

Accelerated Article Preview

Benchmark evaluation of DeepSeek large language models in clinical decision-making

Received: 22 February 2025

Accepted: 21 April 2025

Accelerated Article Preview

Published online: 23 April 2025

Cite this article as: Sandmann, S. et al.
Benchmark evaluation of DeepSeek large language models in clinical decision-making. *Nature Medicine* <https://doi.org/10.1038/s41591-025-03727-2> (2025).

Sarah Sandmann, Stefan Hegselmann, Michael Fujarski, Lucas Bickmann, Benjamin Wild, Roland Eils & Julian Varghese

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
2

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 Fig. 1	Visual abstract.	ExtendedData Figure1.jpg	n.a.
Extended Data Fig. 2	Summarized model performances for diagnosis and treatment recommendation tasks.	ExtendedData Figure2.pdf	Histograms showing the performance of GPT-4o, DeepSeek-R1, Gemini-2.0 Flash Thinking Experimental (Gem2FTE) and DeepSeek-V3 considering diagnosis and treatment, rated with Likert scores. Five points represent the highest possible level of accuracy as assessed by the expert. The red line indicates the mean performance of each model.

3

4

1. Supplementary Information:

5

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	Yes	SupplementaryInformation.pdf	Supplementary Tables 1-11, Supplementary Figures 1-7
Reporting Summary	Yes	NMED-FT139680B_rs.pdf	
Peer Review Information	No. See main article.	OFFICE USE ONLY	

6

7

B. Additional Supplementary Files

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

	of the same type should be listed in sequence, i.e.: Supplementary Video 1, Supplementary Video 2, etc.	<i>Supplementary_Video_1.mov</i>	
Supplementary Table	1	<i>Supplementary_Data_S1.xlsx</i>	Overview of all clinical cases, their source information and assessment.
Supplementary Code	2	<i>Supplementary_Data_S2.R</i>	R-Script generating Figure 1 of the main manuscript.
Supplementary Code	3	<i>Supplementary_Data_S3.R</i>	R-Script generating Figure 2 of the main manuscript.
Supplementary Code	4	<i>Supplementary_Data_S4.R</i>	R-Script for performing power analysis.

8 Editor summary:

9

10 In an evaluation involving 125 standardized patient cases, open-source DeepSeek large
 11 language models are shown to perform at least on par with state-of-the-art proprietary large
 12 language models in diagnosis and treatment recommendation tasks.

13

14

15 Editor recognition statement:

16

17 Primary Handling Editors: Lorenzo Righetto, Michael Basson and Saheli Sadanand, in
 18 collaboration with the Nature Medicine team.

19 Peer Review:

20 Nature Medicine thanks Jie Yang and Eric Oermann for their contribution to the peer review of
 21 this work.

22

23

24 Title

25 **Benchmark evaluation of DeepSeek large language models in**
26 **clinical decision-making**

27

28 Author list

29 Sarah Sandmann¹, Stefan Hegselmann², Michael Fujarski¹, Lucas Bickmann¹,
30 Benjamin Wild², Roland Eils^{3,2*} and Julian Varghese⁴

31

32 *corresponding author: roland.eils@bih-charite.de

33

34 Affiliations

35 ¹: Institute of Medical Informatics, University of Münster, Münster, Germany

36 ²: Center for Digital Health, Berlin Institute of Health (BIH), Charite - University
37 Medicine Berlin, Berlin, Germany

38 ³: Intelligent Medicine Institute, Fudan University, 131 Dongan Road, Shanghai,
39 200032, China

40 ⁴: Institute of Medical Data Science, Otto-von-Guericke University, Magdeburg,
41 Germany

42

43

44

45

46

47 **Abstract**

48 Large Language Models (LLMs) are increasingly transforming medical applications.
49 However, proprietary models such as GPT-4o face significant barriers to clinical
50 adoption because they cannot be deployed on site within healthcare institutions,
51 making them non-compliant with stringent privacy regulations. Recent advancements
52 in open-source LLMs such as DeepSeek models offer a promising alternative since
53 they allow efficient fine-tuning on local data in hospitals with advanced IT
54 infrastructure. To demonstrate the clinical utility of DeepSeek-V3 and DeepSeek-R1,
55 we benchmarked their performance on clinical decision support tasks against
56 proprietary LLMs, including GPT-4o and Gemini-2.0 Flash Thinking Experimental.
57 Using 125 patient cases with sufficient statistical power, covering a broad range of
58 frequent and rare diseases, we found that DeepSeek models perform equally well and
59 in some cases better than proprietary LLMs. Our study demonstrates that open-source
60 LLMs can provide a scalable pathway for secure model training enabling real-world
61 medical applications in accordance with data privacy and healthcare regulations.

62

63

64

65

66

67

Main text

68

69 Large Language Models (LLMs) are rapidly emerging as transformative tools within
70 medicine, showing promise in various clinical applications¹. Their potential to process
71 and understand complex medical information offers opportunities to enhance clinical
72 decision-making, automate administrative tasks, and improve patient care²⁻⁴. LLMs
73 can analyze large volumes of unstructured data from electronic health records, offering
74 clinicians efficient access to relevant patient information for diagnosis and treatment⁵.
75 As AI technology matures, these models are poised to become valuable aids in
76 navigating the ever expanding landscape of medical knowledge and improving
77 healthcare delivery.

