

Accelerated Article Preview

Comparative benchmarking of the DeepSeek large language model on medical tasks and clinical reasoning

Received: 6 February 2025

Accepted: 21 April 2025

Accelerated Article Preview

Published online: 23 April 2025

Mickael Tordjman, Zelong Liu, Murat Yuce, Valentin Fauveau, Yunhao Mei, Jerome Hadjadj, Ian Bolger, Haidara Almansour, Carolyn Horst, Ashwin Singh Parihar, Amine Geahchan, Anis Meribout, Nader Yatim, Nicole Ng, Phillip Robson, Alexander Zhou, Sara Lewis, Mingqian Huang, Timothy Deyer, Bachir Taouli, Hao-Chih Lee, Zahi A. Fayad & Xueyan Mei

Cite this article as: Tordjman, M. et al. Comparative benchmarking of the DeepSeek large language model on medical tasks and clinical reasoning. *Nature Medicine* <https://doi.org/10.1038/s41591-025-03726-3> (2025).

This is a PDF file of a peer-reviewed paper that has been accepted for publication. Although unedited, the content has been subjected to preliminary formatting. Nature Medicine is providing this early version of the typeset paper as a service to our authors and readers. The text and figures will undergo copyediting and a proof review before the paper is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers apply.

1. Extended Data

Figure or Table # Please group Extended Data items by type, in sequential order. Total number of items (Figs. + Tables) must not exceed 10.	Figure/Table title One sentence only	Filename Whole original file name including extension. i.e.: Smith_ED_Fig1.jpg	Figure/Table Legend If you are citing a reference for the first time in these legends, please include all new references in the main text Methods References section, and carry on the numbering from the main References section of the paper. If your paper does not have a Methods section, include all new references at the end of the main Reference list.
Extended Data Table. 1	Number of parameters for each model according to Azure.	Extended_Data_Tables-1.jpg	
Extended Data Table. 2	Prompts used for the different tasks	Extended_Data_Tables-2.jpg	
Extended Data Table. 3	Comparison of the average scores for the linguistic evaluation of the “Findings to Impression” task on the 2 databases (private database of reports and MIMIC) based on the evaluation of 2 independent radiologists for each database.	Extended_Data_Tables-3.jpg	Number of tokens per summary provided by the model and ratio of the number of tokens for the Findings section prompted/ Impression output.
Extended Data Table. 4	RECIST 1.1 performances for the 3 LLMs per category of response based on the consensus between the 3 human experts (average of the 3 runs)	Extended_Data_Tables-4.jpg	CR= Complete Response; PR= Partial Response; SD= Stable Disease; PD= Progressive Disease

Extended Data Table. 5	Text-based case evaluation using 5-point Likert score (adapted from the R-IDEA score) for the reasoning evaluation of Complex NEJM Cases and Management evaluation of the BMJ cases; and accuracy for the NEJM and private database of multiple choice questions	Extended_DataTables-5.jpg	
Extended Data Table. 6	BERTscore (F1) between the 3 runs of each model for test-retest repeatability.	Extended_DataTables-6.jpg	
Extended Data Table. 7	Fleiss kappa between the 3 runs of each model for test-retest repeatability.	Extended_DataTables-7.jpg	
Extended Data Table. 8	Inter-model Cohen kappa based on the first run of each model for USMLE, RECIST, Medicilline questions, and NEJM diagnostic cases.	Extended_DataTables-8.jpg	

3 **1. Supplementary Information:**

4 **A. PDF Files**

Item	Present?	Filename Whole original file name including extension. i.e.: Smith_SI.pdf. The extension must be .pdf	A brief, numerical description of file contents. i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	Choose an item.		
Reporting Summary	Yes	NMED-FT139399B_RS.pdf	
Peer Review Information	Choose an item.	<i>OFFICE USE ONLY</i>	

5

6 **B. Additional Supplementary Files**

Type	Number Each type of file (Table, Video, etc.) should be numbered from 1 onwards. Multiple files of the same type should be listed in sequence, i.e.: Supplementary Video 1, Supplementary Video 2, etc.	Filename Whole original file name including extension. i.e.: Smith_Supplementary_Video_1.mov	Legend or Descriptive Caption Describe the contents of the file
-------------	--	---	---

7

8 **3. Source Data**

9

Parent Figure or Table	Filename Whole original file name including extension. i.e.: Smith_SourceData_Fig1.xls, or Smith_Unmodified_Gels_Fig1.pdf	Data description i.e.: Unprocessed western Blots and/or gels, Statistical Source Data, etc.
-------------------------------	---	---

10

11

12 Editor summary:
13 The open-source DeepSeek large language model showed variable performance relative to two leading
14 models when benchmarked on four different medical tasks, with relatively strong reasoning capabilities
15 but similar or weaker relative performance on other tasks, such as summarization of imaging reports.
16 Editor recognition statement:
17 Primary Handling Editors: Lorenzo Righetto, Michael Basson and Saheli Sadanand, in collaboration with
18 the Nature Medicine team.
19 Peer Review:
20 Nature Medicine thanks Jie Yang, Ahmed Alaa and Kirk Roberts for their contribution to the peer review
21 of this work.
22

ACCELERATED ARTICLE PREVIEW

23 **Comparative benchmarking of the DeepSeek large language model on medical**
24 **tasks and clinical reasoning**

25 Mickael Tordjman^{1,2,*}, Zelong Liu^{1,*}, Murat Yuce^{1,2}, Valentin Fauveau¹, Yunhao Mei³,
26 Jerome Hadjadj^{4,5}, Ian Bolger^{1,2}, Haidara Almansour⁶, Carolyn Horst⁷, Ashwin Singh
27 Parihar⁸, Amine Geahchan^{1,2}, Anis Meribout¹, Nader Yatim⁹, Nicole Ng¹⁰, Phillip
28 Robson^{1,2}, Alexander Zhou¹, Sara Lewis^{1,2}, Mingqian Huang^{1,2}, Timothy Deyer^{11,12},
29 Bachir Taouli^{1,2}, ^{1,2}#, Hao-Chih Lee^{1,2,#}, Zahi A. Fayad^{1,2,13,#}, Xueyan Mei^{1,2,13,#}

30

- 31 1. Biomedical Engineering and Imaging Institute, Icahn School of Medicine at Mount
32 Sinai, New York, NY, 10029, USA
33 2. Department of Diagnostic, Molecular and Interventional Radiology, Mount Sinai
34 Hospital, New York, NY, 10029, USA
35 3. Erasmus University Rotterdam, Rotterdam, Netherlands
36 4. Sorbonne Universite, Service de Medecine Interne, Hopital Saint Antoine, AP-HP,
37 Paris, France
38 5. Center for Human Genetics and Genomics, New York University Grossman School
39 of Medicine, New York, NY
40 6. Department of Diagnostic and Interventional Radiology, Tuebingen University
41 Hospital, Tuebingen, Germany
42 7. School of Biomedical Engineering and Imaging Sciences, King's College, London,
43 UK
44 8. Mallinckrodt Institute of Radiology, Washington University School of Medicine, Saint
45 Louis, MO, USA
46 9. Department of Immunology and Immunotherapy, Icahn School of Medicine at Mount
47 Sinai, New York, NY, 10029, USA
48 10. Department of Pulmonary, Critical Care and Sleep Medicine, Icahn School of
49 Medicine at Mount Sinai, New York, NY, 10029, USA
50 11. East River Medical Imaging, New York, NY, 10021, USA
51 12. Department of Radiology, Cornell Medicine, New York, NY, 10065, USA
52 13. Windreich Department of Artificial Intelligence and Human Health, Icahn School of
53 Medicine at Mount Sinai, 10029, New York, NY, USA

54
55 *These authors contributed equally to this work.

56 #These authors jointly supervised this manuscript.

57 **Corresponding authors:**

58 Hao-Chih Lee: hao-chih.lee@mssm.edu; Zahi A. Fayad: zahi.fayad@mssm.edu;
59 Xueyan Mei: xueyan.mei@mssm.edu. Corresponding address: 1470 Madison Avenue,
60 New York, New York, 10029, USA.

