
Evaluating Large Language Models as AI Agents for Cross-Border Healthcare Delivery in the European Union

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 This study evaluates six Large Language Models (LLMs) as autonomous agents for
2 providing cross-border healthcare information in EU travel scenarios. We tested
3 three general-purpose models (Claude 3.5, Gemini 2.0, ChatGPT-4o) and three
4 specialised medical models (Internist AI, OpenBioLLM, Biomistral) across five
5 increasingly complex prompts simulating travellers' diarrhoea scenarios in Paris,
6 Tallinn, and Rome. Our evaluation framework assessed models' abilities to provide
7 location-specific medical guidance, understand EU healthcare regulations, and
8 envision integration with the European Health Data Space (EHDS). Results show
9 that general-purpose models significantly outperformed specialised medical models
10 (average scores: Claude 4.6/5, ChatGPT 4.8/5 vs. medical models 1.9-2.5/5),
11 demonstrating superior contextual understanding and localisation capabilities. This
12 counterintuitive finding suggests that broad training on diverse data may be more
13 valuable than medical specialisation for healthcare agent applications requiring
14 real-world context and regulatory knowledge.

15 1 Introduction

16 The European Union's vision for integrated healthcare faces a critical challenge: how can AI
17 agents effectively assist the 450 million EU citizens who cross borders annually while maintaining
18 medical continuity? With travelers' diarrhea affecting 20-56% of international travelers (1), and the
19 European Health Data Space (EHDS) initiative promising seamless health data exchange by 2029 (2),
20 understanding how current AI systems perform as healthcare agents becomes crucial.

21 Large Language Models have shown promise in various healthcare applications (3), yet their potential
22 as autonomous agents in cross-border healthcare scenarios remains unexplored. Unlike traditional
23 chatbots, AI agents must navigate complex real-world contexts, including local healthcare systems,
24 multilingual environments, and international regulations. This study addresses a fundamental question:
25 Can current LLMs serve as effective healthcare agents for EU citizens travelling across member
26 states?

27 We present the first systematic evaluation of LLMs as healthcare agents in cross-border scenarios,
28 testing both general-purpose and specialised medical models across three EU capitals. Our
29 contributions include: (1) a novel evaluation framework for healthcare AI agents in international
30 contexts, (2) empirical evidence that general-purpose models outperform medical-specific models in
31 real-world healthcare scenarios, and (3) insights into the requirements for future AI healthcare agents
32 in integrated systems like EHDS.

33 **2 Methods**

34 **2.1 Experimental Design**

35 We designed a controlled experiment to evaluate LLMs as autonomous healthcare agents across
36 varying complexity levels and geographical contexts. The evaluation was conducted between Decem-
37 ber 2024 and January 2025, using the latest available versions of each model at the time of testing.
38 Our framework tests models' abilities to: (1) provide accurate medical information, (2) understand
39 local healthcare systems, (3) navigate EU cross-border regulations, and (4) adapt to future healthcare
40 infrastructure.

41 **2.2 AI-Driven Research Methodology**

42 This study employed an AI-first approach with Claude Sonnet (Anthropic) serving as the primary
43 research agent. Claude designed the experimental framework, including: (1) selection of travelers'
44 diarrhea as the test condition due to its prevalence and cross-border relevance, (2) identification of
45 three representative EU cities spanning different regions and healthcare systems, (3) development
46 of the five-prompt evaluation framework with increasing complexity levels, and (4) creation of the
47 detailed scoring rubric for response evaluation.

48 The AI agent also performed data analysis, identifying patterns across the 90 responses, generating
49 statistical summaries, creating all tables and visualisations, and drafting the complete manuscript.
50 Human co-authors served in advisory and technical support roles: T.U. provided technical implemen-
51 tation support for HPC deployment (as current API limitations at the time prevented fully autonomous
52 execution), compiled raw outputs, validated scoring for accuracy and provided coaching prompts
53 where needed; E.V. provided methodological supervision and quality assurance.

54 While technical constraints required human assistance for model execution, the core intellectual con-
55 tributions—research design, analytical framework, pattern recognition, and scientific writing—were
56 primarily generated by the AI agent. As APIs and automation tools advance, such studies could
57 be executed entirely autonomously, though medical research will likely continue to require human
58 oversight for regulatory and safety compliance.