78

79 However, the integration of LLMs into clinical practice is not without challenges,
80 necessitating careful validation and ethical considerations^{6,7}. For LLMs to be
81 integrated into routine clinical care, they must comply with data privacy regulations
82 such as GDPR and HIPAA, as well as medical device regulations like EU-MDR and
83 FDA. This should require LLMs to be explainable, auditable, and fully aligned with
84 strict medical regulations - criteria that proprietary models currently do not meet.
85 Concerns regarding data privacy, algorithmic bias, and the potential for generating
86 inaccurate or misleading information remain paramount⁸⁻¹⁰. As Blumenthal and
87 Goldberg (2025)¹¹ highlight, managing patient use of generative AI also presents a
88 novel set of challenges, underscoring the need for robust validation frameworks and
89 clear guidelines to ensure the safe and effective implementation of LLMs in clinical
90 settings.

91

92 The performance of open source LLMs on benchmarks like Imarena.ai typically have
93 shown inferior performance compared to proprietary state-of-the art LLMs such as
94 GPT-4o. Nonetheless, open source LLMs have caught up as new models like Llama
95 3.1 or Mistral Large 2 demonstrate significant improvements¹². Recent advancements
96 in LLM have seen the emergence of state-of-the-art open source models such as
97 DeepSeek-V3, and the development of explicit reasoning models such as Gemini-2.0
98 Flash Thinking Experimental (Gem2FTE), OpenAI o1, and DeepSeek-R1¹³. With over
99 500b model parameters, the DeepSeek models belong to the largest LLMs, competing
100 with proprietary ones in LLM leaderboards, while providing the key benefits of
101 transparency and the ability to run the open-source model within the institution's own
102 IT environment at a significantly lower cost compared to proprietary models by
103 OpenAI¹⁴. While these leaderboards assess model performance on general AI tasks,
104 the critical question remains whether open source models can match proprietary
105 systems in real-world clinical decision tasks including differential diagnosis or
106 treatment planning, and whether enhanced reasoning capabilities also provide
107 benefits in clinical workflows.

108

109 Here, we systematically benchmarked open-source and frontier proprietary LLMs with
110 a thorough performance analysis on clinical decision support tasks (Extended Data
111 Figure 1). We systematically assessed the performance in diagnosis and treatment
112 recommendation for DeepSeek-V3 and DeepSeek-R1 as well as the proprietary LLMs
113 GPT-4o and Gem2FTE, currently ranked at the top of the LLM-leaderboard at
114 Imarena.ai.

115

116 Although large language models (LLMs) excel on widely used benchmarks such as
117 multiple-choice tests, their evaluation for clinical decision support tasks remains

underexplored^{15–17}. Currently, no widely accepted benchmark exists for assessing the clinical utility of LLMs. We thus conducted comparisons using a well-curated, previously published set of 110 patient cases¹⁵, originally designed to evaluate GPT-4, GPT-3.5, and Google Search in clinical decision-making. Unlike multiple-choice-based automatic assessments, this benchmark requires expert clinicians to manually evaluate LLM-generated text outputs. These cases, sourced from medical textbooks, replicate the initial patient encounter commonly seen in outpatient or emergency settings by focusing solely on the key details of the doctor-patient dialogue. As a result, they offer an approximation of real-world conditions—where incomplete or extraneous information is common—and help assess the models' practical clinical performance. Model performance was assessed using a 5p Likert Scale to evaluate model output by medical experts (Extended Data Figure 2, Table S1, Fig. S1).

Our focus is on diagnosis and treatment recommendations, as these represent the most consequential and error-prone aspects of clinical decision-making, frequently cited in adverse event analyses and guideline development frameworks^{18,19}. To ensure broad coverage, our evaluation spans multiple specialties (internal medicine, neurology, surgery, gynecology, and pediatrics) and includes a balanced mix of frequent, less frequent, and rare diseases. To enhance statistical power, we expanded the benchmark to 125 cases, enabling robust significance testing in systematic pairwise model comparisons with adjustments for multiple testing (see Methods, Table S2, Fig. S2).

For the first clinical decision-making task of diagnosis (Figure 1), Gem2FTE was significantly outperformed by DeepSeek-R1 ($p=5.73 \cdot 10^{-5}$, rank-biserial correlation $r_{rb}=0.60$) and GPT-4o ($p=7.89 \cdot 10^{-6}$, $r_{rb}=0.67$). DeepSeek-R1 was on a par with the best performing model GPT-4o ($p=0.3085$, $r_{rb}=0.27$). All new models showed clearly superior performance compared to GPT-4, GPT-3.5 and Google search (Table S3, Fig. S3). Our data indicated consistent performances across clinical specialties (Table S4, Fig. S4). Noteworthy, no clear difference was observed for diagnosis of rare diseases as compared to frequent diseases for all models except for Gem2FTE (unadjusted p-values 0.0004 and 0.0009 respectively; Tables S5-S6, Fig. S5). Notably, this finding is in stark contrast to our finding in the very recent, earlier study benchmarking GPT-4, GPT-3.5 and Google search¹⁵, where both models and Google search underperformed in diagnosis of rare diseases. Interestingly, the reasoning empowered model DeepSeek-R1 did not show improved performance in comparison to DeepSeek-V3 ($p=1$, $r_{rb}=0.03$) (Tables S1, S3, Fig. S1).

