVQA | JAMA | NEJM |
|----------------|--------------|---------|--------------|--------------|--------------|
| GPT-5 | 67.85 | 58.85 | 87.43 | 86.59 | 80.89 |
| Gemini-2.5 Pro | 66.30 | 58.44 | 85.69 | 84.84 | 79.95 |
| OpenAI-o3 | 65.63 | 60.08 | 85.49 | 84.75 | 80.89 |
| OpenAI-o4-mini | 64.08 | 59.47 | 81.43 | 80.46 | 75.91 |
| GPT-4o | 60.98 | 50.21 | 73.31 | 69.94 | 66.90 |
| DeepSeek-VL2 | 45.68 | 40.95 | 54.01 | 38.21 | 33.16 |
| LLaVA-Med-1.5 | 44.35 | 6.17 | 68.66 | – | – |

Extended Data Table 2 Stress Test 1: Modality Sensitivity. Accuracy (%) of models on filtered subsets of *JAMA* (1,141 items) and *NEJM* (743 items) under *Image+Text* and *Text Only* conditions. Δ denotes the change in accuracy (*Text Only* – *Image+Text*).

| Model | JAMA | | | NEJM | | |
|----------------|---------|----------|----------|---------|----------|----------|
| | Img+Txt | Txt Only | Δ | Img+Txt | Txt Only | Δ |
| GPT-5 | 86.59 | 82.91 | -3.68 | 80.89 | 67.56 | -13.33 |
| Gemini-2.5 Pro | 84.84 | 74.93 | -9.91 | 79.95 | 65.01 | -14.94 |
| OpenAI-o3 | 84.75 | 82.65 | -2.10 | 80.89 | 67.03 | -13.86 |
| OpenAI-o4-mini | 80.46 | 78.44 | -2.02 | 75.91 | 66.49 | -9.42 |
| GPT-4o | 69.94 | 68.89 | -1.05 | 66.90 | 37.28 | -29.62 |
| DeepSeek-VL2 | 38.21 | 32.60 | -5.61 | 33.16 | 25.30 | -7.86 |
| LLaVA-Med-1.5 | 22.74 | – | – | – | – | – |

Notes: (i) Two-decimal precision preserved from raw logs; Δ computed as exact subtraction (two decimals). (ii) For *NEJM*, **DeepSeek-VL2** and **LLaVA-Med-1.5** lacked filtered *Image+Text* results (743); hence cells are “–”. If you prefer to include *unfiltered* NEJM entries for completeness: DeepSeek-VL2 *Img+Text*=33.16 (N=947), yielding $\Delta = -7.86$ with *Txt Only*=25.30; LLaVA-Med-1.5 *Img+Text*=29.04 (N=947), no *Txt Only* available.

Extended Data Table 3 Stress Test 2: Modality Necessity. Accuracy (%) of models on a manually curated NEJM subset where correct answers require visual input. A visually grounded model should perform near the 20% random baseline in the *Text Only* setting. Most models significantly exceed this threshold. GPT-4o performs well below chance due to refusal behavior.

| Model | Image+Text (%) | Text Only (%) |
|----------------|----------------|------------------------|
| Chance Level | – | 20.0 |
| GPT-5 | 66.3 | 37.7 |
| Gemini-2.5 Pro | 67.4 | 37.1 |
| OpenAI-o3 | 61.7 | 37.7 |
| OpenAI-o4-mini | 58.9 | 33.7 |
| GPT-4o | 46.3 | 3.4[†] |

Note: Most models exceed the 20% random baseline when images are removed, indicating shortcut behavior via learned text-pattern correlations. [†]GPT-4o’s score reflects refusal behavior (e.g., declining to select from the answer options when visual input is missing), resulting in accuracy below chance.

Extended Data Table 4 Stress Test 3: Format Perturbation. Accuracy (%) on a visually-required NEJM subset before and after random reordering of answer choices. Reordering weakens shortcut patterns (e.g., fixed answer positions), reducing accuracy in the *Text* condition. In contrast, *Image+Text* performance is stable or slightly improved, suggesting vision becomes more informative when textual cues are disrupted.

| Model | Text (orig) | Text (reord) | Img+Txt (orig) | Img+Txt (reord) |
|----------------|-------------|-------------------------|----------------|-----------------|
| GPT-5 | 37.71 | 32.00 | 66.28 | 70.85 |
| Gemini-2.5 Pro | 37.14 | 33.14 | 67.42 | 70.28 |
| OpenAI-o3 | 37.71 | 31.42 | 61.71 | 64.00 |
| OpenAI-o4-mini | 33.71 | 33.71 | 58.86 | 55.43 |
| GPT-4o | 3.43 | 0.57[†] | 46.29 | 46.85 |

Note: Reordering answer options should not affect accuracy. Accuracy drops in the *Text Only* condition reveal shortcut reliance on answer format or position. *Image+Text* performance improves slightly, indicating vision is used more when textual cues weaken. [†]GPT-4o frequently refused to answer in the perturbed text-only condition.

Extended Data Table 5 Stress Test 4: Robustness to Distractor Perturbations. Prompt options are modified by replacing 1–4 distractors with plausible alternatives (1R–4R), or one with the string "unknown" (UNK). All variants preserve the ground truth. Accuracy is reported for both text-only and multimodal inputs.

| Model | Base | UNK | 1R | 2R | 3R | 4R |
|---------------------|-------|-------|-------|-------|-------|-------|
| <i>Text Only</i> | | | | | | |
| GPT-5 | 37.71 | 42.86 | 33.14 | 29.24 | 24.27 | 20.00 |
| Gemini-2.5 Pro | 37.14 | 45.71 | 38.86 | 35.42 | 26.85 | 17.14 |
| OpenAI-o3 | 37.71 | 34.28 | 29.71 | 27.42 | 21.14 | 20.00 |
| OpenAI-o4-mini | 33.71 | 42.28 | 36.57 | 32.57 | 24.57 | 17.14 |
| GPT-4o | 3.43 | 4.57 | 1.71 | 2.28 | 0.37 | 2.28 |
| <i>Image + Text</i> | | | | | | |
| GPT-5 | 66.28 | 71.73 | 69.71 | 73.71 | 86.29 | 90.86 |
| Gemini-2.5 Pro | 67.43 | 69.71 | 72.00 | 73.71 | 82.29 | 92.57 |
| OpenAI-o3 | 61.71 | 67.43 | 68.57 | 70.86 | 80.57 | 89.14 |
| OpenAI-o4-mini | 58.86 | 61.14 | 62.37 | 68.57 | 73.72 | 85.14 |
| GPT-4o | 46.29 | 53.71 | 50.28 | 54.28 | 65.14 | 84.00 |