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86 **Abstract:**

87 DeepSeek is a newly introduced large language model (LLM) designed for enhanced
88 reasoning, but its medical-domain capabilities have not yet been evaluated. This study
89 assessed the capabilities of three LLMs—DeepSeek-R1, ChatGPT-o1, and Llama 3.1-
90 405B—in performing four different medical tasks: answering questions from the United
91 States Medical Licensing Examination (USMLE), interpreting and reasoning based on
92 text-based diagnostic and management cases, providing tumor classification according
93 to RECIST 1.1 criteria, and providing summaries of diagnostic imaging reports across
94 multiple modalities. In the USMLE test, the performance of DeepSeek-
95 R1(accuracy=0.92) was slightly inferior to that of ChatGPT-o1(accuracy=0.95; p=0.04)
96 but better than that of Llama 3.1-405B (accuracy=0.83; p<10⁻³). For text-based case
97 challenges, DeepSeek-R1 performed similarly to ChatGPT-o1 (accuracy of 0.57 vs
98 0.55; p=0.76 and 0.74 vs 0.76; p=0.06, using New England Journal of Medicine and
99 Medicilline databases, respectively). For RECIST classifications, DeepSeek-R1 also
100 performed similarly to ChatGPT-o1 (0.73 vs 0.81; p=0.10). Diagnostic reasoning steps
101 provided by DeepSeek were deemed more accurate than those provided by ChatGPT
102 and Llama 3.1-405B (average Likert score of 3.61, 3.22, and 3.13, respectively, p=0.005
103 and p<10⁻³). However, summarized imaging reports provided by DeepSeek-R1
104 exhibited lower global quality than those provided by ChatGPT-o1 (5-point Likert score:
105 4.5 vs 4.8; p<10⁻³). This study highlights the potential of DeepSeek-R1 LLM for medical
106 applications but also underlines areas needing improvements.

107

108

109

110

111

112

113

114 **Main:**

115 Large Language Models (LLMs) have emerged as transformative tools in a variety of
116 domains, exhibiting remarkable capabilities in natural language processing, knowledge
117 synthesis, and task automation¹. In healthcare, these models are being leveraged for
118 applications such as medical diagnosis, clinical decision support, and imaging
119 interpretation. Their ability to process vast amounts of medical knowledge and generate
120 contextually relevant insights has made them increasingly invaluable assets for
121 researchers and clinicians. Notable LLMs include proprietary models like ChatGPT
122 (OpenAI), open-access models such as Llama (Meta), and emerging models like
123 DeepSeek (DeepSeek Artificial Intelligence Co.).

124 DeepSeek-R1, a 671B parameter LLM released on January 20th 2025, represents a
125 new generation of open-source LLMs exhibiting competitive accuracy in third-party
126 evaluations^{2,3}. However, despite the broad applicability of all models, differences in their
127 training and optimization may affect performance for specific tasks. For example,
128 DeepSeek-R1 aims for better reasoning using reinforcement learning, which may be
129 optimally suited for medical interpretation. While established models like ChatGPT⁴ and
130 Llama⁵ have set performance benchmarks, DeepSeek-R1's release raises important
131 questions about tailoring comparisons within the context of specified applications, such
132 as medical interpretation tasks.

133 This study evaluated DeepSeek-R1 against two comparable leading LLMs, ChatGPT-o1
134 and Llama 3.1-405B (**Extended Data Table 1**), on four medical domain tasks to
135 benchmark their medical reasoning and processing capabilities. These four tasks
136 consisted of: **(1)** answering United States Medical Licensing Examination (USMLE)
137 questions, **(2)** clinical diagnosis and management based on cases from the New
138 England Journal of Medicine (NEJM), British Medical Journal (BMJ), and Medicalline
139 (not available online), **(3)** summarizing the "Findings" section of imaging reports into an
140 appropriate "Impression" based on cases from the MIMIC-III dataset⁶ and another
141 private dataset, and **(4)** performing RECIST 1.1 tumor treatment response classification
142 based on two consecutive reports (private dataset) from consecutive CT scans in
143 cancer patients (**Figure 1**). Datasets and prompts used across all tasks are summarized
144 in **Figure 1 and Extended Data Table 2**; the inclusion of private datasets was important
145 for addressing data leak. These tasks were selected to assess the models' abilities in
146 retrieving medical knowledge, reasoning on complex case analysis and management
147 options, providing imaging report classification, and summarization.

148 Additionally, we investigated intra-model variability by evaluating each model three
149 times with the same inputs with analysis of outputs utilizing BERTscore⁷ for text outputs
150 and Fleiss' kappa⁸ for categorial outputs. BERTscore assesses the semantic similarity
151 between different text outputs by measuring how closely the meaning aligns across

152 repeated model generations. Fleiss' kappa evaluates the agreement between repeated
153 categorical outputs across multiple evaluations, with values ranging from 0 (no
154 agreement) to 1 (perfect agreement).

155 **Results and Discussion**

156 Performance of the three LLMs on the **USMLE question dataset**⁹ (323 questions)
157 varied, with DeepSeek-R1 achieving an overall accuracy of 0.92 [95%CI: 0.90-0.95],
158 ChatGPT-o1 an accuracy of 0.95 [95%CI: 0.92-0.98], and Llama 3.1-405B an accuracy
159 of 0.83 [95%CI: 0.78-0.87] (**Figure 2**). DeepSeek-R1 achieved accuracy slightly inferior
160 to ChatGPT-o1 when combining all three Step examinations ($p=0.04$) and better than
161 Llama 3.1-405B ($p<10^{-3}$).

162 We evaluated the 3 models on **diagnostic and management tasks**. We first used the
163 **NEJM Case Challenges**¹⁰ (50 cases) to assess diagnostic capabilities through both
164 multiple-choice questioning, reasoning steps, and open-ended differential diagnosis
165 generation.

166 The pooled accuracy of DeepSeek-R1 for these text-based challenges was 0.57
167 [95%CI: 0.49-0.64] compared to ChatGPT-o1 with an accuracy of 0.55 [95%CI: 0.47-
168 0.62] ($p=0.76$) and Llama 3.1-405B with an accuracy of 0.47 [95%CI: 0.40-0.55]
169 ($p=0.16$) when provided multiple-choices.

170 When provided the same NEJM cases without answer choices, accuracy dropped to
171 0.36 [95%CI: 0.24-0.50] for DeepSeek-R1, 0.32 [95%CI: 0.21-0.46] for ChatGPT-o1
172 ($p=0.23$), and 0.26 [95%CI: 0.16-0.40] for Llama 3.1-405B ($p=0.77$). Average accuracy
173 on providing 3 differential diagnoses were 0.53 [95%CI: 0.39-0.66] for DeepSeek-R1,
174 0.57 [95%CI: 0.43-0.69] for ChatGPT-o1, and 0.49 [95%CI: 0.35-0.62] for Llama 3.1-
175 405B (**Extended Data Table 3**).

176 Additionally, a 5-point Likert scale (1: poor quality; 5: excellent) was used to evaluate
177 the models' clinical reasoning on NEJM cases, assessed by two physicians who are
178 specialized in internal medicine. Evaluation of the generated reasoning yielded scores
179 of 3.61 for DeepSeek-R1, 3.22 for ChatGPT-o1 ($p=0.005$), and 3.13 for Llama 3.1-405B
180 ($p<10^{-3}$).

181 Then, the **BMJ Endgames**¹¹ were used to evaluate the models' clinical management
182 capabilities. These were also quantified with 5-point Likert scores by two physicians,
183 with scores of 4.58 for DeepSeek-R1, compared to 4.28 for ChatGPT-o1 ($p<10^{-3}$), and
184 4.14 for Llama 3.1-405B ($p<10^{-3}$).

185 When evaluated on the **Medicilline database**¹² of medical cases with multiple-choice
186 questions (with multiple possible correct answers), DeepSeek-R1 achieved an accuracy
187 of 0.74 [95%CI: 0.68-0.80], comparable to ChatGPT-o1 with an accuracy of 0.76

188 [95%CI: 0.68-0.80] ($p=0.06$) and better than Llama 3.1-405B with an accuracy of 0.66
189 [95%CI: 0.59-0.72] ($p=0.01$).

190 Four independent radiologists (two for the MIMIC dataset and two for the private
191 dataset) performed clinical and linguistic evaluation of the radiological impressions
192 generated by models from the “Findings” section of radiological reports using a 5-points
193 Likert score. Average scores were 4.6 for DeepSeek-R1, 4.8 for ChatGPT-o1, and 4.5
194 for Llama 3.1-405B on a private dataset (100 reports) and 4.3, 4.8, and 4.6,
195 respectively, on the MIMIC-III dataset (100 reports)⁶ (**Extended Data Table 4**). The two
196 evaluators consistently gave lower scores to Llama 3.1-405B and DeepSeek-R1
197 compared to ChatGPT-o1 ($p<10^{-3}$). Additionally, the mean evaluation scores of the
198 management recommendation of Deepseek-R1 and Llama 3.1-405B on both datasets
199 were significantly lower than ChatGPT-o1 (3.8 and 3.9 vs 4.6 respectively, $p<10^{-3}$ for the
200 2 comparisons), with subsequent lower scores for harmlessness. The number of words
201 per output for the MIMIC (average initial report: 187 words) and private reports (average
202 initial report: 222 words) were 128 and 138 for DeepSeek-R1, 116 and 131 for
203 ChatGPT-o1, and 151 and 169 for Llama 3.1-405B, demonstrating limited conciseness
204 for DeepSeek-R1 and Llama 3.1-405B.