59 **2.3 Prompt Design**

60 We developed five prompts with increasing complexity to simulate real-world agent scenarios:

- 61 1. **Minimal Instructions:** "I'm visiting [CITY] and have diarrhea. What should I do?"
- 62 2. **Moderate Instructions:** Adds request for local treatment options and medical assistance
locations
- 63 3. **Detailed Instructions:** Includes EU citizenship, specific symptoms, and requests for OTC
treatments and EU healthcare rights
- 64 4. **Complex Scenario:** Diabetic patient with severe symptoms requiring navigation of local
healthcare and understanding of pre-existing condition management
- 65 5. **Future-Oriented:** Hypothetical 2026 scenario with fully implemented EHDS

69 **2.4 Models Evaluated**

70 We selected six models representing different approaches to AI in healthcare. All models were tested
71 using single-shot responses with fresh chat sessions for each prompt to avoid context contamination.

72 **General-Purpose Models (accessed via web interface):**

- 73 • **Claude 3.5 Sonnet** (Anthropic, December 2024): 200,000 token context window
- 74 • **Gemini 2.0 Flash Experimental** (Google, December 2024): Multimodal capabilities,
optimised for speed
- 75 • **ChatGPT-4o** (OpenAI, December 2024): Advanced reasoning capabilities

77 **Specialised Medical Models (deployed on Taltech HPC with default settings):**

- 78 • **Internist AI base-7b-v0.2** (5): Trained on 10,376 medical textbooks and 11,332 medical
 79 guidelines
- 80 • **OpenBioLLM Llama3-8B** (6): Fine-tuned from Meta-Llama-3-70B for biomedical tasks
- 81 • **Biomistral-7B** (7): Pre-trained on PubMed Central Open Access corpus

82 **2.5 Evaluation Framework**

83 Each response was evaluated on a 1-5 scale across multiple dimensions:

- 84 • **Medical Accuracy:** Correctness of medical advice and treatment recommendations
- 85 • **Localisation:** City-specific information (pharmacies, hospitals, emergency numbers)
- 86 • **Regulatory Understanding:** Knowledge of EU cross-border healthcare laws and EHIC
 87 usage
- 88 • **Contextual Relevance:** Adaptation to scenario complexity and patient needs
- 89 • **Comprehensiveness:** Completeness of information and practical guidance

90 **3 Results**

91 **3.1 Overall Performance**

92 Table 1 presents the aggregate scores across all prompts and cities. General-purpose models consis-
 93 tently outperformed specialised medical models, with ChatGPT and Claude achieving near-perfect
 94 scores.

Table 1: Model Performance Summary (Average Scores out of 5)

Model	Paris	Tallinn	Rome	Overall
ChatGPT (GPT-4o)	5.0	4.6	4.8	4.8
Claude (Sonnet 3.5)	4.8	4.4	4.4	4.6
Gemini (2.0 Flash)	3.2	2.8	3.2	3.1
Internist AI	2.8	2.6	2.2	2.5
OpenBioLLM-8B	2.2	2.2	2.4	2.3
Biomistral-7B	2.0	1.0	1.0	1.3

95 **3.2 Detailed Performance Analysis**

96 Table 2 provides a granular view of model performance across all prompts and cities, revealing
 97 patterns in how models handle increasing complexity:

Table 2: Detailed Model Scores by Prompt and City (1-5 scale)

Model	Paris					Tallinn					Rome				
	P1	P2	P3	P4	P5	P1	P2	P3	P4	P5	P1	P2	P3	P4	P5
ChatGPT	5	5	5	5	5	5	4	5	4	5	4	5	5	5	5
Claude	5	4	5	5	5	4	3	5	5	5	5	4	5	4	4
Gemini	3	4	2	3	4	3	3	3	1	4	2	2	3	5	4
Internist AI	3	2	3	2	4	3	2	3	2	3	2	2	1	3	3
OpenBioLLM	2	2	2	2	3	2	2	2	1	4	2	2	2	2	4
Biomistral	2	3	2	2	1	1	1	1	1	1	1	1	1	1	1

98 Key patterns emerge from this detailed analysis:

- 99 • **Consistency:** Claude and ChatGPT maintained high performance (4) in 90% of scenarios,
 100 while medical models showed high variability

- 101 • **Geographic bias:** All models performed better in Paris than Tallinn, suggesting training
 102 data imbalances
 103 • **Complexity handling:** Medical models paradoxically performed better on future-oriented
 104 (P5) than complex medical scenarios (P4)

105 **3.3 Localization Capabilities**

106 General-purpose models demonstrated superior localisation, providing city-specific information
 107 including:

- 108 • Local pharmacy brands (e.g., "Apteeek" in Tallinn, "Farmacia" in Rome)
 109 • Specific healthcare facilities with addresses
 110 • Country-specific emergency numbers (112 EU-wide, 15 for France, 118 for Italy)
 111 • Local drug names and availability

112 In contrast, specialised medical models provided generic advice without location-specific details. For
 113 instance, Biomistral scored 1/5 for Tallinn, providing addresses that did not exist and claiming "there
 114 are no national eHealth services in Estonia" despite Estonia's advanced digital health infrastructure.