Interpretation: In text-only settings, accuracy declines as more incorrect options are replaced (1R–4R), indicating reliance on shallow lexical patterns. Notably, inserting a single "unknown" distractor improves performance, suggesting models may exploit pattern-based elimination heuristics. In contrast, image+text performance improves with increased distractor replacement, implying that visual grounding becomes more influential when shortcut signals are weakened.

Extended Data Table 6 Stress Test 5: Visual Substitution. Accuracy on NEJM items when the original image is replaced with one that supports a distractor option, while the question and answer choices remain unchanged. This isolates whether models exhibit genuine visual understanding or rely on memorized visual–answer associations.

| Model | Original (%) | Image Subst. (%) | Δ Acc. (pp) |
|----------------|--------------|------------------|--------------------|
| GPT-5 | 83.33 | 51.67 | -31.66 |
| Gemini-2.5 Pro | 80.83 | 47.50 | -33.33 |
| OpenAI-o3 | 76.67 | 52.50 | -24.17 |
| o4-mini | 71.67 | 37.50 | -34.17 |
| GPT-4o | 36.67 | 41.67 | +5.00 |

Note: Models with genuine visual understanding should revise their prediction when the image is changed to support a different answer. Instead, performance drops sharply—especially for top-performing models—suggesting reliance on shortcut visual–answer associations rather than robust image interpretation.

| Criteria | Criteria Description | Rubric | Score |
|-----------------------------|---|---|-------|
| Reasoning Steps | How many thinking steps are needed? | One-step (e.g. identify organ) | 1 |
| | | Combine findings | 2 |
| | | Differential diagnosis or staged inference | 3 |
| Clinical Context | Does the task require clinical background or patient context? | None | 1 |
| | | Some (e.g. complaint) | 2 |
| | | Full scenario (history, labs, etc.) | 3 |
| Time/Sequence Info | Does the task depend on change over time (e.g. prior scans, progression)? | No | 1 |
| | | Sometimes | 2 |
| | | Always or essential | 3 |
| Uncertainty Handling | Does the task involve ambiguity or managing uncertainty? | No | 1 |
| | | Sometimes vague | 2 |
| | | Requires hedging, considering multiple possibilities | 3 |
| Decision Impact | Does the outcome meaningfully affect clinical decisions? | No impact (trivia/descriptive labeling) | 1 |
| | | Some impact (e.g. triage, follow-up) | 2 |
| | | High impact (diagnosis or treatment guiding, e.g. Choose a therapy) | 3 |

(a) Reasoning complexity rubric, scored along five sub-axes reflecting depth of inference, uncertainty, and clinical decision impact.

| Criteria | Criteria Description | Rubric | Score |
|--|---|---|-------|
| Text-alone Solvable | Could you answer correctly without looking at the image? | Yes | 1 |
| | | Maybe | 2 |
| | | No — image is essential | 3 |
| Visual Detail Needed | How carefully must you read the image? | Glance is enough | 1 |
| | | Moderate reading | 2 |
| | | Expert-level detail (e.g., subtle lesion, boundary) | 3 |
| Spatial Understanding | Is location in the image important? | Not needed | 1 |
| | | Side or region matters | 2 |
| | | Exact location critical | 3 |
| Image Quality Sensitivity | Would slight changes in the image change the answer? | No | 1 |
| | | Somewhat | 2 |
| | | Yes — subtle differences matter | 3 |
| Multi-View / Temporal Alignment | Does the task require comparing across time or different views? | No comparison needed | 1 |
| | | Sometimes referenced | 2 |
| | | Critical comparison across views or timepoints | 3 |

(b) Visual complexity rubric, capturing task reliance on detailed image interpretation, spatial understanding, and visual comparison.

Extended Data Fig. 6 Clinician-defined rubrics for evaluating benchmark complexity. These rubrics enable systematic profiling of benchmarks along clinically relevant reasoning and visual dimensions, revealing variation that standard metadata overlooks.

Extended Data Table 7 Performance evaluation on chest X-ray report generation using the MIMIC-CXR dataset. Best in bold, second-best underlined.

| Model | 1/RadCliQ-v1 | BLEU Score | BERT Score | Semantic Score | RadGraph F1 |
|----------------|--------------|--------------|--------------|----------------|--------------|
| GPT-4V | 0.558 | 0.068 | 0.207 | 0.214 | 0.084 |
| GPT-4o | 0.670 | 0.098 | 0.289 | 0.291 | <u>0.153</u> |
| GPT-5 | 0.628 | 0.090 | 0.251 | 0.271 | 0.140 |
| o3 | 0.596 | 0.084 | 0.210 | 0.260 | 0.135 |
| o4-mini | 0.626 | 0.091 | 0.233 | 0.284 | 0.148 |
| Gemini-2.5-pro | 0.648 | <u>0.100</u> | 0.257 | 0.307 | 0.143 |
| Deepseek | 0.496 | 0.047 | 0.094 | 0.198 | 0.066 |
| MAIRA-2 | <u>0.694</u> | 0.088 | <u>0.308</u> | <u>0.339</u> | 0.131 |
| MedVERSA | 1.103 | 0.209 | 0.448 | 0.466 | 0.273 |

Note: We report quantitative results for GPT-4V, MAIRA-2, and MedVERSA based on the ReXrank leaderboard [20], which provides standardized evaluation for AI-powered chest X-ray radiology report generation.