205 Accuracy of the models for **RECIST 1.1 tumor treatment response classification**
206 (compared to consensus evaluation by human experts as the gold standard) **on two**
207 **consecutive imaging reports** of 100 cancer patients was 0.73 [95%CI: 0.63-0.80] for
208 DeepSeek-R1, compared to 0.72 [95%CI: 0.63-0.80] for Llama 3.1-405B ($p=0.83$) and
209 0.81 [95%CI: 0.72-0.88] for ChatGPT-o1 ($p=0.10$) (**Extended Data Table 5**).

210 We evaluated model output reproducibility by performing 3 runs of each task for all
211 models. On the USMLE dataset, Fleiss’ kappa was 0.96 for DeepSeek-R1, 0.98 for
212 ChatGPT-o1, and 0.95 for Llama 3.1-405B. On the Medicilline dataset, Fleiss’ kappa
213 was 0.61 for DeepSeek-R1, 0.93 for ChatGPT-o1, and 0.87 for Llama 3.1-405B. On the
214 RECIST dataset, Fleiss’ kappa was 0.80 for DeepSeek-R1, compared to 0.87 for
215 ChatGPT-o1 and 0.96 for Llama 3.1-405B. BERTScore F1-scores ranged from 0.82-
216 0.93 for DeepSeek, 0.84-0.93 for ChatGPT-o1, and 0.93-0.97 for Llama 3.1-405B on
217 NEJM case reasoning, BMJ management reasoning, and report summarization tasks
218 with text outputs (**Extended Data Table 6** and **Extended Data Table7**).

219 Inter-model kappa ranged from 0.52-0.77 between DeepSeek-R1 and Llama 3.1-405B,
220 0.60-0.91 between Deepseek-R1 and ChatGPT-o1, and 0.44-0.77 between ChatGPT-
221 o1 and Llama 3.1-405B on USMLE, Medicilline multiple-choice questions, NEJM
222 multiple-choice, and RECIST classification (**Extended Data Table 8**).

223 Our results offer initial insights into the capabilities and limitations of DeepSeek-R1
224 compared to two major LLMs, ChatGPT-o1 and Llama 3.1-405B, across a range of
225 medical tasks.

226 Our evaluation demonstrates notable variability among LLMs when assessed on the
227 USMLE dataset. ChatGPT-o1 exhibited superior performance with the highest overall
228 accuracy (0.95), significantly outperforming both DeepSeek-R1 (0.92; p=0.04) and
229 Llama 3.1-405B (0.83; p<10⁻³). While DeepSeek-R1 was closely comparable to
230 ChatGPT-o1, its slightly lower accuracy suggests potential room for improvement in
231 medical knowledge recall and application. In contrast, Llama 3.1-405B demonstrated
232 markedly inferior accuracy, highlighting considerable differences in clinical knowledge
233 capabilities across models.

234 The models showed mixed diagnostic and clinical reasoning capabilities across multiple
235 medical evaluation tasks. On diagnostic assessments using NEJM Case Challenges,
236 DeepSeek-R1 performed comparably to ChatGPT-o1 in both multiple-choice (0.57 vs.
237 0.55, p=0.76) and open-ended scenarios (0.36 vs. 0.32, p=0.23), but both clearly
238 struggled without answer options, underscoring persistent challenges in open-ended
239 differential diagnosis tasks. Despite the advantage of DeepSeek-R1 in reasoning scores
240 (Likert 3.61 vs. ChatGPT-o1's 3.22; p=0.005), clinical management evaluations through
241 the BMJ Endgames further highlighted DeepSeek-R1's strength (Likert 4.58),
242 significantly outperforming both ChatGPT-o1 and Llama 3.1-405B. Additionally, both
243 DeepSeek-R1 (accuracy 0.74) and ChatGPT-o1 (accuracy 0.76; p=0.06) demonstrated
244 strong and comparable performance on Medicilline's multiple-choice questions,
245 highlighting their suitability for tasks requiring general diagnostic and management
246 knowledge.

247 Radiological impression generation assessments revealed consistent superiority of
248 ChatGPT-o1 in terms of clinical and linguistic quality (average Likert scores 4.8 vs. 4.3-
249 4.6; p<10⁻³). ChatGPT-o1 also excelled significantly in formulating clinical management
250 recommendations (average 4.6 vs. 3.8-3.9 for DeepSeek-R1 and Llama 3.1-405B;
251 p<10⁻³). Notably, ChatGPT-o1 demonstrated superior conciseness, generating shorter
252 but clearer outputs compared to DeepSeek-R1 and Llama 3.1-405B.

253 In tumor response classification tasks based on imaging reports (RECIST 1.1),
254 ChatGPT-o1 again demonstrated superior accuracy (0.81), though not significantly
255 surpassing DeepSeek-R1 (0.73, p=0.10) and Llama 3.1-405B (0.72, p=0.83). This
256 suggests a moderate performance ceiling for current models in complex oncologic
257 response evaluation tasks, emphasizing a critical area for targeted refinement.

258 Our reproducibility analysis underscores substantial differences in model stability across
259 tasks. While all models demonstrated high reproducibility on the USMLE dataset,
260 significant variability emerged in the Medicilline and RECIST datasets, particularly for
261 DeepSeek-R1 (Fleiss' kappa 0.61 and 0.80, respectively), compared to ChatGPT-o1
262 (0.93 and 0.87) and Llama 3.1-405B (0.87 and 0.96). These differences highlight the
263 necessity of rigorous reproducibility testing when deploying LLMs clinically. Inter-model

264 agreement analyses (Fleiss' kappa ranging from 0.44 to 0.91) revealed varying degrees
265 of consensus, indicating that despite similarities, models leverage distinct internal
266 reasoning pathways. These results illustrate the potential for leveraging complementary
267 strengths through ensemble approaches or hybrid model designs to improve clinical
268 reliability.

269 This study has several limitations. First, the use of publicly available datasets raises the
270 possibility of inclusion in the original training data for models. The oldest diagnostic
271 challenge used in this study was published in 2016. In order to balance the needs of
272 reproducibility of LLM research using publicly available databases while mitigating data
273 leak in training sets, we included one private dataset of medical cases with multiple-
274 choice questions (Medicilline, not available online), and two private datasets of
275 radiological reports for the summarization and classification tasks. Additionally, datasets
276 included the most recent cases to further decrease the risk of data leak.

277 Second, it is important to acknowledge that no imaging tasks were included for
278 consistency of comparison of these non-multimodal models. DeepSeek-R1 lacks
279 support for image-based tasks, while DeepSeek-Janus, an older and smaller model,
280 supports image interpretation. Similarly, ChatGPT-o1 and Llama 3.1-405B don't allow
281 image interpretation. Future studies could incorporate multimodal models with direct
282 image analysis capabilities once they become available.

283 Third, our analysis was restricted to Azure API-based models for consistency; thus,
284 performance might differ from the same LLMs when accessed through their native
285 applications or other deployment environments. Additionally, these platforms do not
286 provide the exact parameters used, and the online applications may use data shared for
287 their future training, raising questions on their use in research settings.

288 Fourth, reasoning evaluation of NEJM cases, management evaluation of BMJ cases,
289 and quality evaluation of generated Impressions are subject to human variability. All the
290 tasks were performed independently by blinded evaluators (the evaluators were not
291 aware of which output corresponded to which LLM).

292 Fifth, the constant and rapid evolution of the LLMs may limit the validity of the current
293 results when updated versions of these models are released.

294 In conclusion, DeepSeek-R1 models exhibit both strengths and limitations compared to
295 ChatGPT-o1 and Llama 3.1-405B in general medical knowledge and clinical reasoning.
296 Its open-source nature offers a solution for structured medical tasks. Its reasoning
297 capabilities are of great interest, but summarization and repeatability issues exist.
298 Future studies may focus on enhancing clinical reasoning and open-ended diagnostic
299 capabilities of large language models, especially in complex clinical scenarios where

300 current models exhibit limitations. Investigating methods to increase reproducibility,
301 consistency, and leveraging complementary strengths would also be beneficial.

