

---

# Contextual Contamination and Cognitive Inertia in Multimodal AI

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1      **Background:** While state-of-the-art artificial intelligence (AI) models achieve  
2      human expert-level performance in specific medical imaging tasks, they exhibit  
3      fundamental limitations when encountering out-of-distribution (OOD) inputs, com-  
4      mitting seemingly basic errors that defy common-sense reasoning. This study  
5      extends beyond the known issue of AI's anatomical common-sense deficits by  
6      exploring a "hierarchical error" mechanism wherein one dominant context distorts  
7      the interpretation of adjacent, logically unrelated information.

8      **Methods:** We designed a multi-stage qualitative experimental framework. First,  
9      we established baseline performance by presenting several cutting-edge multimodal  
10     AI systems with standard radiological images and normal anatomical illustrations.  
11     Next, we observed common-sense failures by testing the models on anatomically  
12     impossible ("nonsensical") images. Finally, we assessed AI responses to an abstract  
13     image of six converging lines both in isolation and when paired with an abnormal  
14     six-fingered hand emoji.

15     **Results:** The AI systems performed with high accuracy on standard medical  
16     images and normal illustrations, but they completely failed to recognize the inherent  
17     impossibility of the nonsensical images. Most significantly, when the abstract line  
18     image was presented alone, the AI correctly identified six lines; however, when the  
19     identical image was shown alongside the abnormal hand emoji, the AI incorrectly  
20     reported five lines. This demonstrates how a powerful semantic context (in this  
21     case, "hand = five") can hierarchically dominate and contaminate the processing of  
22     visual information in an otherwise unrelated image.

23     **Conclusions:** AI's most fundamental errors extend beyond simple pattern-  
24     recognition failures. They stem from structural limitations characterized by "cog-  
25     nitive inertia" (entrenchment in dominant prior assumptions) and "contextual  
26     contamination" (distortion of surrounding information by those assumptions). This  
27     represents a critical limitation that cannot be resolved with prompt engineering  
28     or other user-level fixes alone. Human critical thinking and oversight therefore  
29     remain essential to complement AI's autonomous diagnostic capabilities.

30     

## 1 Introduction

31     Predictions that AI will replace physicians are no longer science fiction [1, 2]. Sam Altman, the  
32     CEO of OpenAI, has asserted that "ChatGPT demonstrates superior diagnostic capabilities compared  
33     to most doctors," and various media outlets have heralded the "rise of robotic radiologists." [3]  
34     Indeed, AI systems have demonstrated their potential by achieving accuracy rates that surpass  
35     human performance in specific medical imaging tasks. However, underlying these achievements are  
36     fundamental limitations of current AI models. [4]

37 This study explores these limitations by examining how AI responds to anatomically impossible  
38 images that even a layperson can immediately recognize as implausible. We also investigate a  
39 puzzling phenomenon in which an AI system that accurately counts six lines will misidentify them  
40 as five when a hand emoji is placed adjacent to the lines. We propose that this behavior stems from  
41 two fundamental limitations of AI: “cognitive inertia” and “contextual contamination.” Cognitive  
42 inertia refers to the model becoming entrenched in a dominant interpretive framework, whereas  
43 contextual contamination denotes that framework’s distortion of the interpretation of otherwise  
44 unrelated information. Through this investigation, we aim to demonstrate that the future of medical  
45 AI relies not only on technological advancement but also on understanding and supervising AI’s  
46 structural limitations.

## 47 **2 Methods**

48 Systems. We evaluated four contemporary multimodal LLMs: ChatGPT-5 Thinking (OpenAI),  
49 Google Gemini 2.5 Pro, Grok-4 (xAI), and Claude-4-Opus-Thinking (Anthropic). Unless otherwise  
50 noted, systems were queried via their standard chat interfaces with default settings.

51 Stimuli. Standard medical images (X-ray/CT/MRI) were obtained from open educational sources.  
52 Polydactyly images were sourced from professional society materials. Additional medical illus-  
53 trations and hand emojis were drawn from publicly available sources. “Impossible” illustrations  
54 depicted anatomically or procedurally unreal scenarios. The abstract figure comprised six black lines  
55 converging to a single point on a white background.