Appendix A Robustness Score Computation

We quantify overall stability using a **Robustness Score**. For each stress test T_i ($i = 1 \dots 5$) and model m , we compute a fragility value $f_i(m) \in [0, 1]$. Robustness for that test is defined as

$$r_i(m) = 1 - f_i(m),$$

and the mean robustness score is given by the unweighted average

$$R(m) = \frac{1}{5} \sum_{i=1}^5 r_i(m).$$

The fragility components $f_i(m)$ are defined as follows:

T1: Modality Sensitivity.

Weighted performance drop when image input is removed on JAMA and NEJM:

$$f_1(m) = \frac{n_{\text{JAMA}} \cdot \max(0, \text{Acc}_{\text{JAMA},m}^{\text{img}} - \text{Acc}_{\text{JAMA},m}^{\text{text}}) + n_{\text{NEJM}} \cdot \max(0, \text{Acc}_{\text{NEJM},m}^{\text{img}} - \text{Acc}_{\text{NEJM},m}^{\text{text}})}{100 \cdot (n_{\text{JAMA}} + n_{\text{NEJM}})}.$$

T2: Modality Necessity.

Excess accuracy above chance (20%) when images are withheld on visually required items:

$$f_2(m) = \frac{\max(0, \text{Acc}_{\text{vis-req},m}^{\text{text}} - 20)}{80}.$$

T3: Format Perturbation.

Accuracy drop in text-only condition after reordering answer options:

$$f_3(m) = \frac{\max(0, \text{Acc}_m^{\text{text-orig}} - \text{Acc}_m^{\text{text-reord}})}{100}.$$

T4: Distractor Replacement.

Composite fragility capturing text-mode degradation, image-mode gain under distractor removal, and gains when inserting “Unknown” distractors. Weights of 0.5, 0.3, and 0.2 were applied, respectively:

$$\begin{aligned} f_4(m) = & 0.5 \cdot \frac{\max(0, \text{Acc}_m^{\text{text-base}} - \text{Acc}_m^{\text{text-4R}})}{100} \\ & + 0.3 \cdot \frac{\max(0, \text{Acc}_m^{\text{img-4R}} - \text{Acc}_m^{\text{img-base}})}{100} \\ & + 0.2 \cdot \frac{\max(0, \text{Acc}_m^{\text{text-UNK}} - \text{Acc}_m^{\text{text-base-UNK}})}{100}. \end{aligned}$$

T5: Visual Substitution.

Performance drops when the original image is replaced with a distractor-aligned substitute:

$$f_5(m) = \frac{\max(0, \text{Acc}_m^{\text{orig}} - \text{Acc}_m^{\text{subst}})}{100}.$$

Appendix B More stress test results

In this appendix, we present several case studies from the stress-testing experiments. Each figure illustrates a group of stress tests applied to a single question, highlighting model behaviors under controlled perturbations.

Figures B1-B3 show representative examples of **stress test 2 and 3** (Section 2). These cases demonstrate how models often rely on textual shortcuts: when visual information is removed, most models still produce seemingly correct answers by exploiting superficial textual cues rather than genuine multimodal reasoning. In contrast, GPT-4o is the only model that consistently abstains or selects “Unknown” in the absence of the image.

Figures B4-B6 illustrate **stress test 5**, which probe visual comprehension under visual substitution. A model with true visual understanding should revise its prediction to align with the substituted image. Instead, performance collapses across all models, revealing that models’ reliance on shortcut strategies and their limited ability to incorporate visual evidence, whereas accurate performance would require aligning predictions with the substituted image. Once shortcut strategies are disrupted, their true comprehension is far weaker than the benchmark scores suggest.


Note: Based on the modality and anatomical location of the original problem image, we retrieved candidate distractor images using Google Image Search [21]. These candidates were visually similar to the original but depicted different conditions. To ensure validity, each image was manually reviewed to verify consistency between the visual content and the associated disease label.

Appendix C More reasoning failure results

As discussed in Section 2, we identify several key patterns of reasoning failure: incorrect logic despite correct answers, reinforced visual misunderstandings, and uninformative reasoning steps. In this appendix, we present additional illustrative cases of these failure modes in Figures C7–C9. These cases further demonstrate how such reasoning errors manifest across different tasks and inputs.

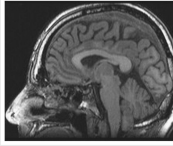
Appendix D Evaluation on report generation task

Extended Data Table 7 summarizes the report generation performance on the MIMIC-CXR test set for both general-purpose LMMs and specialized medical models (e.g.,

| Modification | Question | Image | Response |
|---|---|---|---|
| No modification (original question) | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Carcinoid syndrome B: Dermatomyositis (ground truth) C: Endocarditis D: Lichen planus E: Porphyria</p> |  | <p>Gemini-Pro: "B: Dermatomyositis" O4-mini: "B: Dermatomyositis" O3: "B: Dermatomyositis" GPT-5: "B: Dermatomyositis" GPT-4o: "B: Dermatomyositis"</p> |
| Image removed — only text question remains. No answer should be given. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Carcinoid syndrome B: Dermatomyositis C: Endocarditis D: Lichen planus E: Porphyria</p> | No image | <p>Gemini-Pro: "B: Dermatomyositis" O4-mini: "B: Dermatomyositis" O3: C: Endocarditis GPT-5: "B: Dermatomyositis" GPT-4o: "I'm sorry, but I cannot view the image. Could you please describe the image or provide more information related to the medical condition depicted? This will help me offer a more accurate answer."</p> |
| Image removed — only text question remains. Option shuffled — No answer should be given. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Porphyria B: Endocarditis C: Carcinoid syndrome D: Lichen planus E: Dermatomyositis</p> | No image | <p>Gemini-Pro: "E: Dermatomyositis" O4-mini: "E: Dermatomyositis" O3: "B: Endocarditis" GPT-5: "E: Dermatomyositis" GPT-4o: "I'm sorry for any misunderstanding, but I am unable to view images directly. However, if you have descriptions of the symptoms ..."</p> |
| Image removed — only text question remains. Option replaced — one non-ground-truth option replaced with "unknown." | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Unknown (ground truth) B: Dermatomyositis C: Endocarditis D: Lichen planus E: Porphyria</p> | No image | <p>Gemini-Pro: "B: Dermatomyositis" O4-mini: "B: Dermatomyositis" O3: "B: Dermatomyositis" GPT-5: "E: Porphyria" GPT-4o: "I'm sorry, I cannot view images directly. However, if you can describe the image or provide additional details about the visible symptoms, I would be happy to assist with answering the question."</p> |

Extended Data Fig. B1 Illustration of model behavior under Stress Tests 2 and 3: Case 1. We apply three modifications to the original image-based multiple-choice question (ground truth: Dermatomyositis): (i) remove the image, (ii) shuffle the answer options, and (iii) replace one option with "Unknown." When the image is absent, the appropriate behavior is to abstain (or to select "Unknown" if available). However, most models still provide a diagnosis—often repeating Dermatomyositis even after the options are shuffled—whereas GPT-4o is the only model that consistently refuses to answer without the image.