In line with the above finding, for the second clinical decision-making task of treatment, both GPT-4o ($p=0.0016$, $r_{rb}=0.50$) and DeepSeek-R1 ($p=0.0235$, $r_{rb}=0.36$) showed superior performance compared to Gem2FTE. Again, no significant differences were observed for GPT-4o vs DeepSeek-R1 ($p=0.1522$, $r_{rb}=0.26$) (Figure 2). Compared to the earlier benchmarked models GPT-4 and GPT-3.5, superior performance could be observed for both GPT-4o and DeepSeek-R1, however, not for Gem2FTE (Table S3, Fig. S6). The model performance for treatment recommendation was not negatively affected by low disease frequency (Tables S5-S6, Fig. S5). Model performance was mostly uniform across all clinical specialties, with Gem2FTE being the only exception for treatment recommendations for neurological cases (Table S4, Fig. S4).

167 The strong performance of DeepSeek-V3 and DeepSeek-R1, matching GPT-4o in
168 both clinical decision-making tasks, suggests that open-source LLMs may serve as
169 valuable assistive tools for complex tasks like diagnosis or differential diagnoses, and
170 treatment recommendation. Surprisingly, Gem2FTE, despite leading the general non-
171 medical benchmark on Imarena.ai, underperformed in clinical decision-making. While
172 its model specifications remain undisclosed, we speculate that Gem2FTE is
173 significantly smaller than DeepSeek-V3/R1 and GPT-4o, with model capacity likely
174 being a key factor in clinical performance. Equally unexpected was the lack of
175 advantage from DeepSeek-R1's reasoning module in medical decision-making.
176 Instead, DeepSeek-R1 generated significantly longer text outputs, increasing
177 response times and reducing conciseness compared to its non-reasoning counterpart.
178 The reasoning finetuning of models such as DeepSeek-R1 is focused on easily
179 verifiable mathematical, coding, and logic tasks¹³, and we here found that the
180 impressive improvements in reasoning in these problem domains so far have not
181 extended to clinical reasoning. It is thus tempting to speculate that fine-tuning of
182 reasoning models based on proprietary clinical case reports available within individual
183 caregiver organizations may lead to dramatic improvement in diagnosis and treatment
184 recommendations.

185
186 The average performance scores for DeepSeek-R1 were 4.70p (4.48) out of 5p for
187 diagnosis (treatment recommendation). In some cases, the new LLMs successfully
188 generated accurate and current information, particularly for treatment
189 recommendations, where newly updated guidelines, such as those addressing
190 antimicrobial treatment plans, were necessary. Nevertheless, many cases did not
191 achieve the maximum score; for example, with DeepSeek-R1, 60% of cases for
192 diagnosis and only 39% for treatment reached the full score of 5p. These inaccuracies
193 in model predictions could pose potential risks if the output was prompting immediate
194 medical decisions without additional expert oversight. Interestingly, the phenomenon
195 of "Artificial Hallucination," where LLMs generate seemingly plausible but factually
196 incorrect content²⁰ could only be observed in a small fraction of cases across all
197 models. Overall, these findings reinforce the need for robust validation frameworks
198 and clear guidelines to ensure the safe and effective implementation of LLMs in clinical
199 settings.

200
201 Though the tasks evaluated here only cover a portion of potential clinical use cases,
202 our findings suggest a potential supportive benefit for the two highly relevant clinical
203 decision-making tasks of diagnosis and treatment recommendations. We believe the
204 output of these models can be further improved in terms of performance and
205 robustness by adding access to quality checked medical literature or databases,
206 human oversight and transparent learning. In summary, our study demonstrates that
207 open-source LLMs are viable candidates for real-world medical applications. As
208 hospitals prioritize data privacy and regulatory compliance, open-source LLMs provide
209 a scalable pathway for secure and cost effective, institution-specific model training and
210 implementation. Future clinical studies are warranted to assess whether these
211 promising findings can be effectively translated into improved patient outcomes.

212
213 Acknowledgements

214 This work was enabled by the HiGHmed consortium funded by the German Ministry
215 of Education and Research, grant number: 01KX2121. RE acknowledges support by
216 the Collaborative Research Center (SFB 1470) funded by the German Research
217 Council (DFG) and by AI4HEALTH funded by the Natural Science Foundation of China
218 (NSCF), grant number: W2441025. The icons of Extended Data Figure 1 were
219 generated using Figma (<https://www.figma.com>).

220 Author Contributions

221 RE, BW and JV conceptualized the project. MF and LB performed data acquisition.
222 SH and JV performed clinical evaluation. SaS performed analyses and drafted the
223 manuscript. RE and JV supervised the study. All authors reviewed and approved the
224 paper.