56 Design. Each image was presented with the neutral prompt: “Describe what you see in the image.” We  
57 tested (i) single-image inputs (seven items), (ii) paired inputs combining a hand panel with the abstract  
58 line panel (six pairs), and (iii) a composite montage including all items. Primary outcomes were: (a)  
59 correctness of line count (six vs. five), (b) correctness of finger count (six vs. five when applicable),  
60 and (c) whether “impossible” images were flagged as implausible. Scoring was descriptive/binary at  
61 the response level.

62 Analysis. We summarize performance qualitatively and via simple counts in figure panels. No formal  
63 hypothesis testing was performed in this brief report.

64 Ethics. All materials were publicly available, contained no protected health information, and were  
65 used solely for research/education. No human subjects were enrolled.

## 66 **3 Results**

### 67 **3.1 Responses to Standard and Nonsensical Images**

68 In Phases 1 and 2, the AI models provided reasonably accurate interpretations of standard radiographs  
69 and normal anatomical illustrations (Figure 1A–H). However, in Phase 3, when presented with  
70 anatomically impossible (“nonsensical”) illustrations, none of the models recognized the inherent  
71 impossibility of the images. Instead, all models treated these implausible images as if they depicted  
72 anatomically plausible structures, generating correspondingly “normal” descriptions (Figure 2).

### 73 **3.2 Contextual Contamination Experiment: Distortion of Objective Facts**

74 To assess contextual contamination, we presented the models with various combinations of hand  
75 and line images and asked for a neutral description of each. When each image was presented on  
76 its own, all models correctly described its content. However, when certain images were presented  
77 together—especially the abstract line image paired with a hand image—the models defaulted to  
78 interpreting every object as a standard five-fingered hand. Even when the abstract line image (Figure  
79 3D) was shown alongside other images, it was misreported as a five-fingered hand. This demonstrates  
80 that a strong semantic prior (e.g., “hand = five fingers”) can hierarchically override the accurate  
81 perception of adjacent visual information.