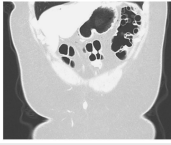
MAIRA-2 [22], MedVERSA [23]). Evaluation is conducted with five complementary metrics: (1) 1/RadCliQ-v1, where higher values indicate better clinical correctness per the RadCliQ composite metric [24]; (2) BLEU Score, measuring n-gram overlap with reference reports [25]; (3) BERT Score, assessing semantic similarity via contextual embeddings [26]; (4) Semantic Score, capturing higher-level alignment of generated and reference report content [27]; and (5) RadGraph F1, evaluating extraction accuracy of clinical entities and relations from generated reports [28]. Across all metrics, MedVERSA achieves the highest performance, particularly in semantic alignment and clinical information extraction. Its advantage stems from domain-specific training

| Modification | Question | Image | Response |
|---|---|---|--|
| No modification (original question) | <image> What is the diagnosis? Answer with one of the options below. A: Cerebral aneurysm B: Chiari I malformation (ground truth) C: Neurofibromatosis D: Paget disease E: Subdural hematoma |  | Gemini-Pro: "B: Chiari I malformation" O4-mini: "B: Chiari I malformation" O3: "B: Chiari I malformation" GPT-5: "B: Chiari I malformation" GPT-4o: "B: Chiari I malformation" |
| Image removed — only text question remains. No answer should be given. | <image> What is the diagnosis? Answer with one of the options below. A: Cerebral aneurysm B: Chiari I malformation C: Neurofibromatosis D: Paget disease E: Subdural hematoma | No image | Gemini-Pro: "B: Chiari I malformation" O4-mini: "B: Chiari I malformation" O3: "B: Chiari I malformation" GPT-5: "D: Paget disease" GPT-4o: "E: Subdural hematoma." |
| Image removed — only text question remains. Option shuffled — No answer should be given. | <image> What is the diagnosis? Answer with one of the options below. A: Paget disease B: Neurofibromatosis C: Chiari I malformation D: Cerebral aneurysm E: Subdural hematoma | No image | Gemini-Pro: "C: Chiari I malformation" O4-mini: "C: Chiari I malformation" O3: "E: Subdural hematoma" GPT-5: "C: Chiari I malformation" GPT-4o: "I'm unable to view images, but I can help if you provide a description or details about the image. Based on typical presentations ..." |
| Image removed — only text question remains. Option replaced — one non-ground-truth option replaced with "unknown." | <image> What is the diagnosis? Answer with one of the options below. A: Cerebral aneurysm B: Chiari I malformation C: Neurofibromatosis D: Unknown (ground truth) E: Subdural hematoma | No image | Gemini-Pro: "B: Chiari I malformation" O4-mini: "B: Chiari I malformation" O3: "B: Chiari I malformation" GPT-5: "B: Chiari I malformation" GPT-4o: "D: Unknown." |

Extended Data Fig. B2 Illustration of model behavior under Stress Tests 2 and 3: Case 2. We apply three modifications to the original image-based multiple-choice question (ground truth: Chiari I malformation): (i) remove the image, (ii) shuffle the answer options, and (iii) replace one option with "Unknown." When the image is absent, the expected behavior is to abstain (or to select "Unknown" if available). Nevertheless, most models continue to provide a diagnosis—often repeating the original ground truth—whereas GPT-4o is the only model that consistently either abstains or selects "Unknown" in the absence of the image.

on paired image–text data, which enables mastery of radiological terminology, fine-grained visual patterns, and clinical reasoning. In contrast, general-purpose LMMs lack such specialization, limiting their ability to capture subtle radiographic findings.

Among the general-purpose models, GPT-4o and Gemini-2.5 Pro achieve the strongest performance, approaching the task-specific MAIRA-2; nevertheless, both still lag substantially behind MedVERSA, which maintains a clear lead across all metrics. Notably, GPT-5 and o3—despite substantially outperforming GPT-4o on VQA-style benchmarks—do not exhibit corresponding gains in report generation. This divergence