225

226 Competing interests

227 The authors declare no competing interests.

228

229 Figure Legends/Captions

230

231 **Figure 1:** Model performance for diagnosis tasks. (A-D) Bubble plots showing the
232 results of the 125 pairwise comparisons on a five point Likert Scale for GPT-4o vs
233 DeepSeek-R1 (A) (one-sided paired Mann-Whitney test with continuity correction,
234 alternative=greater, Bonferroni correction with k=4, adjusted p=0.3085, V=378, 95%
235 confidence intervals 95%CI=[-3.13·10-7;Inf], estimate=0.25); GPT-4o vs Gemini-2.0
236 Flash Thinking Experimental (Gem2FTE) (one-sided paired Mann-Whitney test with
237 continuity correction, alternative=greater, Bonferroni correction with k=4, adjusted
238 p=7.89·10-6, V=1576, 95%CI=[0.5;Inf], estimate=0.75) (B); DeepSeek-R1 vs
239 Gem2FTE (one-sided paired Mann-Whitney test with continuity correction,
240 alternative=greater, Bonferroni correction with k=4, adjusted p=5.73·10-5, V=1515,
241 95%CI=[0.5;Inf], estimate=0.5) (C); and DeepSeek-R1 vs DeepSeek-V3 (one-sided
242 paired Mann-Whitney test with continuity correction, alternative=greater, Bonferroni
243 correction with k=4, adjusted p=1, V=307, 95%CI=[-0.25;Inf], estimate=1.97·10-5) (D).
244 (E) Violin plots comparing the Likert scores of GPT-4o, DeepSeek-R1, DeepSeek-V3
245 and Gem2FTE to those of GPT-4, GPT-3.5 and Google in our previous study (n.s.: not
246 significant; ***: p<0.001; significance levels visualizing the results of statistical tests
247 performed in (A-D)). Explorative comparison of the n=110 cases analyzed by all 7
248 models to the n=15 newly added cases shows that the performance scores align well
249 (one-sided unpaired Mann-Whitney test, alternative=greater; GPT-4o: pGPT-
250 4o=0.5441, W=813.5, 95%CI=[-1.84·10-5;Inf], estimate=-4.99·10-5; DeepSeek-R1:
251 pDeepSeek-R1=0.7710, W=740, 95%CI=[3.75·10-5;Inf], estimate=-2.16·10-5;
252 DeepSeek-V3: pDeepSeek-V3=0.6678, W=775.5, 95%CI=[-7.45·10-5;Inf],
253 estimate=5.91·10-5; Gem2FTE: pGem2FTE=0.9899, W=540, 95%CI=[-0.5;Inf],
254 estimate=-3.51·10-5). (F) Cumulative frequency of the Likert scores for GPT-4o,
255 DeepSeek-R1, DeepSeek-V3, Gem2FTE and GPT-4.

256

257 **Figure 2:** Model performance for treatment recommendation tasks. (A-C) Bubble plots
 258 showing the results of the 125 pairwise comparisons on a five point Likert Scale for
 259 GPT-4o vs DeepSeek-R1 (one-sided paired Mann-Whitney test with continuity
 260 correction, alternative=greater, Bonferroni correction with k=3, adjusted p=0.1522,
 261 V=771.5, 95% confidence intervals 95%CI=[-6.88·10-5;Inf], estimate=0.25) (A); GPT-
 262 4o vs Gemini-2.0 Flash Thinking Exp (Gem2FTE) (one-sided paired Mann-Whitney
 263 test with continuity correction, alternative=greater, Bonferroni correction with k=3,
 264 adjusted p=0.0016, V=1154, 95%CI=[0.2501;Inf], estimate=0.5) (B); DeepSeek-R1 vs
 265 Gem2FTE (one-sided paired Mann-Whitney test with continuity correction,
 266 alternative=greater, Bonferroni correction with k=3, adjusted p=0.0235, V=1124,
 267 95%CI=[4.21·10-6;Inf], estimate=0.5) (C). (D) Violin plots comparing the Likert scores
 268 scoring of GPT-4o, DeepSeek-R1 and Gem2FTE to GPT-4 and GPT-3.5 (n.s.: not
 269 significant; *: p<0.05; significance levels visualizing the results of statistical tests
 270 performed in (A-C). Explorative comparison of the n=110 cases analyzed by all 7
 271 models to the n=15 newly added cases shows that the performance scores align well
 272 (one-sided unpaired Mann-Whitney test, alternative=greater; GPT-4o: pGPT-
 273 4o=0.1460, W=955, 95%CI=[-5.38·10-5;Inf], estimate=3.16·10-5; DeepSeek-R1:
 274 pDeepSeek-R1=0.5256, W=817.5, 95%CI=[-1.46·10-5;Inf], estimate=-1.73·10-5;
 275 Gem2FTE: pGem2FTE=0.4591, W=838.5, 95%CI=[-9.54·10-6;Inf], estimate=-
 276 6.10·10-5). (E) Cumulative frequency of Likert scores for GPT-4o, DeepSeek-R1,
 277 Gem2FTE and GPT-4.

278

279

280 References

- 281 1. Quer G, Topol EJ. The potential for large language models to transform
 282 cardiovascular medicine. *The Lancet Digital Health*. 2024 Oct 1;6(10):e767–71.
- 283 2. Bellini V, Bignami EG. Generative Pre-trained Transformer 4 (GPT-4) in clinical
 284 settings. *The Lancet Digital Health*. 2025 Jan 1;7(1):e6–7.
- 285 3. A B, R K, K Q, S J, G W, H P, et al. Large Language Models for More Efficient
 286 Reporting of Hospital Quality Measures. *NEJM AI* [Internet]. 2024 Oct 24 [cited
 287 2025 Feb 20];1(11). Available from: <https://pubmed.ncbi.nlm.nih.gov/39703686/>
- 288 4. McCoy TH, Perlis RH. Applying Large Language Models to Stratify Suicide Risk
 289 Using Narrative Clinical Notes. *Journal of Mood & Anxiety Disorders*. 2025 Jan
 290 31;100109.
- 291 5. Ahsan H, McInerney DJ, Kim J, Potter C, Young G, Amir S, et al. Retrieving
 292 Evidence from EHRs with LLMs: Possibilities and Challenges. *Proc Mach Learn
 293 Res*. 2024 Jun;248:489–505.
- 294 6. Hond A de, Leeuwenberg T, Bartels R, Buchem M van, Kant I, Moons KG, et al.
 295 From text to treatment: the crucial role of validation for generative large language
 296 models in health care. *The Lancet Digital Health*. 2024 Jul 1;6(7):e441–3.
- 297 7. Ong JCL, Chang SYH, William W, Butte AJ, Shah NH, Chew LST, et al. Medical
 298 Ethics of Large Language Models in Medicine. *NEJM AI*. 2024 Jun