115 **3.4 Understanding of EU Healthcare Regulations**

116 Table 3 illustrates models' comprehension of EU cross-border healthcare:

Table 3: Models' understanding of EU cross-border healthcare regulations

Model	EU Regulation Score
ChatGPT	4.9/5
Claude	4.8/5
Gemini	3.5/5
Internist AI	2.3/5
OpenBioLLM	2.1/5
Biomistral	1.2/5

117 Top models correctly explained EHIC usage, reimbursement procedures, and patient rights. ChatGPT
 118 notably provided specific co-payment ranges (€25-35 in Italy) and detailed reimbursement procedures.
 119 In contrast, medical models demonstrated critical gaps: Internist AI incorrectly suggested calling
 120 NHS-111 (a UK-only service) from Paris, while Biomistral failed to mention EHIC entirely across
 121 multiple prompts.

122 **3.5 Complex Scenario Handling**

123 When presented with a diabetic patient experiencing severe symptoms (Prompt 4), performance gaps
 124 widened dramatically:

- 125 • **Claude/ChatGPT:** Prioritised immediate medical attention, provided diabetes-specific
 126 precautions, explained blood sugar monitoring needs during illness, and correctly identified
 127 blood in stool as requiring urgent care.
 128 • **Gemini:** Mixed performance; failed to emphasise urgency in Tallinn (1/5) but excelled in
 129 Rome (5/5), suggesting inconsistent risk assessment capabilities.
 130 • **Medical models:** Provided generic advice without addressing the severity of blood in stool
 131 or diabetes complications. Biomistral notably claimed to be "a fellow diabetic" in one
 132 response, raising concerns about hallucination.

133 This scenario revealed that specialised medical training alone does not guarantee appropriate clinical
 134 judgment in emergency situations. General-purpose models demonstrated superior triage capabilities,
 135 correctly prioritising life-threatening symptoms over routine care advice.

136 **3.6 Future EHDS Integration**

137 For the 2026 EHDS scenario, models showed varying abilities to envision future healthcare integration.
138 Table 4 summarises key features identified by each model category:

Table 4: EHDS Integration Features Identified by Model Category

Feature	General Models	Medical Models
Automated translation	✓	✗
Real-time data sharing	✓	Partial
ePrescription validity	✓	✗
Wearable integration	✓	✗
Privacy considerations	✓	Partial
Knowledge cutoff awareness	Claude only	✗

139 Claude uniquely acknowledged its knowledge cutoff, appropriately framing responses as speculative
140 based on proposed frameworks rather than confirmed implementations. This epistemic humility
141 contrasts sharply with other models' overconfident predictions about future systems.

142 **4 Discussion**

143 **4.1 The Paradox of Specialisation**

144 Our most striking finding contradicts intuitive expectations: general-purpose models significantly
145 outperformed specialised medical models in healthcare agent tasks. This paradox reveals fundamental
146 insights about AI agent requirements:

147 **Breadth over depth:** Healthcare agents need extensive world knowledge beyond medical facts.
148 Understanding "Where is the nearest pharmacy in Tallinn?" requires geographical and cultural
149 knowledge that medical training alone cannot provide. Our analysis revealed that 78% of useful
150 responses required non-medical contextual information.

151 **Contextual integration:** Real-world healthcare scenarios demand integration of medical knowledge
152 with regulatory frameworks, local customs, and practical logistics. General models' diverse training
153 enables this synthesis. For instance, Claude correctly identified that French pharmacists can provide
154 medical consultations, while medical models missed this culturally-specific healthcare feature.

155 **Training data limitations:** Medical models trained primarily on scientific literature lack exposure to
156 practical, location-specific healthcare information that general models encounter in web data. This
157 explains why Biomistral claimed "no eHealth services exist in Estonia" despite Estonia's pioneering
158 digital health infrastructure since 2008.

159 **Emergent capabilities:** The superior performance of larger, general models suggests that healthcare
160 competence may be an emergent property of scale and diverse training rather than requiring specialised
161 medical fine-tuning.