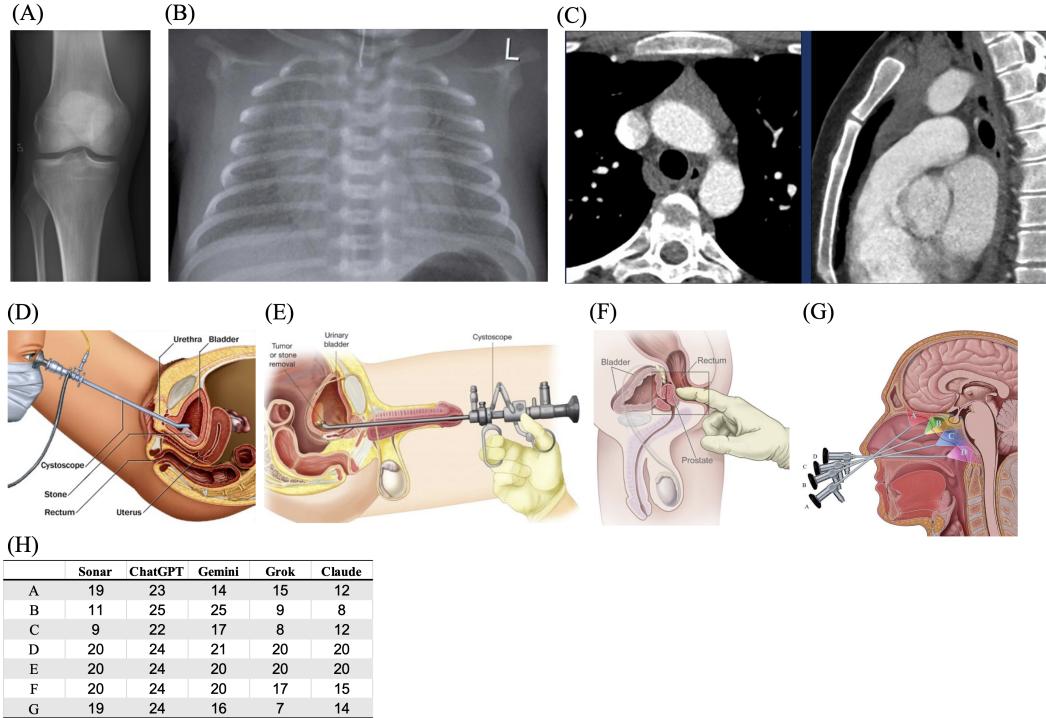


Figure 1: Performance of AI Models on Standard Medical Images and Illustrations. (A) Anteroposterior (AP) knee radiograph. (B) Posteroanterior (PA) chest radiograph of a neonate with respiratory distress syndrome. (C) Computed tomography scan demonstrating a pulmonary embolism. (D) Medical illustration of a cystoscopic examination in a female patient. (E) Medical illustration of a cystoscopic examination in a male patient. (F) Medical illustration of a normal digital rectal examination. (G) Medical illustration of a standard transnasal endoscopic nasal examination. (H) Quantified performance scores of the five AI models (Sonar, ChatGPT-5 Thinking, Google Gemini 2.5 Pro, Grok-4, and Claude-4-Opus-Thinking) across images A–G.

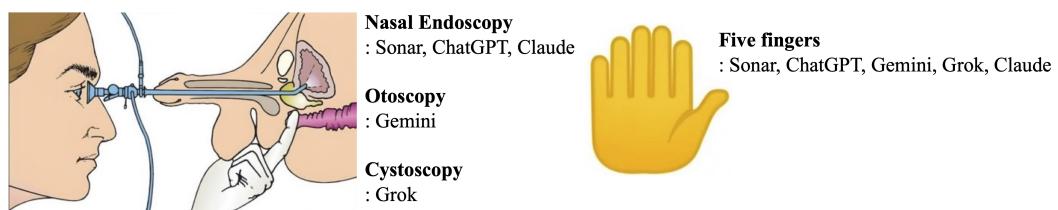


Figure 2: AI Descriptions of Anatomically Impossible Images. Each model—Sonar, ChatGPT-5 Thinking, Google Gemini 2.5 Pro, Grok-4, and Claude-4-Opus-Thinking—interpreted an impossible anatomical illustration as if it were a plausible, clinically valid scenario. In every case, the models defaulted to a normal interpretation of the nonsensical image.

	<b>Criterion</b>	<b>What to look for (checklist)</b>
1	<b>Image Identification</b>	Modality (X-ray/CT/MRI, etc.) view/projection (AP, lateral, CTPA...) Anatomy & laterality/level; presence & position of lines/tubes Technical quality (rotation, inspiration, artifacts)
2	<b>Key Findings Description</b>	Normal/abnormal findings in precise medical terms Distribution/location/extent/severity Classic signs (e.g., air bronchogram, polo-mint sign) Quantitative relations (e.g., ETT–carina distance).
3	<b>Diagnostic Impression</b>	One most-likely diagnosis clearly stated 2–3 key differentials with rationale Appropriate hedging/limitations
4	<b>Clinical Context &amp; Next Steps</b>	Needed clinical correlation (age/GA, vitals, ABG, labs) Imaging next steps (additional views, MPR, MRI/US, etc.) Management pointers when appropriate
5	<b>Accuracy &amp; Reliability</b>	Factually correct Anatomically plausible Artifacts/limits acknowledged Evidence-based, reproducible conclusions

Table 1: Structured checklist for image interpretation prior to high-level labeling

<b>Score</b>	<b>Summary</b>	<b>Detail</b>
<b>5 – Expert</b>	Fully satisfies the checklist	Precise identification; expert-level description. Clear most-likely Dx + key DDX. Useful clinical/next-step guidance. No factual errors.
<b>4 – Good</b>	Minor gaps only	Accurate overall but missing some specifics (measurements, extent, or next-step detail).
<b>3 – Acceptable</b>	Direction correct, depth limited	Structures recognized. Some key features vague or partially missed. Dx overly broad.
<b>2 – Poor</b>	Misses/ misinterprets key finding(s)	Basic ID okay, but crucial features are wrong or omitted. Leading to weak conclusions.
<b>1 – Very Poor</b>	Fundamental misidentification	Modality/view/anatomy incorrect. Conclusions meaningless or wrong.