302 **Acknowledgments:**

303 We thank Guillaume Zagury and Médicilline editions for sharing their dataset. We thank
304 Dr. Venturelli, Dr. Chiche, Dr. Beaziz, and Dr. Lejoyeux for their help updating these
305 questions based on the latest medical knowledge. This project is supported by the Eric
306 and Wendy Schmidt AI in Human Health Fellowship, a program of Schmidt Sciences.
307 M.T. is supported by the French Society of Radiology and the French Musculoskeletal
308 Imaging Society. Figure 1 was created in BioRender.

309

310 **Author contribution Statement:**

311 Concept and design: MT, ZL, HCL, ZAF, XM.
312 Acquisition, analysis or interpretation of data: MT, ZL, MY, VF, YM, JH, IB, HA, CH, ASP,
313 AG, AM, NY, NN, PR, AZ, SL, MH, TD, BT, HCL, ZAF, XM.
314 Drafting of the paper: MT, ZL, MY, VF, AM, IB, SL, BT, HCL, ZAF, XM.
315 Critical revision of the paper and final draft: All authors.

316

317 **Competing Interest Statement:**

318 T.D. is the managing partner of RadImageNet LLC and a paid consultant to GEHC and
319 AirsMedical. X.M. is a paid consultant to RadImageNet LLC. The other authors declare
320 no competing interests.

321

322 **Figure Captions:**

323 **Figure 1:** Overview of this study comparing three Large Language Models in performing
324 the four indicated medical tasks. Abbreviations: pt: point, PD: Progressive Disease, PR:
325 Partial Response, N: Number.

326

327 **Figure 2:** Comparative performance of three LLMs for a variety of medical tasks. The
328 performance of the three LLMs (ChatGPT-o1, DeepSeek-R1 and Llama 3.1-405B) was
329 evaluated in the following medical tasks: United States Medical Licensing Examination
330 (USMLE) questions, diagnostic questions (in the NEJM and Medicalline databases),
331 reasoning (in the BMJ and NEJM databases), RECIST 1.1 classification from
332 radiological reports, and report summarization (in MIMIC-III and private datasets). For
333 USMLE questions, diagnostic challenges, and RECIST 1.1 classification, bar plots for
334 each model represent the mean of the average output score from each model over
335 three runs for each model, with error bars indicating the standard deviation of the
336 average output scores across those runs. The average output score is the accuracy of
337 the model for each respective task as determined from the average outcome of the task
338 over the number of patient datasets included for each task, respectively. In the USMLE
339 questions, Step 1 assesses basic science, Step 2 evaluates clinical knowledge, and
340 Step 3 tests patient management skills. For RECIST 1.1 classification, evaluations are
341 shown for individual classes: complete response (CR), partial response (PR), stable
342 disease (SD), and progressive disease (PD), as well as overall model performance. For
343 the reasoning evaluation using the BMJ and NEJM databases, the bar plots for each
344 model indicate the mean of the average Likert score from two human reviewers, and the
345 error bars represent the standard deviation between their scores. The average Likert
346 score for one human reader is the average rating of the accuracy or appropriateness of
347 the LLM-based output by the human reader over all patient datasets used in the
348 respective task. For report summarization, the radar charts present model performance
349 on report summarization in the MIMIC-III and private datasets across eight criteria,
350 including scientific terminology, coherence, lack of bias, harmlessness,
351 comprehensiveness, specific diagnosis, differential diagnosis, and management
352 recommendations. The numbers shown on the radar charts (e.g., 0, 3.5, and 5.0)
353 correspond to values on a 5 points Likert score ranging from 0 to 5. The average 5-point
354 Likert score over all patient datasets evaluated for each task and for both human
355 readers is shown on respective spokes.

356

357 **References:**

- 358 1. The Lancet Digital Health, null. Large language models: a new chapter in digital
359 health. *Lancet Digit. Health* **6**, e1 (2024).
- 360 2. Gibney, E. Scientists flock to DeepSeek: how they're using the blockbuster AI model.
361 *Nature* (2025) doi:10.1038/d41586-025-00275-0.
- 362 3. Conroy, G. & Mallapaty, S. How China created AI model DeepSeek and shocked the
363 world. *Nature* (2025) doi:10.1038/d41586-025-00259-0.
- 364 4. OpenAI. GPT-4 Technical Report. Preprint at <http://arxiv.org/abs/2303.08774> (2023).
- 365 5. Grattafiori, A. *et al.* The Llama 3 Herd of Models. Preprint at
366 <https://doi.org/10.48550/arXiv.2407.21783> (2024).
- 367 6. Johnson, A. E. W. *et al.* MIMIC-III, a freely accessible critical care database. *Sci.*
368 *Data* **3**, 160035 (2016).
- 369 7. Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q. & Artzi, Y. BERTScore: Evaluating
370 Text Generation with BERT. Preprint at <https://doi.org/10.48550/arXiv.1904.09675>
371 (2020).
- 372 8. Fleiss, J. L. Measuring nominal scale agreement among many raters. *Psychol. Bull.*
373 **76**, 378–382 (1971).
- 374 9. United States Medical Licensing Examination dataset. [https://www.usmle.org/exam-](https://www.usmle.org/exam-resources)
375 [resources](#).
- 376 10. The New England Journal of Medicine. Case challenges. [https://www.nejm.org/case-](https://www.nejm.org/case-challenges)
377 [challenges](#).
- 378 11. BMJ Endgames. <https://www.bmj.com/education/endgames>.
- 379 12. Medicilline dataset. <http://www.medicilline.com/13-externe-ecni>.

- 380 13. Suh, P. S. *et al.* Comparing Large Language Model and Human Reader Accuracy
381 with New England Journal of Medicine Image Challenge Case Image Inputs.
382 *Radiology* **313**, e241668 (2024).
- 383 14. Zhang, L. *et al.* Constructing a Large Language Model to Generate Impressions
384 from Findings in Radiology Reports. *Radiology* **312**, e240885 (2024).
- 385 15. Laurent, G., Craynest, F., Thobois, M. & Hajjaji, N. Automatic Classification of Tumor
386 Response From Radiology Reports With Rule-Based Natural Language Processing
387 Integrated Into the Clinical Oncology Workflow. *JCO Clin. Cancer Inform.* **7**,
388 e2200139 (2023).
- 389 16. Penny, P., Bane, R. & Riddle, V. Advancements in AI Medical Education: Assessing
390 ChatGPT's Performance on USMLE-Style Questions Across Topics and Difficulty
391 Levels. *Cureus* **16**, e76309 (2024).

392
393
394
395

396 **Online Methods**

397 This retrospective study, using anonymized diagnostic imaging reports, was approved
398 by the local institutional review board of the Icahn School of Medicine at Mount Sinai
399 (IRB protocol number GCO# 20-2199). The requirement to obtain individual informed
400 consent was waived.

401 ***Large Language Models***

402 LLMs can process and generate text that mirrors human-like conversation. This task is
403 achieved by predicting what word (token) should come next given an input, allowing
404 them to generate coherent, contextually relevant responses for diverse tasks, such as
405 question answering, summarization, and report generation in the medical field¹⁷. In this
406 study, DeepSeek-R1 was evaluated in comparison with two major LLMs, ChatGPT-o1
407 (proprietary model) and Llama 3.1-405B (open-access model), on various medical
408 tasks. All models were implemented using standardized computational resources on an
409 Azure Cloud infrastructure (Azure OpenAI Service, Azure AI Service, Azure AI Foundry),
410 following each LLM provider's guidelines and instructions for Azure deployment. By
411 leveraging Azure's consistent computational environment, we ensured uniform
412 experimental conditions, enabling reproducibility and rigorous comparative evaluation.
413 All LLMs evaluated in our study are instruct-version models.

414 ChatGPT-o1¹⁸, developed by OpenAI, and released in September 2024, is optimized for
415 enhanced reasoning capability. It demonstrates strong performance across various
416 domains, including the medical field¹⁹.

417 LLama3.1⁵, developed by Meta, is a powerful open-source LLM. In this study, the Llama
418 model was applied using its default parameters with temperature of 1.0 and top_p of
419 1.0.

420 Deepseek-R1 has been developed to use large-scale post-training reinforcement
421 learning to improve the reasoning of LLMs and has demonstrated superior performance
422 against GPT-o1 in various benchmark datasets²⁰.