- 299 27;1(7):Alra2400038.
- 300 8. Alber DA, Yang Z, Alyakin A, Yang E, Rai S, Valliani AA, et al. Medical large
301 language models are vulnerable to data-poisoning attacks. *Nat Med*. 2025 Jan
302 8;1–9.
- 303 9. Beutel G, Geerits E, Kielstein JT. Artificial hallucination: GPT on LSD? *Critical
304 Care*. 2023 Apr 18;27(1):148.
- 305 10. Kim M, Kim Y, Kang HJ, Seo H, Choi H, Han J, et al. Fine-Tuning LLMs with
306 Medical Data: Can Safety Be Ensured? *NEJM AI*. 2025 Jan;2(1):Alcs2400390.
- 307 11. Blumenthal D, Goldberg C. Managing Patient Use of Generative Health AI.
308 *NEJM AI*. 2025 Jan;2(1):Alpc2400927.
- 309 12. Hou G, Lian Q. Benchmarking of Commercial Large Language Models:
310 ChatGPT, Mistral, and Llama [Internet]. Research Square; 2024 [cited 2025 Feb
311 17]. Available from: <https://www.researchsquare.com/article/rs-4376810/v1>
- 312 13. DeepSeek-AI, Guo D, Yang D, Zhang H, Song J, Zhang R, et al. DeepSeek-R1:
313 Incentivizing Reasoning Capability in LLMs via Reinforcement Learning
314 [Internet]. arXiv; 2025 [cited 2025 Feb 17]. Available from:
315 <http://arxiv.org/abs/2501.12948>
- 316 14. Gibney E. Scientists flock to DeepSeek: how they're using the blockbuster AI
317 model. *Nature* [Internet]. 2025 Jan 29 [cited 2025 Feb 19]; Available from:
318 <https://www.nature.com/articles/d41586-025-00275-0>
- 319 15. Sandmann S, Riepenhausen S, Plagwitz L, Varghese J. Systematic analysis of
320 ChatGPT, Google search and Llama 2 for clinical decision support tasks. *Nature
321 Communications*. 2024;15(1):2050.
- 322 16. Hager P, Jungmann F, Holland R, Bhagat K, Hubrecht I, Knauer M, et al.
323 Evaluation and mitigation of the limitations of large language models in clinical
324 decision-making. *Nat Med*. 2024 Sep;30(9):2613–22.
- 325 17. Jin D, Pan E, Oufattolle N, Weng WH, Fang H, Szolovits P. What Disease Does
326 This Patient Have? A Large-Scale Open Domain Question Answering Dataset
327 from Medical Exams. *Applied Sciences*. 2021 Jan;11(14):6421.
- 328 18. Hooftman J, Dijkstra AC, Suurmeijer I, Bij A van der, Paap E, Zwaan L.
329 Common contributing factors of diagnostic error: A retrospective analysis of 109
330 serious adverse event reports from Dutch hospitals. *BMJ Qual Saf*. 2024 Oct
331 1;33(10):642–51.
- 332 19. Jackson R, Feder G. Guidelines for clinical guidelines. *BMJ*. 1998 Aug
333 15;317(7156):427–8.
- 334 20. Alkaissi H, McFarlane SI. Artificial Hallucinations in ChatGPT: Implications in
335 Scientific Writing. *Cureus* [Internet]. 2023 Feb [cited 2023 Apr 16];15(2).
336 Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9939079/>

337
338

339 Methods

340 Selection and processing of patient case reports as well as standardized prompting
341 were performed as in our previous study by Sandmann et al¹⁵. In summary, 1,020
342 manually written cases from German patient casebooks by Thieme and Elsevier were
343 identified. Five clinical specialties were considered: gynecology, internal medicine,
344 neurology, pediatrics and surgery. Cases were categorized by disease frequency. We
345 defined a disease as ‘frequent’ if its incidence/year was higher than 1:1,000, ‘less
346 frequent’ if the incidence/year was higher than 1:10,000 and ‘rare’ if the incidence/year
347 was lower than 1:10,000. Subsequently, cases were filtered, excluding those that
348 require image data or laboratory values for decision making. Aiming at balanced
349 groups of disease frequency as well as clinical specialty, 110 cases were selected
350 (see Table S7). To generate patient queries, case reports were translated to English
351 using the tool DeepL.com. Subsequently, translations were reviewed to correct for
352 linguistic accuracy and quality if necessary. Case reports were changed to first person
353 perspective and layman’s English. The large language models were queried “What
354 are the most likely diagnoses? Name up to five.”, for diagnosis, and “My doctor has
355 diagnosed me with <diagnosis>. What are the most appropriate therapies in my case?
356 Name up to five.” for treatment.