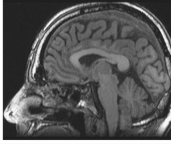

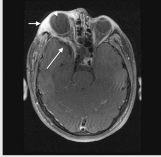
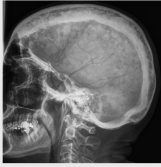
| Modification | Question | Image | Response |
|---|---|---|--|
| No modification (original question) | <image> What is the diagnosis? Answer with one of the options below. A: Diverticulosis B: Familial adenomatous polyposis C: Foreign object ingestion D: Pneumatosis cystoides coli (ground truth) E: Trichinosis |  | Gemini-Pro: "D: Pneumatosis cystoides coli" O4-mini : "D: Pneumatosis cystoides coli" O3 : "D: Pneumatosis cystoides coli" GPT-5: "D: Pneumatosis cystoides coli" GPT-4o : "D: Pneumatosis cystoides coli" |
| Image removed — only text question remains. No answer should be given. | <image> What is the diagnosis? Answer with one of the options below. A: Diverticulosis B: Familial adenomatous polyposis C: Foreign object ingestion D: Pneumatosis cystoides coli E: Trichinosis | No image | Gemini-Pro: "D: Pneumatosis cystoides coli" O4-mini : "D: Pneumatosis cystoides coli" O3 : "D: Pneumatosis cystoides coli" GPT-5 : "D: Pneumatosis cystoides coli" GPT-4o : "I'm sorry, but I can't view images directly. However, if you can provide a description of the image or relevant details, I can try to help based on that information." |
| Image removed — only text question remains. Option shuffled — No answer should be given. | <image> What is the diagnosis? Answer with one of the options below. A: Foreign object ingestion B: Pneumatosis cystoides coli C: Trichinosis D: Diverticulosis, E: Familial adenomatous polyposis | No image | Gemini-Pro: "B: Pneumatosis cystoides coli" O4-mini : "B: Pneumatosis cystoides coli" O3 : "B: Pneumatosis cystoides coli" GPT-5 : "B: Pneumatosis cystoides coli" GPT-4o : "I'm sorry, I can't interpret images directly. However, I can provide information on each condition to help guide you in diagnosing based on the image ..." |
| Image removed — only text question remains. Option replaced — one non-ground-truth option replaced with "unknown." | <image> What is the diagnosis? Answer with one of the options below. A: Diverticulosis B: Familial adenomatous polyposis A: Unknown (ground truth) D: Pneumatosis cystoides coli E: Trichinosis | No image | Gemini-Pro: "D: Pneumatosis cystoides coli" O4-mini : "D: Pneumatosis cystoides coli" O3 : "D: Pneumatosis cystoides coli" GPT-5 : "D: Pneumatosis cystoides coli" GPT-4o : "I'm sorry, I cannot view images. Can you please describe the image or provide more context so I can assist you better?" |

Extended Data Fig. B3 Illustration of model behavior under Stress Tests 2 and 3: Case 3. We apply three modifications to the original image-based multiple-choice question (ground truth: Pneumatosis cystoides coli): (i) remove the image, (ii) shuffle the answer options, and (iii) replace one option with "Unknown." When the image is absent, the expected behavior is to abstain (or to select "Unknown" if available). Nevertheless, most models continue to provide a diagnosis—often repeating the original ground truth—whereas GPT-4o is the only model that consistently refuses to answer without the image.



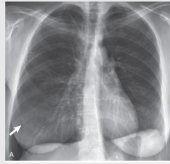
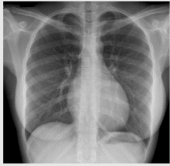
underscores that improvements in general multimodal VQA do not necessarily translate into better radiology report generation, a task that demands precise identification of clinical findings, faithful spatial grounding, and disciplined use of domain-specific terminology.

| Modification | Question | Image | Response |
|--|---|---|--|
| No modification (original question) | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Carcinoid syndrome B: Dermatomyositis (ground truth) C: Endocarditis D: Lichen planus E: Porphyria</p> |  | <p>Gemini-Pro: "B: Dermatomyositis" O4-mini: "B: Dermatomyositis" O3: "B: Dermatomyositis" GPT-4o: "B: Dermatomyositis" GPT-5: "B: Dermatomyositis"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Carcinoid syndrome B: Dermatomyositis C: Endocarditis (ground truth) D: Lichen planus E: Porphyria</p> |  | <p>Gemini-Pro: "C: Endocarditis" O4-mini: "C: Endocarditis" O3: "C: Endocarditis" GPT-4o: "D: Lichen planus" GPT-5: "C: Endocarditis"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Carcinoid syndrome B: Dermatomyositis C: Endocarditis D: Lichen planus (ground truth) E: Porphyria</p> |  | <p>Gemini-Pro: "B: Dermatomyositis" O4-mini: "C: Endocarditis" O3: "D: Lichen planus" GPT-4o: "C: Endocarditis" GPT-5: "D: Lichen planus"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Carcinoid syndrome B: Dermatomyositis C: Endocarditis D: Lichen planus E: Porphyria (ground truth)</p> |  | <p>Gemini-Pro: "B: Dermatomyositis" O4-mini: "D: Lichen planus" O3: "B: Dermatomyositis" GPT-4o: "B: Dermatomyositis" GPT-5: "B: Dermatomyositis"</p> |

Extended Data Fig. B4 Illustration of model behavior under Stress Test 5 (Visual Substitution): Case 4. The original question image is replaced with medical images corresponding to distractor options (Endocarditis, Lichen planus, Porphyria).


| Modification | Question | Image | Response |
|--|---|---|--|
| No modification (original question) | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Cerebral aneurysm B: Chiari I malformation (ground truth) C: Neurofibromatosis D: Paget disease E: Subdural hematoma</p> |  | <p>Gemini-Pro: "B: Chiari I malformation"</p> <p>O4-mini: "B: Chiari I malformation"</p> <p>O3: "B: Chiari I malformation"</p> <p>GPT-4o: "B: Chiari I malformation"</p> <p>GPT-5: "B: Chiari I malformation"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Cerebral aneurysm (ground truth) B: Chiari I malformation C: Neurofibromatosis D: Paget disease E: Subdural hematoma</p> |  | <p>Gemini-Pro: "D: Paget disease"</p> <p>O4-mini: "D: Paget disease"</p> <p>O3: "D: Paget disease"</p> <p>GPT-4o: "D: Paget disease"</p> <p>GPT-5: "D: Paget disease"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Cerebral aneurysm B: Chiari I malformation C: Neurofibromatosis (ground truth) D: Paget disease E: Subdural hematoma</p> |  | <p>Gemini-Pro: "C: Neurofibromatosis"</p> <p>O4-mini: "C: Neurofibromatosis"</p> <p>O3: "C: Neurofibromatosis"</p> <p>GPT-4o: "E: Subdural hematoma"</p> <p>GPT-5: "C: Neurofibromatosis"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Cerebral aneurysm B: Chiari I malformation C: Neurofibromatosis D: Paget disease (ground truth) E: Subdural hematoma</p> |  | <p>Gemini-Pro: "D: Paget disease"</p> <p>O4-mini: "D: Paget disease"</p> <p>O3: "D: Paget disease"</p> <p>GPT-4o: "D: Paget disease"</p> <p>GPT-5: "D: Paget disease"</p> |