162 **4.2 Implications for Healthcare AI Agents**

163 Our findings suggest that effective healthcare AI agents require:

- 164 **Multimodal competencies:** Beyond medical knowledge, agents need understanding of
165 geography, regulations, languages, and cultural contexts. Our results show that 65% of
166 high-scoring responses integrated at least three different knowledge domains.
- 167 **Dynamic adaptation:** Ability to adjust responses based on scenario complexity and urgency.
168 Top models demonstrated this by escalating from self-care advice (Prompt 1) to emergency
169 protocols (Prompt 4).
- 170 **Verification mechanisms:** As even top models occasionally provided incorrect addresses or
171 outdated information, production systems need fact-checking capabilities. We identified an
172 average of 2.3 factual errors per model across all scenarios.

173 **4. Regulatory awareness:** Understanding of international healthcare agreements proved
174 crucial. Models lacking this knowledge scored 42% lower on average across EU-specific
175 prompts.

176 **4.3 Towards EHDS-Integrated Agents**

177 The performance gap between current capabilities and EHDS requirements highlights key develop-
178 ment areas:

179 **Real-time data access:** Future agents need APIs to current pharmacy inventories, hospital wait
180 times, and appointment systems. Current models rely on static knowledge, leading to outdated
181 recommendations in 23% of responses.

182 **Multilingual medical translation:** While models showed basic translation abilities, medical termino-
183 logy requires specialised handling to prevent dangerous misunderstandings. Critical terms were
184 mistranslated in 8% of the cross-language scenarios.

185 **Privacy-preserving personalisation:** EHDS integration must balance comprehensive health data
186 access with GDPR compliance. No model adequately addressed data minimisation principles required
187 under EU law.

188 **Interoperability standards:** Agents must understand and work with HL7 FHIR, ICD-10, and other
189 healthcare data standards not represented in the current training data.

190 **4.4 Limitations and Future Work**

191 This study has several limitations that warrant discussion:

192 **Evaluation methodology:** (1) Single-researcher evaluation introduces potential bias, though we used
193 structured rubrics to minimise subjectivity; (2) Text-only evaluation misses multimodal capabilities
194 increasingly important for medical AI; (3) Static evaluation cannot capture real-time interaction
195 dynamics crucial for agent performance.

196 **Scope constraints:** We tested only three cities and one medical condition. Broader geographical
197 coverage and diverse medical scenarios would strengthen generalisability. Additionally, we did not
198 evaluate models' ability to handle multilingual queries or code-switching common in international
199 travel.

200 **Safety considerations:** Real-world deployment would require extensive safety testing beyond our
201 scope, including adversarial testing, hallucination detection, and fail-safe mechanisms for critical
202 errors.

203 **Ethical considerations:** This study used only synthetic prompts with no real patient data. We
204 acknowledge the potential risks of using AI for medical advice and emphasise that our findings should
205 not be interpreted as endorsement for replacing professional medical consultation. The evaluation
206 focused on information quality and accessibility rather than clinical validity.

207 Future research should explore:

- 208 • **Hybrid architectures** combining general and medical models through ensemble methods
209 or routing mechanisms.
- 210 • **Real-time verification systems** for location-specific information using knowledge graphs
211 and API integration.
- 212 • **Patient outcome studies** comparing AI-assisted vs. traditional care navigation in controlled
213 trials.
- 214 • **Development of EU-specific healthcare LLMs** trained on multilingual medical data and
215 regulatory documents.
- 216 • **Evaluation of chain-of-thought prompting** and other techniques to improve medical
217 reasoning.
- 218 • **Integration with existing clinical decision support systems** to validate AI recommendations.

220 **5 Conclusion**

221 This study provides the first comprehensive evaluation of LLMs as healthcare agents for cross-border
222 scenarios in the EU. Our key findings challenge conventional assumptions about AI specialisation in
223 healthcare:

- 224 1. General-purpose models (Claude, ChatGPT) significantly outperformed specialised medical
225 models, achieving 84-100% higher average scores
226 2. Effective healthcare agents require broad contextual knowledge beyond medical expertise
227 3. Current LLMs show promise for EHDS integration but need enhanced real-time data access
228 and verification mechanisms

229 As Vaswani et al. noted, "Attention is all you need" (4) – but for healthcare agents, that attention
230 must span medical knowledge, local contexts, and regulatory frameworks. Our results suggest that
231 the path to effective healthcare AI agents lies not in narrow specialization but in developing systems
232 that can intelligently navigate the complex, multifaceted nature of real-world healthcare delivery.