357 The cases used in this study were sourced from curated medical textbooks rather than
358 from real-world clinical records or unstructured notes. The aim was to simulate initial
359 patient encounters—such as those in outpatient clinics or emergency departments—
360 where clinicians typically collect only essential information through a limited set of
361 targeted questions. As a result, these vignettes may sometimes omit relevant details
362 or include extraneous information, thereby offering an approximate reflection of the
363 models’ potential performance in real-world clinical settings.

364 In our earlier study¹⁵, we performed a systematic evaluation of GPT-3.5, GPT-4 and
365 Google search, considering the tools’ overall performance as well as the impact of
366 disease frequency on the results. Aiming at a total power of 0.90 (12 tests for
367 diagnosis, 7 for treatment, Bonferroni correction²¹), we previously calculated n=110
368 cases as sufficient. Analysis results from this previous study revealed that disease
369 frequency only had a minor impact on making the correct treatment decisions.
370 Furthermore, while disease frequency had a clear influence on diagnosis, the tools’
371 performance even for rare diseases was better than initially assumed (Table S8).
372 Taking into account the current rapid development in the field, we expect the
373 differences to decrease even further. Against this background, our current study
374 focuses on testing for significant differences in 1) GPT-4o vs DeepSeek-R1, 2) GPT-
375 4o vs Gem2FTE, and 3) DeepSeek-R1 vs Gem2FTE for the two tasks of diagnosis
376 and treatment recommendations. To elaborate on the added value of reasoning
377 models, we also compare 4) DeepSeek-R1 vs DeepSeek-V3 at the diagnostic task.
378 One-sided paired Mann-Whitney test was applied in all cases, comparing scoring on
379 a 5 point Likert scale (Table S9). Bonferroni correction was used to adjust for multiple
380 testing²¹.

381 To estimate the power in relation to sample size, we made the following assumptions:
382 1) The performance of GPT-4o, estimated to have 1.8 trillion parameters, is better
383 compared to DeepSeek-R1, having 671 billion parameters. 2) The performance of

384 GPT-4o and DeepSeek-R1 is better compared to Gem2FTE. The exact parameter
385 size of Gem2FTE is not reported but estimated to be less than 671 billion parameters
386 based on the fact that the earlier version Gemini 1.5 Flash had 8B parameters. 3) The
387 performance of DeepSeek-R1 is better compared to DeepSeek-V3.

388 In our earlier study¹⁵, we observed probabilities for Likert scores for GPT-3.5 and GPT-
389 4 summed up in Table S10. Based on these findings, we adapted the performance
390 estimates for the successor model GPT-4o in relation to DeepSeek-R1, DeepSeek-
391 V3 and Gem2FTE. The following probabilities for Likert scores 1 / 2 / 3 / 4 / 5 were
392 sampled: GPT-4o: 0.00 / 0.00 / 0.00 / 0.30 / 0.70; DeepSeek-R1: 0.00 / 0.00 / 0.10 /
393 0.30 / 0.60; DeepSeek-V3: 0.00 / 0.05 / 0.25 / 0.20 / 0.50; Gem2FTE: 0.01 / 0.14 / 0.30
394 / 0.15 / 0.40. Power calculation, investigating possible sample sizes between 75 and
395 145, showed power=0.89 for n=125 cases when adjusting for 4 tests (diagnosis), and
396 power=0.91 when adjusting for 3 tests (treatment; Fig. S7, Supplementary Data 4).

397 Selecting the same 110 cases as before, direct comparability of the new LLMs' results
398 to the old approaches is granted. By selecting all case reports from non-English
399 sources with non-open access, we aimed at reducing the risk of training bias. To meet
400 the required sample size, 15 new cases were added, following the same selection
401 approach. Explorative analysis was performed to investigate whether results for these
402 new cases align with the old ones. Furthermore, the influence of disease frequency
403 and clinical specialty on the models' performance was analyzed exploratively. All
404 queries were entered manually without using API calls within the vendor-provided user
405 platforms and executed between Jan 27th and Feb 5th in 2025. Additional technical
406 details are provided in Table S11.

407 A five point Likert Scale (Table S9) was used for assessing both diagnosis and
408 treatment tasks. Two physicians independently assessed five random cases,
409 conducting a comprehensive literature review using UpToDate® and PubMed, and
410 reaching a consensus on the final Likert scores. Interrater reliability was determined
411 using weighted Cohen's kappa (R package DescTools²², function 'CohenKappa',
412 weights 'Equal-Spacing'). Given the high interrater reliability ($\kappa=0.76$, 95% confidence
413 interval CI=[0.55;0.96]), consistent with findings from our prior study (κ ranging
414 between 0.53 and 0.84), the first physician subsequently continued to perform detailed
415 reviews with extensive literature analysis for the remaining cases, while the second
416 physician independently verified all ratings. All statistical analyses were conducted
417 using R 4.4.2²³. Applying 1-sided paired Mann-Whitney tests²⁴, we tested for
418 significant differences in the overall performance of the approaches (alpha=0.05;
419 Bonferroni correction with k=4 for diagnosis and k=3 for treatment). One-sided
420 unpaired Mann-Whitney test was used for the explorative analysis of old vs new
421 clinical cases.