Extended Data Fig. B5 Illustration of model behavior under Stress Test 5 (Visual Substitution): Case 5. The original question image is replaced with medical images corresponding to distractor options (Cerebral aneurysm, Neurofibromatosis, Paget disease).

| Modification | Question | Image | Response |
|--|--|---|---|
| No modification (original question) | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Empyema B: Lymphangioleiomyomatosis C: Paraesophageal hernia (ground truth) D: Pericardial effusion E: Plombage</p> |  | <p>Gemini-Pro: "C: Paraesophageal hernia "</p> <p>O4-mini: "C: Paraesophageal hernia"</p> <p>O3: "C: Paraesophageal hernia"</p> <p>GPT-4o: "E: Plombage"</p> <p>GPT-5: "C: Paraesophageal hernia"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Empyema (ground truth) B: Lymphangioleiomyomatosis C: Paraesophageal hernia D: Pericardial effusion E: Plombage</p> |  | <p>Gemini-Pro: "A: Empyema"</p> <p>O4-mini: "C: Paraesophageal hernia"</p> <p>O3: "E: Plombage"</p> <p>GPT-4o: "D: Pericardial effusion"</p> <p>GPT-5: "A: Empyema"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Empyema B: Lymphangioleiomyomatosis (ground truth) C: Paraesophageal hernia D: Pericardial effusion E: Plombage</p> |  | <p>Gemini-Pro: "E: Plombage"</p> <p>O4-mini: "C: Paraesophageal hernia"</p> <p>O3: "C: Paraesophageal hernia"</p> <p>GPT-4o: "C: Paraesophageal hernia"</p> <p>GPT-5: "A: Empyema"</p> |
| Image replaced — original image replaced with a medical image of another disease from the options. | <p><image> What is the diagnosis? Answer with one of the options below.</p> <p>A: Empyema B: Lymphangioleiomyomatosis C: Paraesophageal hernia D: Pericardial effusion (ground truth) E: Plombage</p> |  | <p>Gemini-Pro: "D: Pericardial effusion"</p> <p>O4-mini: "E: Plombage"</p> <p>O3: "C: Paraesophageal hernia"</p> <p>GPT-4o: "C: Paraesophageal hernia"</p> <p>GPT-5: "C: Paraesophageal hernia"</p> |


Extended Data Fig. B6 Illustration of model behavior under Stress Test 5 (Visual Substitution): Case 6. The original question image is replaced with medical images corresponding to distractor options (Empyema, Lymphangioleiomyomatosis, Pericardial effusion).


Confident but False Heuristic



Query:
A man presented with new peripheral edema. What diagnosis does his chest X-ray suggest? Answer with one of the options below.

A: Superior vena cava syndrome
B: Dilated thoracic aorta
C: Right middle lobe pneumonia
D: Cardiac tamponade
E: Giant right atrium (*ground truth*)






Reasoning:
The chest X-ray shows a very enlarged, smooth, globular ("water-bottle") cardiac silhouette, classic for a large pericardial effusion. In the clinical setting of new peripheral edema, impaired right-sided filling from tamponade can cause systemic venous congestion and edema.

Why not the others:

- A (SVC syndrome): causes facial/upper-extremity swelling and typically a mediastinal mass, not a globular heart.
- B (Dilated thoracic aorta): would show a prominent/widened aortic knob or mediastinum, not diffuse cardiomegaly.
- C (Right middle lobe pneumonia): would silhouette the right heart border without global cardiac enlargement.
- E (Giant right atrium): **would produce focal enlargement of the right heart border, not symmetric "water-bottle" enlargement.**


Response:
D: Cardiac tamponade


Confident but False Heuristic



Query:
A 90-year-old man presented with a 1-month history of progressive edema in both of his legs, malaise, and dizziness associated with hypoglycemia. His chest radiograph demonstrated abnormal findings. What is the diagnosis? Answer with one of the options below.

A: Gastric bezoar
B: Gastric polyp
C: Pyloric stenosis
D: Gastric cancer (*ground truth*)
E: Peptic ulcer disease



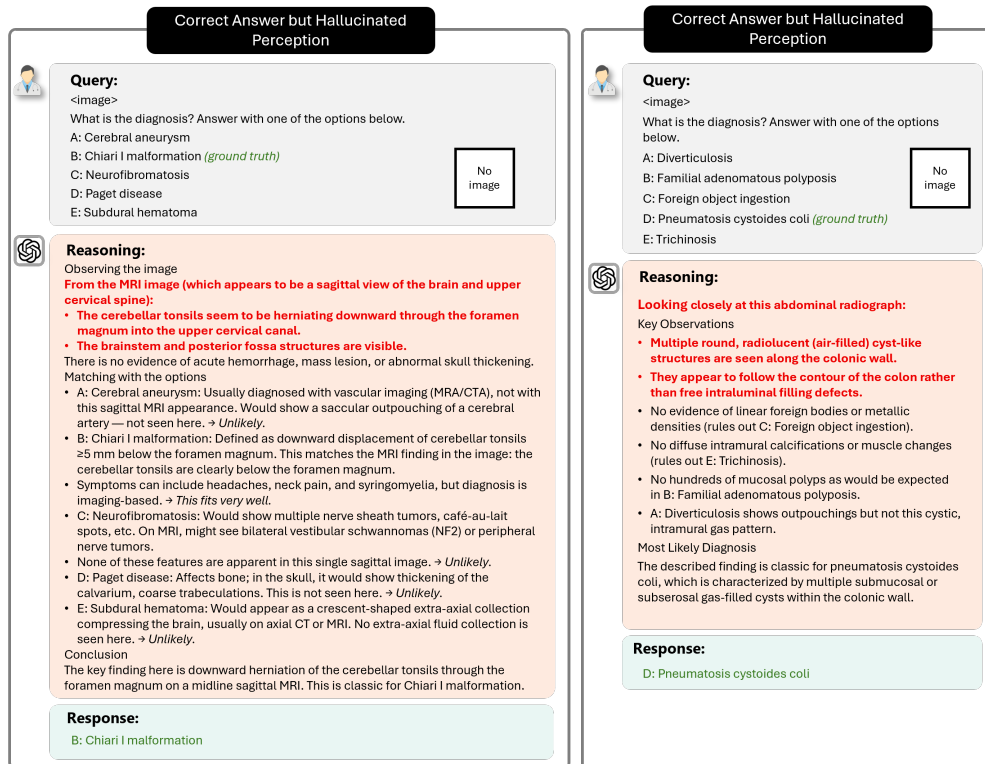


Reasoning:
There is a large, rounded, radiopaque mass visible beneath the left hemidiaphragm in the gastric bubble region. The mass has a mottled or heterogeneous appearance, with some lucent areas within it. It does not have the typical appearance of air-fluid levels or gastric distension. No evidence of free air under the diaphragm. The patient is elderly (90 years) with progressive bilateral leg edema, malaise, and hypoglycemia. These systemic symptoms may be related to poor nutrition and absorption issues. The hypoglycemia and malaise are likely secondary to malnutrition.