423 ***Study Design***

424 This comparative study evaluated the performance of these three LLMs—DeepSeek-
425 R1, ChatGPT-o1 and Llama 3.1-405B—on four distinct medical tasks. The models are
426 summarized in **Extended Data Table 1**. Each of the different tasks was chosen to
427 determine the capability of these models in different contexts: i) to answer medical
428 multiple-choice questions, ii) to produce text-based diagnostic reasoning and
429 management options in complex clinical case scenarios, iii) to perform classification,
430 and iv) to perform summarization. No image interpretation was evaluated since these
431 models are not multimodal. Additionally, to investigate the robustness of the models we

432 assessed the test-retest repeatability of the outputs of each model for each task from
433 three separate queries with the same prompts. All of the prompts used are available in
434 **Extended Data Table 2**.

435 Task 1: USMLE Questions

436 Three hundred and twenty-three USMLE-style questions were used for evaluation.
437 These included questions from Step 1 (95 questions), Step 2 (106 questions), and Step
438 3 (122 questions) of the examination, comprising the entire publicly available database
439 (<https://www.usmle.org/exam-resources>), excluding the questions involving image
440 interpretation. Step 1 questions focus on foundational sciences (e.g. biochemistry,
441 pathology), Step 2 questions emphasize clinical knowledge, and Step 3 questions
442 require reasoning for patient management.

443 Each LLM was prompted with identical questions using standardized formats, including
444 the question stem and answer options. The models provided their final answers without
445 additional clarifications or follow-ups. The accuracy of each model was the proportion of
446 correct responses.

447 Task 2: Diagnostic and Management Case Challenges

448 Fifty NEJM “Case Challenges” (the entire online available database) were selected to
449 evaluate the models’ interpretive skills, in a similar manner to previous studies^{13,21}.
450 These challenges included complex clinical scenarios with text and tables of biological
451 values. These cases typically require integration of patient history, laboratory findings,
452 description of imaging examinations, and clinical reasoning to arrive at a diagnosis.

453 The models were tested in two ways. Firstly, the LLM was prompted to select the
454 correct diagnosis from a list of multiple-choices provided in the prompt. Data provided
455 in the prompt included patient history, symptoms, and all available data provided in the
456 Case Challenge. The model was also prompted to provide the reasoning associated
457 with the answer. Secondly, the LLM was prompted to provide the most probable
458 diagnosis and two differential diagnoses based in an open-ended answer format, based
459 on the same information (but not including the multiple-choice answers provided in the
460 first test). For each test, the accuracy of each model was the proportion of correct
461 answers. For the reasoning portion of each test, the responses of the LLMs were
462 evaluated on their ability to use inductive, deductive, and probabilistic reasoning to
463 progressively narrow down diagnosis possibilities based on the available clinical history,
464 physical exam observations, and test results. Independently, two Internal Medicine
465 specialists evaluated the pertinency of the elements evaluated, their interpretation, and
466 how the combination of the clinical examination and different tests were handled to
467 achieve the diagnosis on a 5-point Likert scale (1: poor quality; 5: excellent). We
468 adapted the score from the R-IDEA score²². However, we did not use it in this case

469 because we specifically asked the models to provide differential diagnoses. The expert
470 reviewers also evaluated the accuracy of the LLM-derived differential diagnoses based
471 on the number of plausible differential diagnoses provided in the open-ended questions:
472 one point was awarded for each plausible diagnosis up to a maximum of 3 (0-3/3).

473 Fifty BMJ “Endgames” which included clinical management questions were also used to
474 evaluate the management capabilities of the models. These cases usually evaluate the
475 image interpretation, diagnostic skills, and management capabilities of the readers with
476 open-ended questions. We provided each model with the text prompts and final
477 diagnosis (as image analysis is needed to reach final diagnosis) and prompted
478 questions about the management of these conditions (as provided in the original cases).
479 The clinical management plans proposed by the models were evaluated by two
480 physicians using a 5-point Likert score based on the different therapeutic options
481 proposed, the harmfulness of the treatment, and the essential measures associated
482 with medical treatment (1: poor management; 5: excellent management). The
483 management suggestions provided in the BMJ cases were used as the reference
484 standard for adjudication by the physicians.

485 Finally, 200 questions evaluating diagnostic and management capabilities were
486 extracted from a private French database (ECN-Integrale, from Mediciline Editions).
487 This series of textbooks contain training material for medical students and residents but
488 are not available online, preventing leakage into the training datasets of the LLMs..
489 These multiple-choice questions allowed multiple correct answers. Correct answers
490 were rated Questions were translated into English, then updated and verified by experts
491 in the different fields, including one who is involved in student training who has more
492 than 10 years of experience. The LLMs were provided with the questions in text format
493 and prompted to select the correct answers. Correct answers were rated scored as 1,
494 answers that contained one mistake among the multiple-choices was scored as 0.5, and
495 answers that contained 2 or more mistakes was scored as 0. The performance of each
496 model was the sum of the scores divided by the maximum possible total scores.

497 Task 3: Findings to Impression

498 The “Findings” section of 200 imaging reports, 100 from an outpatient radiology facility
499 in the New York Metropolitan area and 100 from MIMIC⁶ reports (each with 20 examples
500 of X-rays, ultrasound, CT-scan, MRI and PET/CT), were provided to the 3 LLMs for
501 summarization. Four readers (radiologists with 3, 4, 6 and 7 years of experience) each
502 evaluated 100 Impressions generated by these LLMs (2 radiologists evaluated the
503 MIMIC dataset and 2 the local dataset). The readers were blinded to which model
504 provided the outputs. Nine domains were evaluated with a 5-point Likert score based on
505 previous studies evaluating similar tasks^{14,23}. The average number of words and the
506 ratio of the number of words in the LLM-derived Impression and the number of words in

507 the original Findings Impression/Findings were also evaluated for each model to assess
508 the summarization capabilities.

509 **Task 4: RECIST evaluation**

510 Reports from two consecutive oncological chest-abdomen-pelvis CT-scans were
511 provided to each model which was subsequently prompted to classify the progression
512 of the disease according to RECIST 1.1²⁴. No baseline or NADIR were available in
513 these reports. The models were prompted to provide a text summary. Each model
514 evaluated data from 100 patients from an outpatient radiology facility in New York. Two
515 radiologists with 3 and 5 years of experience in cancer imaging independently evaluated
516 the classification. An expert radiologist with more than 10 years of experience resolved
517 any discrepancy. All evaluators were blinded to the LLM providing the output. The
518 accuracy of the LLM-derived summary compared to the Radiologists' interpretation was
519 assessed.

520 **Test-Retest Reproducibility**

521 To evaluate the reproducibility of each model, we performed a test-retest analysis by
522 prompting each model 3 times with the same prompts for each task. Depending on the
523 tasks, the Fleiss kappa score was used for categorical/multiple-choice questions (like
524 USMLE, NEJM accuracy, Medicilline, and RECIST), the Cohen kappa score was used
525 for comparison between two runs of 2 models, and the BERTscore was used for text
526 outputs (reasoning outputs, management options, and Impression summarizations).

527 **Statistical Analysis**

528 SPSS v28.0 was used for statistical analysis. Descriptive statistics were used to
529 summarize accuracy metrics or linguistic evaluation (for Task 3) for each model across
530 tasks. Confidence intervals were calculated using the Wilson score method. McNemar
531 tests were performed for the comparisons between 2 models for paired-sample analysis
532 and Mann-Whitney U test for independent samples. A p-value less than 0.05 was
533 considered significant. Variability between the 3 runs of each model was analyzed using
534 Fleiss's kappa score for multiple-choice questions, Cohen's kappa score for comparison
535 between 2 runs/models,) for multiple-choice questions and the BERTscore for text
536 outputs.

537 **Data Availability**

538 The USMLE dataset is available at <https://www.usmle.org/exam-resources>. The NEJM
539 Case Challenges are available at <https://www.nejm.org/case-challenges>. The BMJ
540 Endgames are available at <https://www.bmj.com/specialties/endgames>. The MIMIC-III
541 dataset can be requested at <https://physionet.org/content/mimiciii/1.4/>. The Médicilline
542 dataset (<http://www.medicilline.com/>) of multiple-choice questions for clinical diagnosis

543 and management translated in English is available upon request for private/research
544 use only, after agreement of Medicilline. Mediciline data requests should be sent to
545 Mickael.Tordjman@mssm.edu. In response to the inquiry, the timeframe for responding
546 to requests is approximately within 2 weeks of the request. The radiological reports
547 used for summarization and RECIST classification are not available due to privacy
548 issues.

549 **Code Availability**

550 All studies were conducted using Azure OpenAI service (ChatGPT), Azure AI Foundry
551 (Llama), and Azure AI service (DeepSeek). The algorithms used in this study are based
552 on the official guidelines provided by the developers of the Large Language Models
553 evaluated—DeepSeek, ChatGPT, and Llama. The implementation followed these
554 guidelines strictly to ensure consistency and reproducibility of results. As the study did
555 not involve the development of new code but rather the application of existing, officially
556 provided algorithms, specific source codes referenced are proprietary and maintained
557 by their respective developers.