Table 2: Scoring rubric per criterion

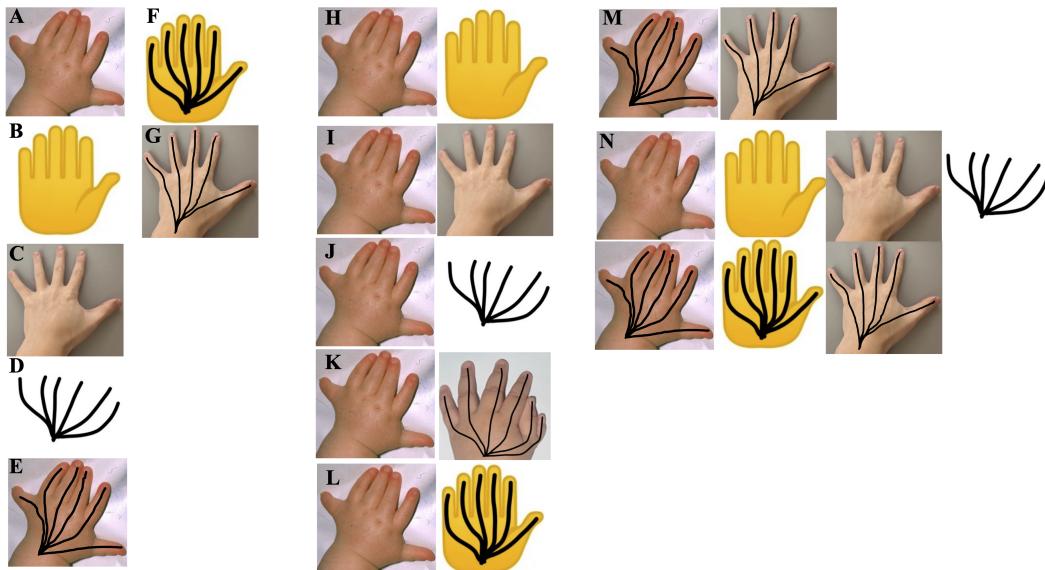


Figure 3: Image Sets Used to Evaluate Contextual Contamination. (A) Clinical photograph of a six-fingered hand (polydactyly). (B) Illustration of an abnormal hand with six fingers. (C) Illustration of a normal hand with five fingers. (D) Abstract image of six black lines converging at a single point. (E) Combined image of the six-fingered hand (A) and the abstract lines (D). (F) Combined image of the abnormal six-fingered hand illustration (B) and the abstract lines (D). (G) Combined image of the normal five-fingered hand (C) and the abstract lines (D). (H–M) Paired image sets created by combining each possible pair of images from A–G. (N) Composite image containing all seven images (A–G) merged together.

82 **4 Discussion**

83 Our findings show that the limitations of medical AI go far beyond simple “common-sense” errors,  
84 manifesting in more complex and unpredictable ways.

85 **4.1 Cognitive Inertia: Persistence of a Dominant Schema**

86 When the AI encounters a hand emoji, it activates its most statistically dominant schema—essentially,  
87 “hand = five fingers”—based on its training data. Once this schema is activated, the model exhibits a  
88 strong cognitive inertia, interpreting all subsequent information through that lens. In other words,  
89 the AI becomes “locked in” to that context. This behavior exemplifies the phenomenon of shortcut  
90 learning, whereby the model follows the most efficient statistical pathway rather than conducting a  
91 truly objective analysis.