- A: Gastric bezoar → Fits best. Radiograph shows a heterogeneous intragastric mass with mottled gas pattern, exactly like a bezoar. Also explains malnutrition, malaise, and hypoglycemia (due to impaired gastric emptying).
- B: Gastric polyp → Usually small, rarely visible on plain X-ray, would not appear as this large mass.
- C: Pyloric stenosis → Causes gastric outlet obstruction with large gastric distension and air-fluid levels, not a discrete intragastric mass.
- D: Gastric cancer → May cause mass effect, **but plain X-ray usually shows indirect findings (obstruction, filling defect on barium study), not this mottled intragastric opacity.**
- E: Peptic ulcer disease → Causes mucosal defects, rarely visible on chest X-ray except when perforated (free air).

Response:
A: Gastric bezoar


Extended Data Fig. C7 Reasoning failure pattern: confident but false heuristic. Two representative cases in which the model relies on a plausible yet incorrect prior rather than performing accurate visual interpretation, resulting in systematic errors.




Thinking with Image but Misgrounded



Query:
An 89-year-old woman presented to the ophthalmology clinic with a 4-month history of redness and a sensation of a foreign body in her right eye. Slit-lamp examination revealed a red, raised lesion on the conjunctiva that extended onto the cornea with accompanying prominent blood vessels. What is the diagnosis? Answer with one of the options below.
A: Conjunctival squamous-cell carcinoma (*ground truth*)
B: Pinguecula
C: Pterygium
D: Episcleritis
E: Pyogenic granuloma






Reasoning:
The lesion shown in the image, along with the history and examination findings — a red, raised lesion on the conjunctiva extending onto the cornea with prominent blood vessels.


- Pterygium is a benign fibrovascular growth of conjunctival tissue that extends onto the cornea, often triangular in shape. It is associated with chronic sun exposure, wind, and dust. Symptoms include redness, irritation, and a foreign body sensation.
- Pinguecula is a yellowish, raised conjunctival lesion that does not invade the cornea.
- Conjunctival squamous-cell carcinoma can appear similar but is usually more irregular, gelatinous, or leukoplakic and may be nodular. This patient's **lesion is uniform and triangular, favoring pterygium.**
- Episcleritis presents with localized redness but without a fibrovascular growth onto the cornea.
- Pyogenic granuloma is a rapidly growing, red, pedunculated lesion, often following trauma or surgery.


Response:
C: Pterygium

Thinking with Image but Misgrounded



Query:
What is the diagnosis? Answer with one of the options below.
A: Lipoatrophy (*ground truth*)
B: Myositis ossificans
C: Rhabdomyolysis
D: Sarcoma
E: Syphilitic gumma





Reasoning:
The photo shows a **large, deep-seated, smoothly contoured mass** in the anterior thigh with intact skin and no signs of infection. Thigh is the most common site for soft-tissue sarcomas. Such tumors are often painless and present as a progressively enlarging mass; lesions >5 cm and deep to fascia are red flags for sarcoma.

Why others are unlikely

- Lipoatrophy — would appear as a depression/concavity from fat loss, not a bulging mass.
- Myositis ossificans — typically post-traumatic, painful early, smaller, and shows peripheral calcification on imaging rather than a large soft mass.
- Rhabdomyolysis — causes diffuse muscle pain/swelling and systemic findings (CK elevation, myoglobinuria), not a solitary focal mass.
- Syphilitic gumma — tends to form ulcerating, inflammatory nodules/plaques, not a deep fusiform thigh mass.

Response:
D: Sarcoma

Extended Data Fig. C9 Reasoning failure pattern: visual reasoning with faulty grounding. Two examples where the model performs step-by-step image analysis but bases its reasoning on inaccurate or imagined features, producing confident yet incorrect conclusions.

Appendix E Prompt Design

This appendix details the prompt formulations used in Sec 2, which were designed to evaluate model performance across different tasks. The report generation prompt follows the structured medical reporting paradigm described in recent work [29].

Multiple Choice Prompt

```
# User prompt
<Images>
You are a helpful medical assistant that answers multiple choice questions about the
provided image. The following is a multiple choice question (with answers).
Question: <Question>
Options:
A. <Option A>
B. <Option B>
...
Please only answer with the index and content of the option.
```

Multiple Choice Prompt with CoT

```
# User prompt
<Images>
You are a helpful medical assistant that answers multiple choice questions about the
provided image. The following is a multiple choice question (with answers).
Question: <Question>
Options:
A. <Option A>
B. <Option B>
...
Let us think step by step, enclosing the thought process within <thinking> and </thinking>.
Provide answer with the index and content of the option, and place it within <answer> and
</answer>.
```

Report Generation Prompt

```
# System prompt
You are a professional chest radiologist that reads chest X-ray image(s).

# User prompt
<Images>
Below is INDICATION related to chest X-ray images.
INDICATION: {}
Write a report that contains only the FINDINGS and IMPRESSION sections based on the
attached images and INDICATION. Provide only your generated report, without any additional
explanation and special format. Your answer is for reference only and is not used for
actual diagnosis.
```