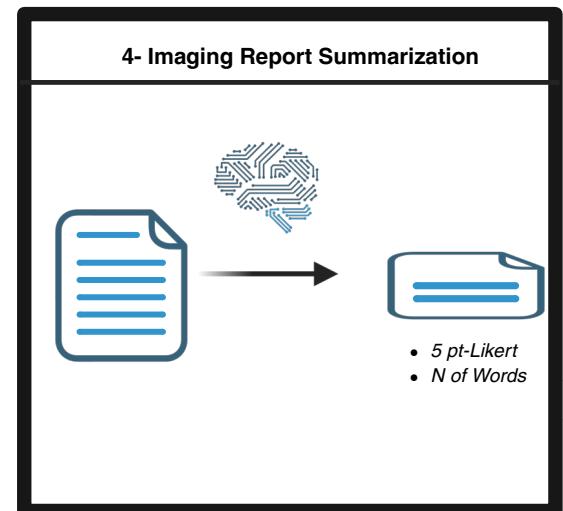
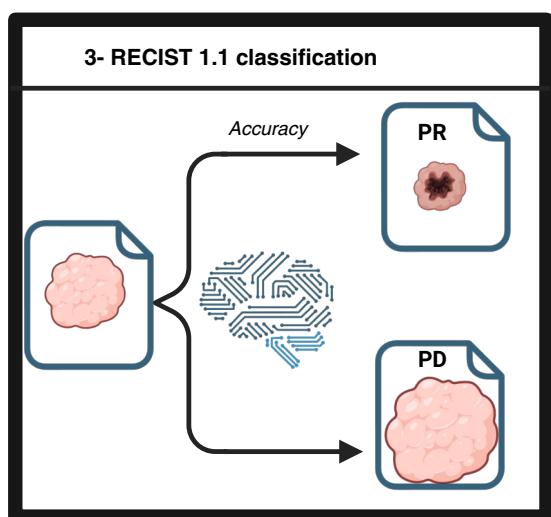
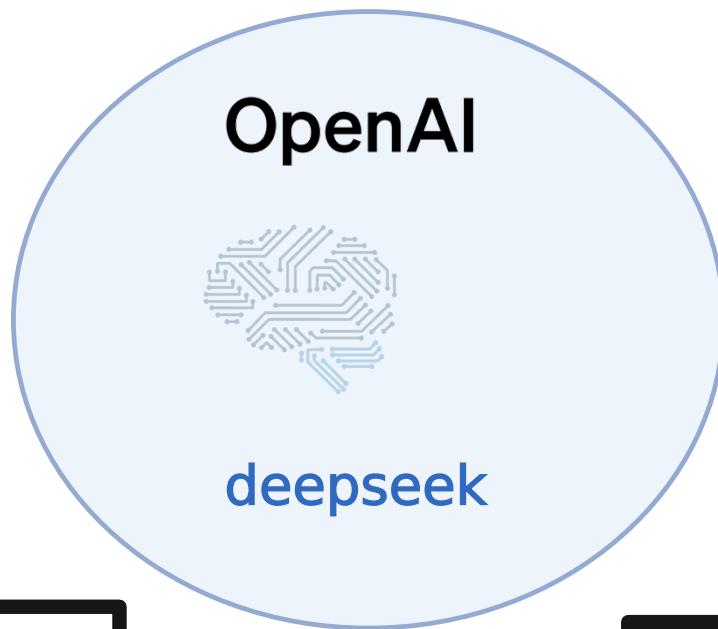
558 **References:**

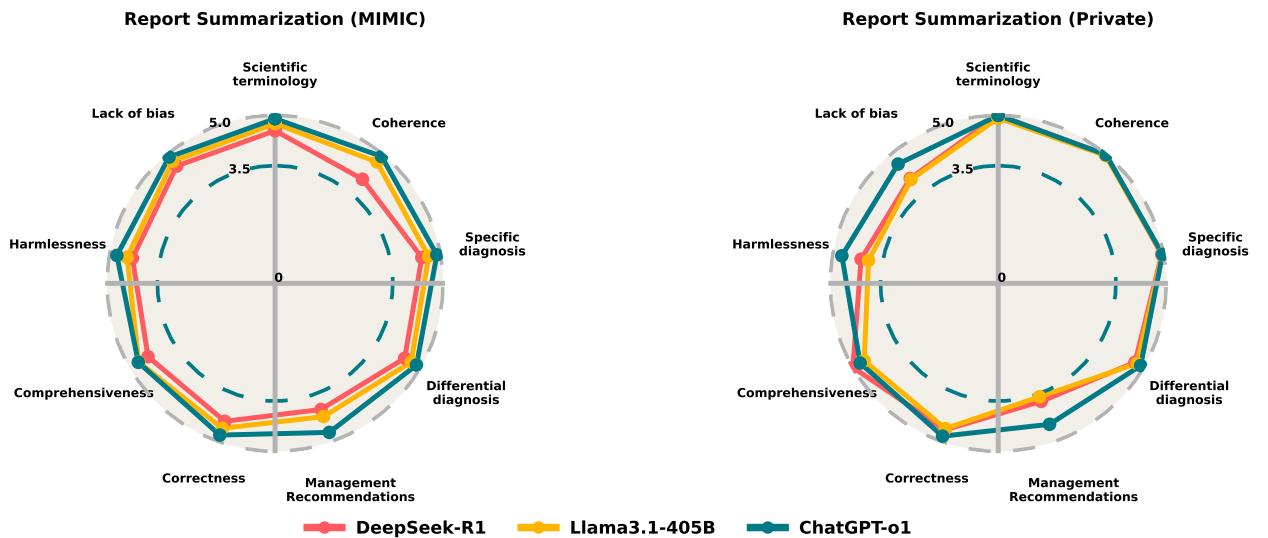
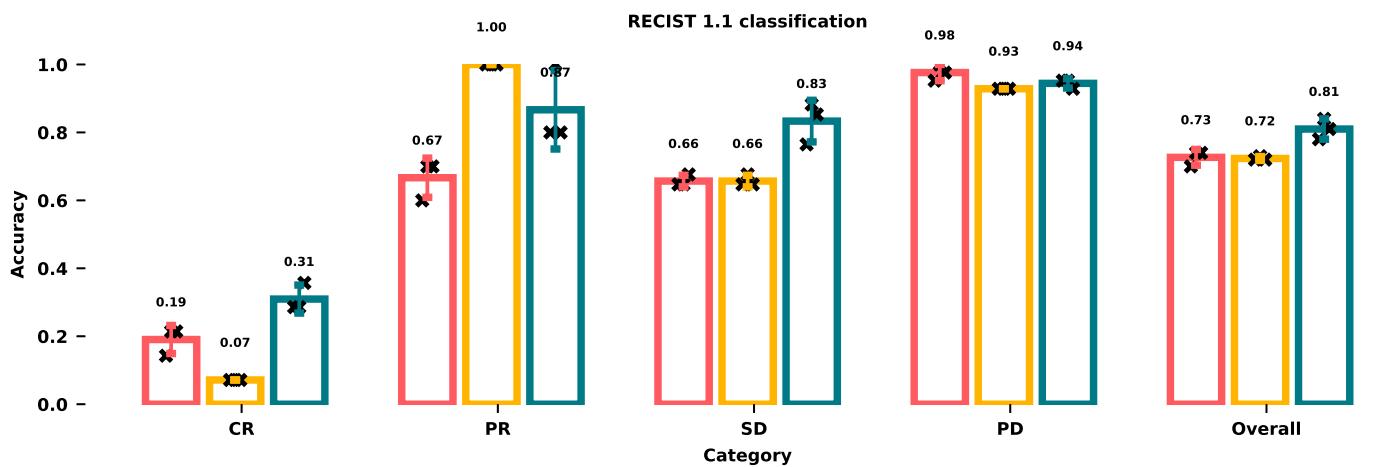
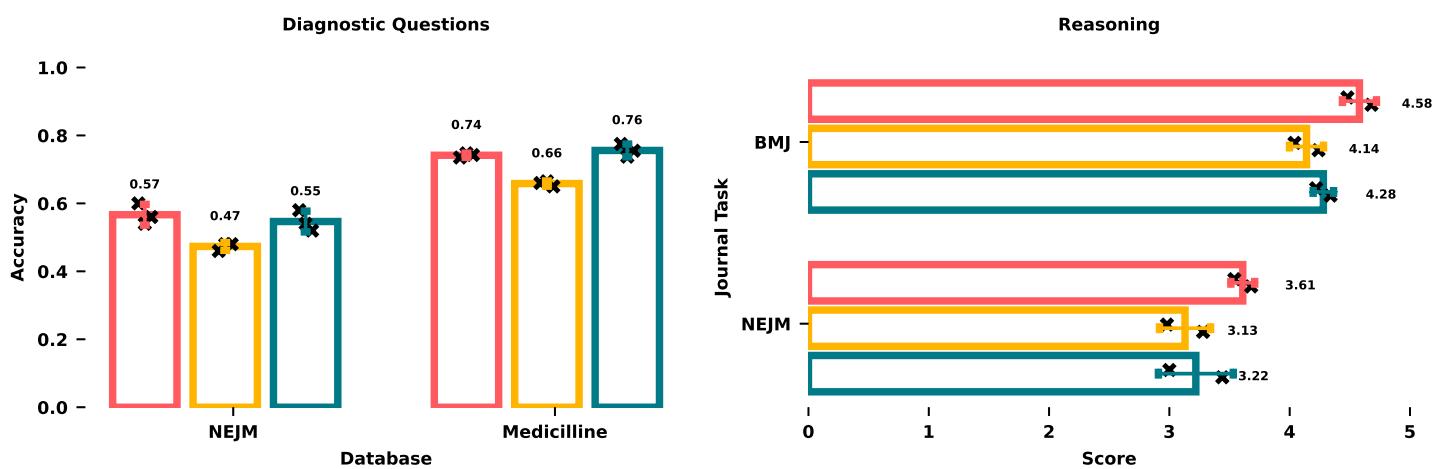
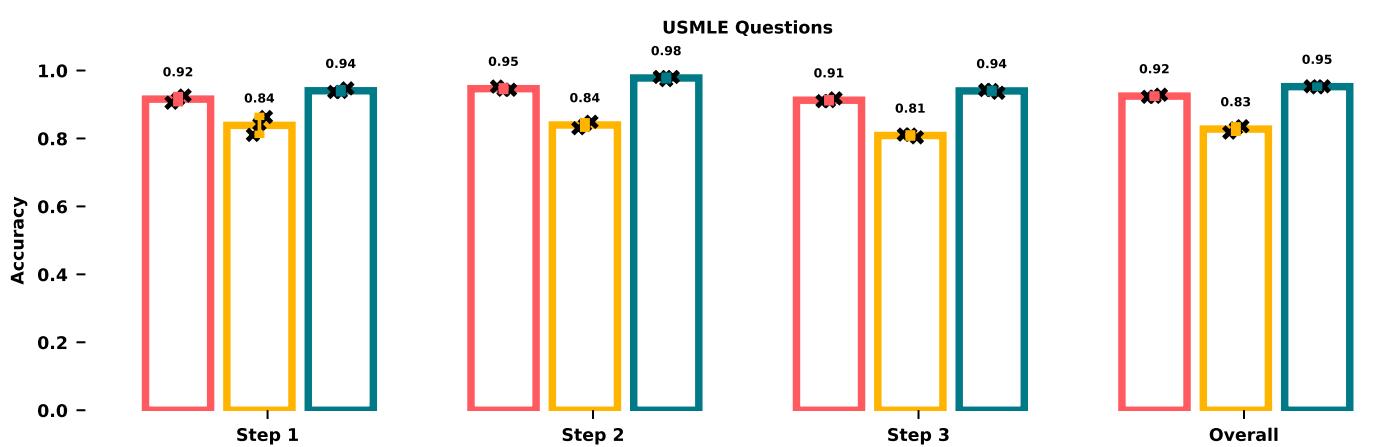
- 559 17. Zhao, W. X. *et al.* A Survey of Large Language Models. Preprint at
560 <https://doi.org/10.48550/arXiv.2303.18223> (2024).
- 561 18. Learning to reason with LLMs. <https://openai.com/index/learning-to-reason-with-llms/>.
- 563 19. Temsah, M.-H., Jamal, A., Alhasan, K., Temsah, A. A. & Malki, K. H. OpenAI o1-
564 Preview vs. ChatGPT in Healthcare: A New Frontier in Medical AI Reasoning.
565 *Cureus* **16**, e70640 (2024).
- 566 20. DeepSeek-AI *et al.* DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via
567 Reinforcement Learning. Preprint at <https://doi.org/10.48550/arXiv.2501.12948>
568 (2025).
- 569 21. Tenner, Z. M., Cottone, M. C. & Chavez, M. R. Harnessing the open access version
570 of ChatGPT for enhanced clinical opinions. *PLOS Digit. Health* **3**, e0000355 (2024).