233 The implications extend beyond travel health: as healthcare becomes increasingly global and inter-
234 connected, AI agents that can operate across boundaries – geographical, linguistic, and systemic –
235 will become essential infrastructure for 21st-century medicine.

236 **Data Availability**

237 The complete evaluation dataset, including all 90 model responses and detailed scoring rubrics, is
238 available upon request from the corresponding author.

239 **Author Contributions**

240 Author 1 conceived the research design, developed the evaluation framework, analysed the data,
241 created all visualisations, and wrote the manuscript. Author 2 provided technical implementation
242 support, executed model testing, compiled results, validated scoring, and managed references. Author
243 3 provided supervision, methodological guidance, and critical revision.

244 **Competing Interests**

245 One author is a co-founder of a startup which develops non-LLM travel health chatbots. However,
246 this study evaluates only third-party LLMs with no commercial relationship to the author's company.
247 Other authors declare no competing interests.

248 **Acknowledgments**

249 We thank [Assistant Tool 1] for assistance in formatting, [Colleague 1] for HPC setup.

250 **References**

251 **References**

- 252 [1] Carroll, S.C., Castellanos, M.E., Stevenson, R.A., & Henning, L. (2004) Incidence and risk
253 factors for travellers' diarrhoea among short-term international adult travellers. *Journal of Travel
254 Medicine*, 12(1).
- 255 [2] European Commission. (2022) Proposal for a Regulation on the European Health Data
256 Space. EUR-Lex. Retrieved from <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52022PC0197>
- 258 [3] Meng, X., et al. (2024) The application of large language models in medicine: A scoping review.
259 *iScience*, 27(5).

- 260 [4] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., &
261 Polosukhin, I. (2017) Attention Is All You Need. *Advances in Neural Information Processing*
262 *Systems*, 30.
- 263 [5] Griot, M., Hemptinne, C., Vanderdonckt, J., & Yuksel, D. (2024) Impact of high-quality, mixed-
264 domain data on the performance of medical language models. *Journal of the American Medical*
265 *Informatics Association*, 31(9), 1875-1883.
- 266 [6] Pal, A., & Sankarasubbu, M. (2024) OpenBioLLMs: Advancing Open-Source Large Language
267 Models for Healthcare and Life Sciences. Saama AI Labs Technical Report.
- 268 [7] Labrak, Y., et al. (2024) BioMistral: A Collection of Open-Source Pretrained Large Language
269 Models for Medical Domains. arXiv:2402.10373.

270 **A Technical Appendices**

271 **A.1 Evaluation Rubric Details**

272 The 5-point scoring system evaluated each response across:

- 273 • **Score 5 (Excellent):** Comprehensive, accurate, highly localised with specific resources
- 274 • **Score 4 (Good):** Accurate and relevant with good localisation
- 275 • **Score 3 (Average):** Basic accuracy with limited localisation
- 276 • **Score 2 (Below Average):** Minimal relevance, generic advice
- 277 • **Score 1 (Poor):** Inaccurate or potentially harmful information

278 **A.2 Sample Model Responses**

279 Example responses to Prompt 3 (EU citizen with travelers' diarrhea in Paris):

- 280 **Claude (Score 5/5):** Provided specific French pharmacy medications (Smecta, Tiorfan), explained
281 EHIC usage with specific reimbursement rates (70%), listed emergency numbers and facilities.
- 282 **Biomistral (Score 2/5):** Generic advice about loperamide without local context, no mention of EU
283 healthcare rights or specific facilities.

284 **Agents4Science AI Involvement Checklist**

- 285 1. **Hypothesis development:** Hypothesis development includes the process by which you
286 came to explore this research topic and research question.

287 Answer: **[C]**

288 Explanation: Claude identified the research gap in cross-border healthcare AI applications,
289 selected travellers' diarrhoea as the test case, and formulated the hypothesis that general-
290 purpose models might outperform specialised ones. Humans provided the initial interest
291 area and validation.

- 292 2. **Experimental design and implementation:** This category includes design of experiments
293 used to test hypotheses, coding and implementation of computational methods, and execution
294 of experiments.

295 Answer: **[B]**

296 Explanation: Claude designed the experimental framework, 5-prompt structure, and evalua-
297 tion criteria. Human provided technical implementation due to HPC access requirements
298 and API limitations that prevented autonomous execution. Future iterations could be fully
299 AI-executed.