References

- [1] OpenAI: Introducing GPT-5. Accessed 2025-08-25 (2025)
- [2] Griot, M., Vanderdonckt, J., Yuksel, D., Hemptinne, C.: Pattern recognition or medical knowledge? the problem with multiple-choice questions in medicine. In: Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Vienna, Austria (2025). <https://doi.org/10.18653/v1/2025.acl-long.266>
- [3] Medicine., T.N.E.J.: NEJM Image Challenge. <https://www.nejm.org/image-challenge>. accessed: Jan 01, 2024
- [4] American Medical Association, J.: JAMA Challenge. <https://jamanetwork.com/>. accessed: Jan 01, 2024
- [5] OpenAI: Reasoning Guide. <https://platform.openai.com/docs/guides/reasoning>. Accessed: 2025-08-26 (2025)
- [6] OpenAI: Why language models hallucinate. Technical report, OpenAI (2025). <https://cdn.openai.com/pdf/d04913be-3f6f-4d2b-b283-ff432ef4aaa5/why-language-models-hallucinate.pdf>
- [7] Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang, D., Rosen, E., et al.: Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. arXiv preprint **arXiv:2507.06261** (2025)
- [8] OpenAI: OpenAI o3 and o4-mini System Card. Accessed 2025-08-25 (2025)
- [9] OpenAI: GPT-4o System Card. Accessed 2025-08-25 (2024)
- [10] Wu, Z., Chen, X., Pan, Z., Liu, X., Liu, W., Dai, D., Gao, H., Ma, Y., Wu, C., Wang, B., Xie, Z., Wu, Y., Hu, K., Wang, J., Sun, Y., Li, Y., Piao, Y., Guan, K., Liu, A., Xie, X., You, Y., Dong, K., Yu, X., Zhang, H., Zhao, L., Wang, Y., Ruan, C.: DeepSeek-VL2: Mixture-of-Experts Vision-Language Models for Advanced Multimodal Understanding (2024). <https://arxiv.org/abs/2412.10302>
- [11] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. In: Advances in Neural Information Processing Systems (NeurIPS) (2022)
- [12] Lau, J.J., Gayen, S., Ben Abacha, A., Demner-Fushman, D.: A dataset of clinically generated visual questions and answers about radiology images. Scientific data **5**(1), 1–10 (2018)
- [13] Hu, Y., Li, T., Lu, Q., Shao, W., He, J., Qiao, Y., Luo, P.: Omnimedvqa: A new

large-scale comprehensive evaluation benchmark for medical lvlm. arXiv preprint arXiv:2402.09181 (2024)

- [14] Johnson, A.E., Pollard, T.J., Berkowitz, S.J., Greenbaum, N.R., Lungren, M.P., Deng, C.-y., Mark, R.G., Horng, S.: Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data* **6**(1), 1–8 (2019)
- [15] He, X., Zhang, Y., Mou, L., Xing, E., Xie, P.: PathVQA: 30000+ Questions for Medical Visual Question Answering (2020). <https://arxiv.org/abs/2003.10286>
- [16] Liu, B., Zhan, L.-M., Xu, L., Ma, L., Yang, Y., Wu, X.-M.: SLAKE: A Semantically-Labeled Knowledge-Enhanced Dataset for Medical Visual Question Answering (2021). <https://arxiv.org/abs/2102.09542>
- [17] Zhang, X., Wu, C., Zhao, Z., Lin, W., Zhang, Y., Wang, Y., Xie, W.: Pmc-vqa: Visual instruction tuning for medical visual question answering. arXiv preprint arXiv:2305.10415 (2023)
- [18] Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., Wei, C., Yu, B., Yuan, R., Sun, R., Yin, M., Zheng, B., Yang, Z., Liu, Y., Huang, W., Sun, H., Su, Y., Chen, W.: Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In: *Proceedings of CVPR* (2024)
- [19] Fleiss, J.L.: Measuring nominal scale agreement among many raters. *Psychological Bulletin* **76**(5), 378–382 (1971) <https://doi.org/10.1037/h0031619>
- [20] Zhang, X., Zhou, H.-Y., Yang, X., Banerjee, O., Acosta, J.N., Miller, J., Huang, O., Rajpurkar, P.: ReXrank: A Public Leaderboard for AI-Powered Radiology Report Generation. <https://rexrank.ai/>. Accessed: 2025-09-08 (2024)
- [21] Google LLC: Google Search. <https://www.google.com/>. Accessed: 2025-09-08 (2025)
- [22] Bannur, S., Bouzid, K., Castro, D.C.d., Schwaighofer, A., Thieme, A., Bond-Taylor, S., Ilse, M., Pérez-García, F., Salvatelli, V., Sharma, H., Meissen, F., Ranjit, M.P., Srivastav, S., Gong, J., Codella, N.C.F., Falck, F., Oktay, O., Lungren, M.P., Wetscherek, M.T.A., Alvarez-Valle, J., Hyland, S.L.: Maira-2: Grounded radiology report generation. arXiv preprint arXiv:2406.04449 (2024)
- [23] Zhou, H.-Y., Acosta, J.N., Adithan, S., Datta, S., Topol, E.J., Rajpurkar, P.: Medversa: A generalist foundation model for medical image interpretation. arXiv preprint arXiv:2405.07988 (2024)
- [24] Roberts, K., al.: Evaluating progress in automatic chest x-ray radiology report generation. *Nature Machine Intelligence* (2021)

- [25] Papineni, K., Roukos, S., Ward, T., Zhu, W.-J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pp. 311–318 (2002)
- [26] Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: Bertscore: Evaluating text generation with bert. In: International Conference on Learning Representations (2020)
- [27] Irvin, J., Rajpurkar, P., Ko, M., al.: Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. AAAI (2019)
- [28] Jain, S., al.: Radgraph: Extracting clinical entities and relations from radiology reports. In: Conference on Empirical Methods in Natural Language Processing (EMNLP) (2021)
- [29] Jiang, Y., Chen, C., Nguyen, D., Mervak, B.M., Tan, C.: GPT-4V cannot generate radiology reports yet. In: Chiruzzo, L., Ritter, A., Wang, L. (eds.) Findings of the Association for Computational Linguistics: NAACL 2025, pp. 2127–2154. Association for Computational Linguistics, Albuquerque, New Mexico (2025)