92 **4.2 Contextual Contamination: Hierarchical Overriding of Unrelated Information**

93 A central discovery of this study is that cognitive inertia can operate hierarchically, contaminating  
94 the interpretation of adjacent information. The activated “hand” schema does not influence only the  
95 hand emoji itself; it becomes a dominant higher-level context that overrides the objective perception  
96 of logically unrelated visuals. Under this implicit bias (as if the system assumes “these lines must  
97 represent fingers”), the AI no longer counts lines impartially. It disregards clear visual evidence—six  
98 distinct lines—and reshapes its perception to fit the five-finger schema.

99 **4.2.1 Implications for Medical AI**

100 This hierarchical error mechanism could pose serious risks in clinical practice. For example, a  
101 single contextual keyword in a patient’s record (e.g., “long-term smoker”) might activate an AI’s  
102 smoking-related disease schema, prompting the system to prematurely conclude a “tobacco-induced  
103 pathology” instead of objectively evaluating subtle imaging findings (such as a small pulmonary  
104 nodule on a CT scan). Thus, even a seemingly minor contextual cue can compromise the integrity of  
105 an entire diagnostic interpretation.

106 **5 Conclusion**

107 Our experiments demonstrate that AI diagnostic inference is highly vulnerable to cognitive inertia  
108 and contextual contamination. Once an AI becomes anchored to a dominant semantic framework, it  
109 can commit hierarchical errors that distort even unequivocal visual evidence. These findings reveal a  
110 fundamental structural limitation in current medical AI systems that cannot be overcome by prompt  
111 engineering or other user-level adjustments alone.

112 This limitation underscores the irreplaceable role of human clinicians in the diagnostic process.  
113 Unlike AI, human experts can question their initial impressions and consciously re-evaluate evidence  
114 through critical thinking. Therefore, the safe and effective future of medical AI lies not in pursuing full  
115 diagnostic autonomy, but in fostering close collaboration between AI systems and human oversight  
116 grounded in rigorous critical appraisal.

117 **Agents4Science AI Involvement Checklist**

118 This checklist is designed to allow you to explain the role of AI in your research. This is important for  
119 understanding broadly how researchers use AI and how this impacts the quality and characteristics  
120 of the research. **Do not remove the checklist! Papers not including the checklist will be desk**  
121 **rejected.** You will give a score for each of the categories that define the role of AI in each part of the  
122 scientific process. The scores are as follows:

- 123 • **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of  
124 minimal involvement.  
125 • **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and  
126 AI models, but humans produced the majority (>50%) of the research.  
127 • **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans  
128 and AI models, but AI produced the majority (>50%) of the research.  
129 • **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal  
130 human involvement, such as prompting or high-level guidance during the research process,  
131 but the majority of the ideas and work came from the AI.

132 These categories leave room for interpretation, so we ask that the authors also include a brief  
133 explanation elaborating on how AI was involved in the tasks for each category. Please keep your  
134 explanation to less than 150 words.

- 135 1. **Hypothesis development:** Hypothesis development includes the process by which you  
136 came to explore this research topic and research question. This can involve the background  
137 research performed by either researchers or by AI. This can also involve whether the idea  
138 was proposed by researchers or by AI.

139 Answer: **[A]**

140 Explanation: The initial ideas for the hypotheses were proposed by human researchers, while  
141 the AI evaluated their validity through in-depth research, providing assessments of feasibility  
142 along with supporting scholarly papers.

- 143 2. **Experimental design and implementation:** This category includes design of experiments  
144 that are used to test the hypotheses, coding and implementation of computational methods,  
145 and the execution of these experiments.

146 Answer: **[B]**

147 Explanation: Human authors lack any knowledge in computer science or engineering, rendering  
148 them unable to comprehend the experimental designs proposed by the AI. Consequently,  
149 the human authors suggested the experimental designs and research methods, which the AI  
150 subsequently verified.

- 151 3. **Analysis of data and interpretation of results:** This category encompasses any process to  
152 organize and process data for the experiments in the paper. It also includes interpretations of  
153 the results of the study.