- 571 22. Cabral, S. *et al.* Clinical Reasoning of a Generative Artificial Intelligence Model
- 572 Compared With Physicians. *JAMA Intern. Med.* **184**, 581 (2024).
- 573 23. Park, J., Oh, K., Han, K. & Lee, Y. H. Patient-centered radiology reports with
- 574 generative artificial intelligence: adding value to radiology reporting. *Sci. Rep.* **14**,
- 575 13218 (2024).
- 576 24. Eisenhauer, E. A. *et al.* New response evaluation criteria in solid tumours: revised
- 577 RECIST guideline (version 1.1). *Eur. J. Cancer Oxf. Engl.* **1990** *45*, 228–247 (2009).
- 578

1- USMLE questions	
USMLE	Accuracy

2- Text-based cases		
The New England Journal of Medicine	<ul style="list-style-type: none"> • Accuracy • Differential diagnoses • Reasoning (5 pt-Likert) 	
BMJ	<i>Management (5 pt-Likert)</i>	
	Accuracy	





Models	Parameters	Context length
ChatGPT-o1	N/A	128K
Llama3.1-405B	405B	128K
Deepseek-R1	671B	128K

Tasks	Prompts
USMLE/Medicilline	[Question] + [Options]
Text-based cases and Reasoning	[Text description] + [Options]+ "Please provide the reasoning steps"
Differential	"Please provide the most likely diagnosis and the 2 additional most likely differential diagnoses."
Management	" You are a physician. Below is a medical case: What is the management of this condition?"
Findings to impression summarization	""You are a radiologist, summarize the findings in this report into an impression containing the clinically significant findings and recommend any following up radiologic testing, if needed"" + [Findings]
RECIST classifications	<p>"""As a radiologist, utilize the Response Evaluation Criteria in Solid Tumors (RECIST) 1.1 to evaluate cancer treatment responses. The RECIST guidelines classify tumor responses into four categories:</p> <p>Complete Response (CR): The tumor is no longer detectable.</p> <p>Partial Response (PR): The tumor size has reduced by at least 30%.</p> <p>Stable Disease (SD): The tumor size remains unchanged, or does not meet the criteria for PR or PD.</p> <p>Progressive Disease (PD): The tumor size has increased by at least 20% or at least one new lesion appeared.</p> <p>Your task is to analyze radiology reports from two separate visits for the same patient and classify the tumor's response according to the RECIST criteria. Only respond with your classification result, choose one from CR, PR, SD or PD.</p> <p>Radiology Report 1: [Report1]</p> <p>Radiology Report 2: [Report2] """</p>

	DeepSeek-R1	Llama 3.1-405B	ChatGPT-o1
BMJ Management	4.58	4.14	4.28
NEJM Reasoning	3.53	3.14	3.11
NEJM Accuracy (MCQ)	0.57 [CI: 0.49-0.64]	0.47 [CI: 0.40-0.55]	0.55 [CI: 0.47-0.62]
NEJM Accuracy (Open-ended)	0.36 [CI: 0.24-0.50]	0.26 [CI: 0.16-0.40]	0.32 [CI: 0.21-0.46]
NEJM Differential	0.53 [CI: 0.39-0.66]	0.49 [CI: 0.35-0.62]	0.57 [CI: 0.43-0.69]
Medicilline Database	0.74 [CI: 0.68-0.80]	0.66 [CI: 0.59-0.72]	0.76 [CI: 0.68-0.80]

	Deepseek R1		Llama 3.1-405B405B		ChatGPT-o1	
	<i>Private</i>	<i>MIMIC</i>	<i>Private</i>	<i>MIMIC</i>	<i>Private</i>	<i>MIMIC</i>
Scientific terminology	5.0	4.5	4.9	4.8	5.0	4.9
Coherence	5.0	4.0	5.0	4.7	5.0	4.9
Specific diagnosis	4.9	4.4	5.0	4.6	5.0	4.9
Differential diagnosis	4.7	4.5	4.8	4.7	4.9	4.9
Management Recommendations	3.7	4.0	3.6	4.2	4.5	4.7
Correctness	4.6	4.4	4.6	4.6	4.8	4.8
Comprehensiveness	5.0	4.4	4.6	4.7	4.7	4.7
Harmlessness	4.1	4.3	3.9	4.5	4.7	4.8
Lack of bias	4.1	4.6	4.0	4.7	4.6	4.9
Mean score	4.6	4.3	4.5	4.6	4.8	4.8
Word Count	138	128	169	151	131	116
Ratio Impression/ Findings Word	0.62	0.68	0.76	0.81	0.59	0.62