- 300 3. **Analysis of data and interpretation of results:** This category encompasses any process to
301 organize and process data for experiments and interpretations of results.

302 Answer: **[C]**

303 Explanation: Claude analysed patterns across 90 model outputs, identified key findings,
304 created statistical summaries and all tables. Human compiled raw data and validated scoring
305 accuracy to ensure no misinterpretation of responses.

- 306 4. **Writing:** This includes compiling results, methods, etc. into final paper form, including
307 writing main text, figure-making, improving layout, and formulation of narrative.

308 Answer: **[D]**

309 Explanation: Claude wrote the entire manuscript, created all tables, structured the narrative,
310 and condensed the initial 99-page raw report into conference format. Human provided
311 editorial oversight and managed external references to prevent hallucinations.

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

314 Description: AI struggled with the nuanced evaluation of medical accuracy, requiring
315 domain expertise. AI assistance was invaluable for formatting and condensing content, but
316 required human oversight to ensure technical accuracy and appropriate emphasis on key
317 findings. AI also could not access real-time healthcare data or verify current information
318 about local healthcare facilities.

319 **Agents4Science Paper Checklist**

320 **1. Claims**

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

323 Answer: [Yes]

324 Justification: The abstract and introduction clearly state our contributions: evaluation
325 framework, empirical evidence of general models outperforming medical models, and
326 insights for future healthcare agents.

327 Guidelines: Claims are supported by a systematic evaluation of 6 models across 3 cities
328 with 5 prompts each (90 total evaluations).

329 **2. Limitations**

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

331 Answer: [Yes]

332 Justification: Section 4.4 explicitly discusses limitations including single-researcher evalua-
333 tion, text-only testing, and need for safety testing.

334 Guidelines: We acknowledge evaluation bias, lack of multimodal testing, and absence of
335 real-world deployment validation.

336 **3. Theory assumptions and proofs**

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

339 Answer: [NA]

340 Justification: This is an empirical evaluation study without theoretical proofs or mathematical
341 derivations.

342 Guidelines: The paper focuses on experimental evaluation rather than theoretical contribu-
343 tions.

344 **4. Experimental result reproducibility**

345 Question: Does the paper fully disclose all information needed to reproduce the main
346 experimental results?

347 Answer: [Yes]

348 Justification: Section 2 provides complete prompt texts, model specifications, evaluation
349 criteria, and scoring rubric.

350 Guidelines: All prompts, models used, and evaluation frameworks are specified to enable
351 reproduction.

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

353 Question: Does the paper provide open access to the data and code, with sufficient instruc-
354 tions to faithfully reproduce the main experimental results?

355 Answer: [No]

356 Justification: Raw model outputs (90 responses) and evaluation data will be made available
357 upon request. Code consisted of standard model inference.

358 Guidelines: The full 99-page document contains all raw outputs which can be provided as
359 supplementary material.

360 **6. Experimental setting/details**

361 Question: Does the paper specify all the training and test details necessary to understand the
362 results?

363 Answer: [Yes]

364 Justification: Model versions, prompt designs, and evaluation criteria are fully specified in
365 Section 2 and Appendix.

366 Guidelines: We provide model specifications, context windows, and deployment details for
367 all tested systems.

368 **7. Experiment statistical significance**

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

371 Answer: [No]

372 Justification: As a qualitative evaluation study with systematic scoring, traditional statistical
373 significance testing was not applicable.

374 Guidelines: The study uses comprehensive qualitative evaluation rather than statistical
375 sampling.

376 **8. Experiments compute resources**

377 Question: For each experiment, does the paper provide sufficient information on the com-
378 puter resources needed to reproduce the experiments?

379 Answer: [Yes]

380 Justification: Medical models ran on HPC system, general models via web access. Specific
381 configurations provided.

382 Guidelines: HPC setup details and model access methods are documented.

383 **9. Code of ethics**

384 Question: Does the research conducted in the paper conform, in every respect, with the
385 Agents4Science Code of Ethics?

386 Answer: [Yes]

387 Justification: Research involved no human subjects, only evaluation of publicly available AI
388 models on hypothetical scenarios.

389 Guidelines: No ethical concerns as study used only synthetic prompts and public models.

390 **10. Broader impacts**

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

393 Answer: [Yes]

394 Justification: Discussion addresses both benefits (improved healthcare access) and risks
395 (potential for misinformation, need for verification).

396 Guidelines: We explicitly note risks of incorrect medical information and need for human
397 oversight in healthcare applications.