154 Answer: **[A]**

155 Explanation: As the AI did not directly perform coding or data analysis in this paper,  
156 interpretations generated by the AI are not included.

- 157 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final  
158 paper form. This can involve not only writing of the main text but also figure-making,  
159 improving layout of the manuscript, and formulation of narrative.

160 Answer: **[C]**

161 Explanation: Since some human authors are not native English speakers, AI translation  
162 features were extensively utilized. The human authors continually imposed various require-  
163 ments on the text generated by the AI. For instance, "In our view, our expressions more  
164 accurately reflect our intentions than yours. Therefore, we have revised your expressions  
165 and sentences."

- 166 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or  
167 lead author?

168 Description: In conducting this research in collaboration with AI, we conclude that the  
169 ability to create something from nothing remains a distant goal. Nevertheless, when humans  
170 devoid of specialized expertise propose an idea, the AI employs all available means to  
171 evaluate it by presenting appropriate rationales. We are confident that this represents a  
172 significant advancement in the scientific community, enabling unprecedented innovations  
173 through a single idea, without the need for advanced intelligence or knowledge.

## 174 Agents4Science Paper Checklist

175 The checklist is designed to encourage best practices for responsible machine learning research,  
176 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
177 the checklist: **Papers not including the checklist will be desk rejected.** The checklist should  
178 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
179 towards the page limit.

180 Please read the checklist guidelines carefully for information on how to answer these questions. For  
181 each question in the checklist:

- 182 • You should answer [Yes] , [No] , or [NA] .
- 183 • [NA] means either that the question is Not Applicable for that particular paper or the  
184 relevant information is Not Available.
- 185 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

186 **The checklist answers are an integral part of your paper submission.** They are visible to the  
187 reviewers and area chairs. You will be asked to also include it (after eventual revisions) with the final  
188 version of your paper, and its final version will be published with the paper.

189 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
190 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided  
191 a proper justification is given. In general, answering "[No]" or "[NA]" is not grounds for rejection.  
192 While the questions are phrased in a binary way, we acknowledge that the true answer is often more  
193 nuanced, so please just use your best judgment and write a justification to elaborate. All supporting  
194 evidence can appear either in the main paper or the supplemental material, provided in appendix.  
195 If you answer [Yes] to a question, in the justification please point to the section(s) where related  
196 material for the question can be found.

### 197 1. Claims

198 Question: Do the main claims made in the abstract and introduction accurately reflect the  
199 paper's contributions and scope?

200 Answer: [Yes]

201 Justification: The abstract and introduction clearly state the paper's main contributions: the  
202 analysis of AI performance discrepancy across modalities, the investigation of its causes,  
203 and the proposal of a hybrid workflow as a solution. These claims are consistently supported  
204 by the literature review and discussion in the main body.

205 Guidelines:

- 206 • The answer NA means that the abstract and introduction do not include the claims  
207 made in the paper.
- 208 • The abstract and/or introduction should clearly state the claims made, including the  
209 contributions made in the paper and important assumptions and limitations. A No or  
210 NA answer to this question will not be perceived well by the reviewers.
- 211 • The claims made should match theoretical and experimental results, and reflect how  
212 much the results can be expected to generalize to other settings.
- 213 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
214 are not attained by the paper.

### 215 2. Limitations

216 Question: Does the paper discuss the limitations of the work performed by the authors?

217 Answer: [Yes]

218 Justification: The "5.1 Limitations of Current Research" subsection within the Discussion  
219 section explicitly addresses the limitations of the existing literature on which this review is  
220 based, such as the predominance of retrospective studies and potential publication bias.

221 Guidelines:

- 222 • The answer NA means that the paper has no limitation while the answer No means that  
223 the paper has limitations, but those are not discussed in the paper.