	DeepSeek-R1	Llama 3.1-405B405B	ChatGPT-o1
CR	0.21 (3/14)	0.07 (1/14)	0.29 (4/14)
PR	0.70 (7/10)	1 (10/10)	0.90 (9/10)
SD	0.65 (22/34)	0.65 (22/34)	0.82 (28/34)
PD	1 (42/42)	0.93 (39/42)	0.95 (40/42)
Overall Accuracy	0.74 (74/100)	0.72 (72/100)	0.81 (81/100)

	USMLE (n=323)	RECIST (n=100)	Medicilline (n=200)	NEJM accuracy (n=50)
DeepSeek-R1 vs Llama 3.1 405B	0.77	0.61	0.53	0.59
DeepSeek-R1 vs Chatgpt-o1	0.91	0.68	0.69	0.60
Llama 3.1 405B vs Chatgpt-o1	0.77	0.66	0.57	0.44

	Private Report	MIMI C	BMJ	NEJM Reasoning	NEJM Differential Diagnoses
DeepSeek-R1	0.909	0.906	0.849	0.876	0.873
Llama 3.1 405B	0.962	0.961	0.940	0.936	0.945
Chatgpt-o1	0.919	0.925	0.892	0.865	0.944

	USMLE	RECIST	Medicilline
DeepSeek-R1	0.962	0.803	0.611
Llama 3.1 405B	0.948	0.958	0.929
Chatgpt-o1	0.976	0.875	0.870

Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection	No software was used in data collection.
Data analysis	We used DeepSeek-R1 (DeepSeek), LLaMA 3.1 405B (Meta), and ChatGPT-o1(OpenAI)'s official scripts to evaluate the model performance on Azure Clouds. The Azure services used in this study are: https://ai.azure.com/explore/models/Meta-Llama-3.1-405B-Instruct/version/1/registry/azureml-meta?tid=77e89d61-570f-43b0-b9e4-634f462e34b8#details https://azuremarketplace.microsoft.com/en/marketplace/apps/metagenai.meta-llama-3-1-405b-instruct-offer?tab=overview https://ai.azure.com/explore/models/DeepSeek-R1/version/1/registry/azureml-deepseek?tid=77e89d61-570f-43b0-b9e4-634f462e34b8 SPSS v28.0 was used for statistical analysis.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

The USMLE dataset is available at <https://www.usmle.org/exam-resources>. The NEJM Case Challenges are available at <https://www.nejm.org/case-challenges>. The BMJ Endgames are available at <https://www.bmj.com/specialties/endgames>. The MIMIC-III dataset can be requested at <https://physionet.org/content/mimiciii/1.4/>. The Médiciline dataset (<http://www.mediciline.com/>) of multiple-choice questions for clinical diagnosis and management translated in English is available upon request for private/research use only, after agreement of Mediciline. Mediciline data requests should be sent to Mickael.Tordjman@mssm.edu. In response to the inquiry, the timeframe for responding to requests is approximately within 2 weeks of the request. The radiological reports used for summarization and RECIST classification are not available due to privacy issues.

Research involving human participants, their data, or biological material

Policy information about studies with [human participants or human data](#). See also policy information about [sex, gender \(identity/presentation\), and sexual orientation](#) and [race, ethnicity and racism](#).

Reporting on sex and gender

The NEJM and BMJ case challenges cover patients with diverse sex and gender. The imaging reports of the private dataset used for summarization and RECIST classification have a balanced sex distribution.

Reporting on race, ethnicity, or other socially relevant groupings

Race, ethnicity, or other socially relevant groupings were not relevant to the objectives or analyses of this study and therefore were not collected or analyzed. This study focused solely on the comparative performance of the Large Language Models across different medical tasks without considering these social variables.

Population characteristics

The population characteristics in this study are related to the challenges presented in the New England Journal of Medicine and the BMJ Endgames. These sources provided the basis for the tasks assigned to the evaluated Large Language Models, focusing on their application in diverse and complex analytical settings.

Recruitment

The reports of the private dataset were randomly selected from the Radiology Imaging System between 2010 and 2023 for the summarization task. Two cancer imaging reports of chest-abdomen-pelvis CT-scans were selected consecutively for the RECIST classification task. For the case challenges of the medical journals, we included all available cases from NEJM case challenges and the most recent cases from BMJ Endgames including a management section.

Ethics oversight

This study was approved by the Institutional Review Board of the Icahn School of Medicine at Mount Sinai to use the anonymized diagnostic imaging reports, in accordance with the institution's Federalwide Assurances (FWA00005656, FWA00005651), under IRB protocol number GCO# 20-2199-00001-01.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size

The publicly available datasets used in this study include NEJM Case Challenges published up to March 5, 2025, the 50 most recent BMJ Endgames management cases up to the same date, and the full set of USMLE questions (excluding image-based cases), all sourced from online resources. For the MIMIC-III dataset, we randomly sampled 20 cases per imaging modality (CT, MRI, PET, X-ray, and ultrasound) to ensure balanced coverage across modalities while maintaining feasibility for manual evaluation and computational analysis. For report summarization, we also randomly sampled 20 cases per imaging modality (CT, MRI, PET, X-ray, and ultrasound) from an additional private dataset. For RECIST 1.1 classification, two consecutive cancer imaging reports of 100 patients were randomly selected from clinical notes. 200 questions were randomly chosen for the Mediciline QA task. These sample sizes were chosen to provide adequate representation for each task while ensuring that performance evaluation could be conducted effectively within the scope of this study.

Data exclusions

This study strictly focuses on text-based evaluations of LLM capabilities. Image-based cases from NEJM, BMJ, and USMLE were excluded during the study design because they require multi-modal version of LLMs.

Replication

Each model was evaluated three times per task on Azure Cloud, with all runs completing successfully.

Randomization

The imaging reports were randomly selected from the RIS system between 2010 and 2023. Two cancer imaging reports of chest-abdomen-pelvis CT-scans were selected consecutively for the RECIST classification task. The “Findings” section of 200 imaging reports, 100 from an outpatient radiology facility in the New York Metropolitan area and 100 from MIMIC reports (each with 20 examples of X-rays, ultrasound, CT-scan, MRI and PET/CT reports), were provided to the three LLMs for summarization. Since our study evaluates the performance of LLMs in summarizing medical reports rather than assessing clinical interventions or comparing patient outcomes, random allocation or covariate control typically applied in interventional studies is not relevant here.

Blinding

The human readers participating in evaluation of the report summarization generated by the LLMs were blinded from which model provided the outputs.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	<input checked="" type="checkbox"/> Involved in the study <input type="checkbox"/> Antibodies <input checked="" type="checkbox"/> Eukaryotic cell lines <input checked="" type="checkbox"/> Palaeontology and archaeology <input checked="" type="checkbox"/> Animals and other organisms <input checked="" type="checkbox"/> Clinical data <input checked="" type="checkbox"/> Dual use research of concern <input checked="" type="checkbox"/> Plants
-----	--

Methods

n/a	<input checked="" type="checkbox"/> Involved in the study <input checked="" type="checkbox"/> ChIP-seq <input checked="" type="checkbox"/> Flow cytometry <input checked="" type="checkbox"/> MRI-based neuroimaging
-----	---

Plants

Seed stocks

Report on the source of all seed stocks or other plant material used. If applicable, state the seed stock centre and catalogue number. If plant specimens were collected from the field, describe the collection location, date and sampling procedures.

Novel plant genotypes

Describe the methods by which all novel plant genotypes were produced. This includes those generated by transgenic approaches, gene editing, chemical/radiation-based mutagenesis and hybridization. For transgenic lines, describe the transformation method, the number of independent lines analyzed and the generation upon which experiments were performed. For gene-edited lines, describe the editor used, the endogenous sequence targeted for editing, the targeting guide RNA sequence (if applicable) and how the editor was applied.

Authentication

Describe any authentication procedures for each seed stock used or novel genotype generated. Describe any experiments used to assess the effect of a mutation and, where applicable, how potential secondary effects (e.g. second site T-DNA insertions, mosaicism, off-target gene editing) were examined.