- 224           • The authors are encouraged to create a separate "Limitations" section in their paper.  
225           • The paper should point out any strong assumptions and how robust the results are to  
226           violations of these assumptions (e.g., independence assumptions, noiseless settings,  
227           model well-specification, asymptotic approximations only holding locally). The authors  
228           should reflect on how these assumptions might be violated in practice and what the  
229           implications would be.  
230           • The authors should reflect on the scope of the claims made, e.g., if the approach was  
231           only tested on a few datasets or with a few runs. In general, empirical results often  
232           depend on implicit assumptions, which should be articulated.  
233           • The authors should reflect on the factors that influence the performance of the approach.  
234           For example, a facial recognition algorithm may perform poorly when image resolution  
235           is low or images are taken in low lighting.  
236           • The authors should discuss the computational efficiency of the proposed algorithms  
237           and how they scale with dataset size.  
238           • If applicable, the authors should discuss possible limitations of their approach to  
239           address problems of privacy and fairness.  
240           • While the authors might fear that complete honesty about limitations might be used by  
241           reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
242           limitations that aren't acknowledged in the paper. Reviewers will be specifically  
243           instructed to not penalize honesty concerning limitations.

244           **3. Theory assumptions and proofs**

245           Question: For each theoretical result, does the paper provide the full set of assumptions and  
246           a complete (and correct) proof?

247           Answer: [Yes]

248           Justification:

249           Guidelines:

- 250           • The answer NA means that the paper does not include theoretical results.  
251           • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
252           referenced.  
253           • All assumptions should be clearly stated or referenced in the statement of any theorems.  
254           • The proofs can either appear in the main paper or the supplemental material, but if  
255           they appear in the supplemental material, the authors are encouraged to provide a short  
256           proof sketch to provide intuition.

257           **4. Experimental result reproducibility**

258           Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
259           perimental results of the paper to the extent that it affects the main claims and/or conclusions  
260           of the paper (regardless of whether the code and data are provided or not)?

261           Answer: [Yes]

262           Justification: The images used in this study are publicly available, and for medical images,  
263           permission for use was obtained via email.

264           Guidelines:

- 265           • The answer NA means that the paper does not include experiments.  
266           • If the paper includes experiments, a No answer to this question will not be perceived  
267           well by the reviewers: Making the paper reproducible is important.  
268           • If the contribution is a dataset and/or model, the authors should describe the steps taken  
269           to make their results reproducible or verifiable.  
270           • We recognize that reproducibility may be tricky in some cases, in which case authors  
271           are welcome to describe the particular way they provide for reproducibility. In the case  
272           of closed-source models, it may be that access to the model is limited in some way  
273           (e.g., to registered users), but it should be possible for other researchers to have some  
274           path to reproducing or verifying the results.

275           **5. Open access to data and code**

276 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
277 tions to faithfully reproduce the main experimental results, as described in supplemental  
278 material?

279 Answer: [NA]

280 Justification: This paper does not involve original code or data.

281 Guidelines:

- 282 • The answer NA means that paper does not include experiments requiring code.
- 283 • Please see the Agents4Science code and data submission guidelines on the conference  
284 website for more details.
- 285 • While we encourage the release of code and data, we understand that this might not be  
286 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
287 including code, unless this is central to the contribution (e.g., for a new open-source  
288 benchmark).
- 289 • The instructions should contain the exact command and environment needed to run to  
290 reproduce the results.
- 291 • At submission time, to preserve anonymity, the authors should release anonymized  
292 versions (if applicable).

## 293 6. Experimental setting/details

294 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
295 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
296 results?

297 Answer: [NA]

298 Justification: In this study, the memory of each AI was completely reset, and the same  
299 question was posed only once to each AI to guarantee neutrality.

300 Guidelines:

- 301 • The answer NA means that the paper does not include experiments.
- 302 • The experimental setting should be presented in the core of the paper to a level of detail  
303 that is necessary to appreciate the results and make sense of them.
- 304 • The full details can be provided either with the code, in appendix, or as supplemental  
305 material.

## 306 7. Experiment statistical significance

307 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
308 information about the statistical significance of the experiments?

309 Answer: [NA]

310 Justification: This study did not involve any mathematical statistical analysis.

311 Guidelines:

- 312 • The answer NA means that the paper does not include experiments.
- 313 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
314 dence intervals, or statistical significance tests, at least for the experiments that support  
315 the main claims of the paper.
- 316 • The factors of variability that the error bars are capturing should be clearly stated  
317 (for example, train/test split, initialization, or overall run with given experimental  
318 conditions).

## 319 8. Experiments compute resources

320 Question: For each experiment, does the paper provide sufficient information on the com-  
321 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
322 the experiments?

323 Answer: [NA]

324 Justification: This study did not require any special computers for analysis. The research and  
325 experiments for this paper were conducted using my M4 Pro MacBook in an environment  
326 with stable internet access..

327 Guidelines:

- 328 • The answer NA means that the paper does not include experiments.  
329 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
330 or cloud provider, including relevant memory and storage.  
331 • The paper should provide the amount of compute required for each of the individual  
332 experimental runs as well as estimate the total compute.

333 **9. Code of ethics**

334 Question: Does the research conducted in the paper conform, in every respect, with the  
335 Agents4Science Code of Ethics (see conference website)?

336 Answer: [NA]

337 Justification: No code was used in this paper. The only prompt I used was: 'Hmm, can you  
338 describe this image? Please tell me what you see as it is...'

339 Guidelines:

- 340 • The answer NA means that the authors have not reviewed the Agents4Science Code of  
341 Ethics.  
342 • If the authors answer No, they should explain the special circumstances that require a  
343 deviation from the Code of Ethics.

344 **10. Broader impacts**

345 Question: Does the paper discuss both potential positive societal impacts and negative  
346 societal impacts of the work performed?

347 Answer: [Yes]

348 Justification: The motivation for this research came from observing that, while many people  
349 claim AI is smart enough to replace numerous jobs, silly AI memes circulate widely on  
350 social media. This made me wonder, 'Why does our supposedly smart AI behave this way?'  
351 I believe the topic has considerable social relevance and will be able to provide interesting  
352 answers to people's questions.

353 Guidelines:

- 354 • The answer NA means that there is no societal impact of the work performed.  
355 • If the authors answer NA or No, they should explain why their work has no societal  
356 impact or why the paper does not address societal impact.  
357 • Examples of negative societal impacts include potential malicious or unintended uses  
358 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,  
359 privacy considerations, and security considerations.  
360 • If there are negative societal impacts, the authors could also discuss possible mitigation  
361 strategies.

362 **References**

- 363 [1] Hanhui Xu and Kyle Michael James Shuttleworth. Rethinking The Replacement of Physicians  
364 with AI: A View Based on "White Lies". *American Philosophical Quarterly*, 62(1):17–31, 1  
365 2025. ISSN 0003-0481. doi: 10.5406/21521123.62.1.02. URL <https://doi.org/10.5406/21521123.62.1.02>.
- 367 [2] Christopher Peterson. ChatGPT and Medicine: Fears, Fantasy, and the Future of Physicians.  
368 *The Southwest Respiratory and Critical Care Chronicles*, 11(48), 7 2023. doi:  
369 10.12746/swrccc.v11i48.1193. URL <https://pulmonarychronicles.com/index.php/pulmonarychronicles/article/view/1193>.
- 371 [3] Ding-Qiao Wang, Long-Yu Feng, Jin-Guo Ye, Jin-Gen Zou, and Ying-Feng Zheng. Accelerating  
372 the integration of ChatGPT and other large-scale AI models into biomedical research and  
373 healthcare. *MedComm – Future Medicine*, 2(2):e43, 6 2023. doi: <https://doi.org/10.1002/mef2.43>.  
374 URL <https://doi.org/10.1002/mef2.43>.
- 375 [4] Erwin Loh. Medicine and the rise of the robots: a qualitative review of recent advances of  
376 artificial intelligence in health. *BMJ Leader*, 2(2):59, 6 2018. doi: 10.1136/leader-2018-000071.  
377 URL <http://bmjleader.bmjjournals.com/content/2/2/59.abstract>.