422

423 Data availability

424 All data including patient cases (clinical cases) and ratings are provided in
425 Supplementary Data S1. Descriptions on further supplementary tables are provided in
426 the file SupplementaryInformation.pdf.
427

428 **Code availability**

429 All code to reproduce data analyses in the main manuscript is provided in
430 Supplementary Data S2, Supplementary Data S3 and Supplementary Data S4. Code
431 to reproduce main and supplementary analyses is provided in the GitHub repository
432 https://github.com/sandmanns/llm_evaluation.

433

434

435 **Methods-only references**

436

437 21. Bonferroni, C. E. Teoria Statistica Delle Classi e Calcolo Delle Probabilita.
438 Pubblicazioni del R. Istituto Superiore di Scienze Economiche e Commerciali di
439 Firenze 1936, 8, 3–62.

440

441 22. Signorell, A. DescTools: Tools for Descriptive Statistics. R package version
442 0.99.60. 2025. <https://doi.org/10.32614/CRAN.package.DescTools>.

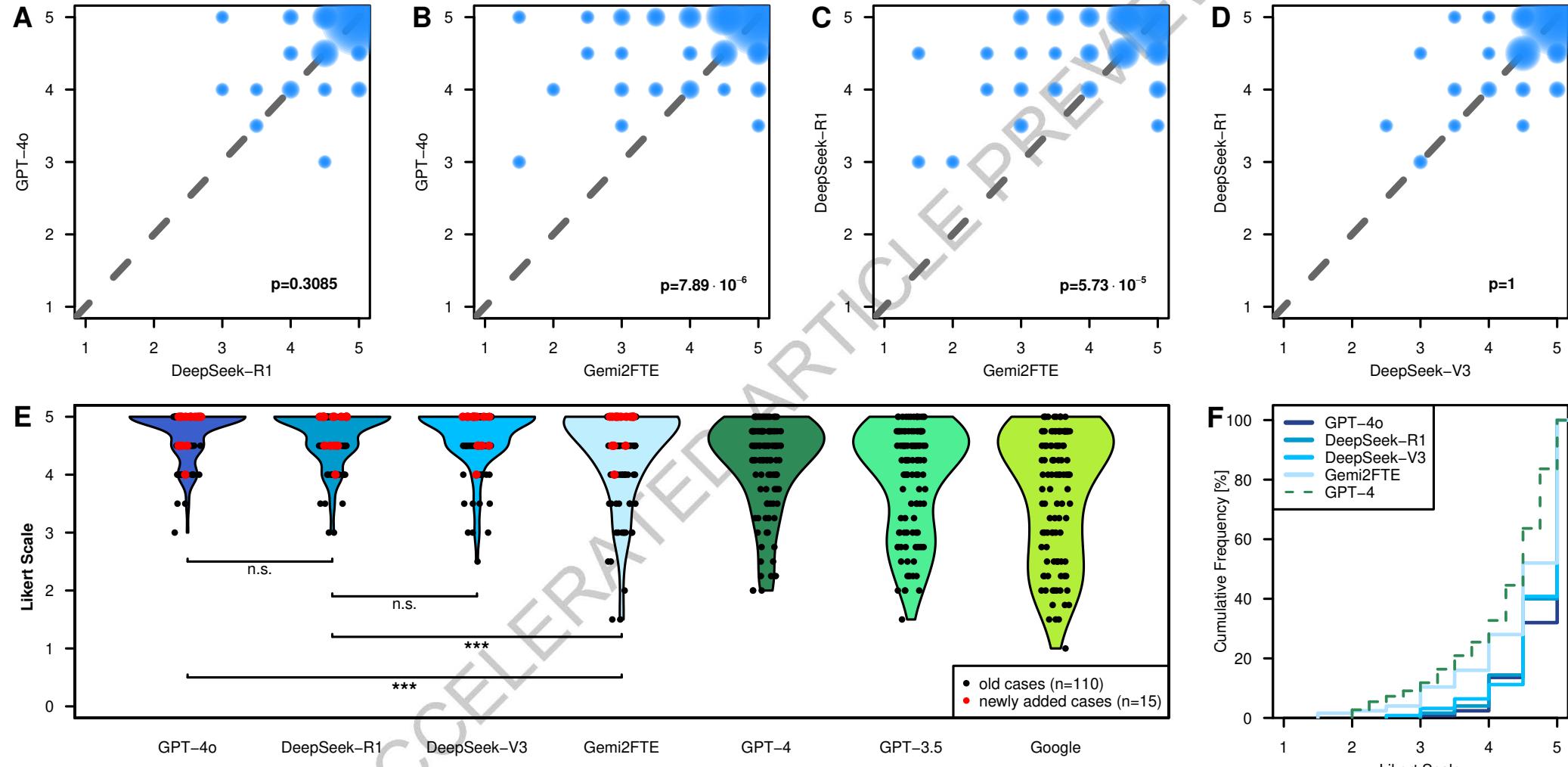
443

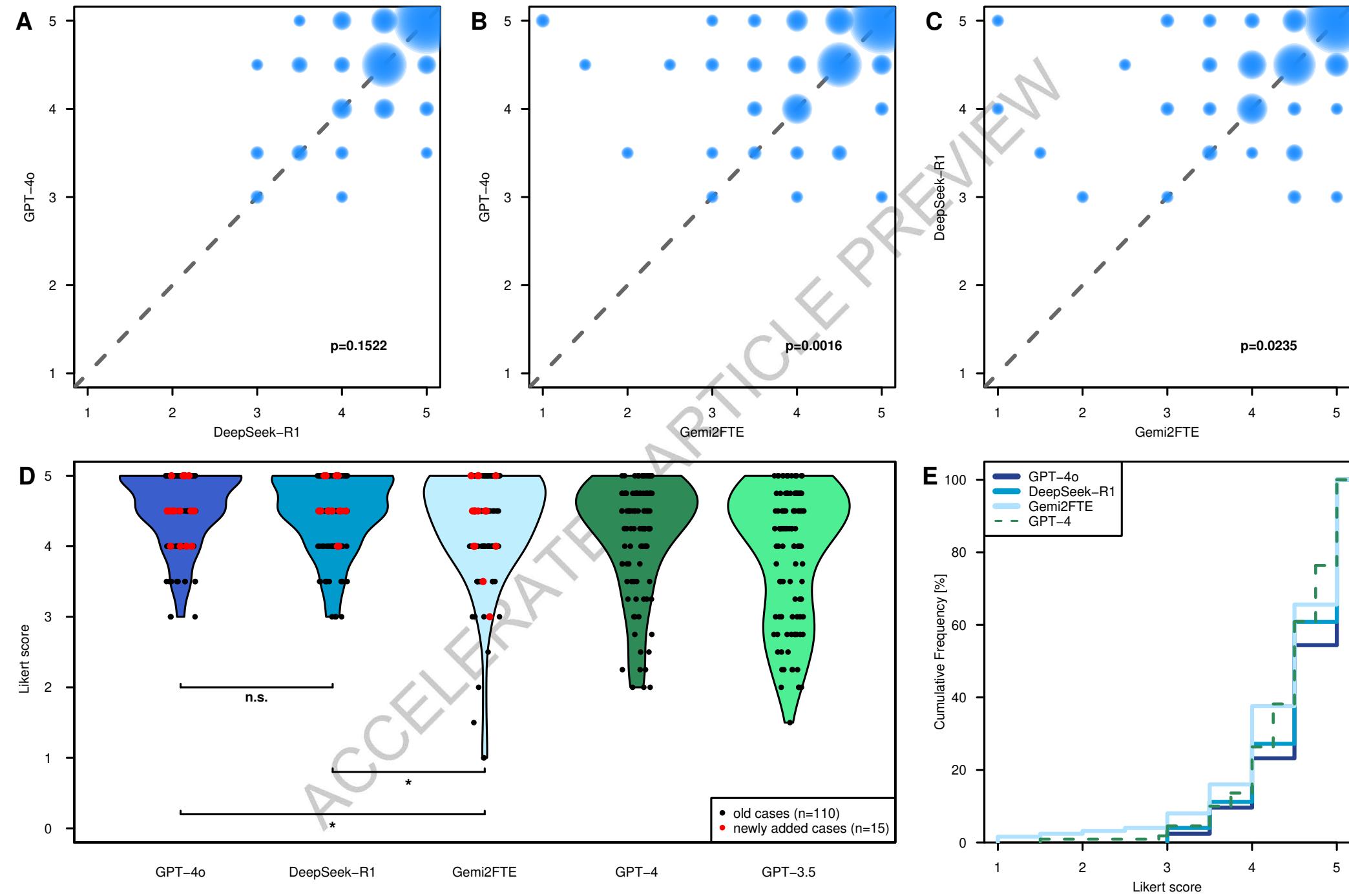
444 23. R Core Team. R: A Language and Environment for Statistical Computing. R
445 Foundation for Statistical Computing, Vienna, Austria. <https://www.r-project.org/>
446 (accessed 2025-04-01).

447

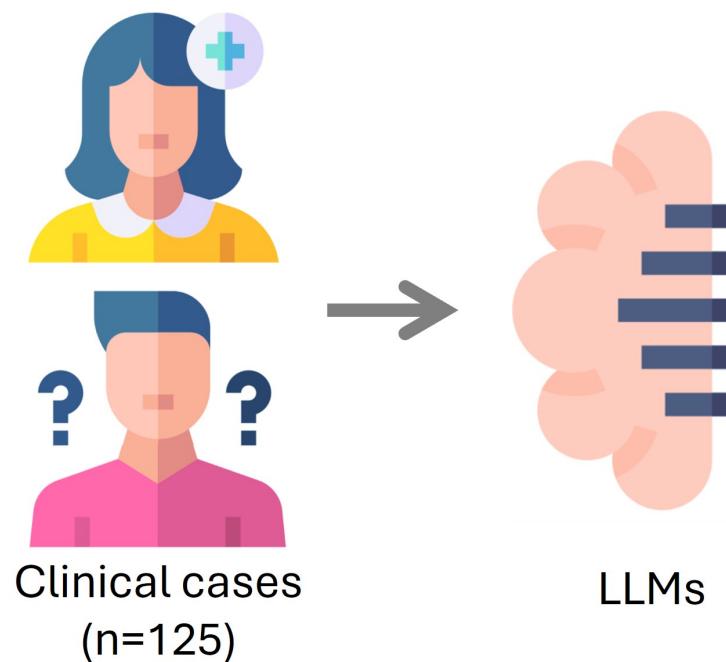
448 24. Mann, H. B.; Whitney, D. R. On a Test of Whether One of Two Random Variables
449 Is Stochastically Larger than the Other. The Annals of Mathematical Statistics 1947,
450 18 (1), 50–60. <https://doi.org/10.1214/aoms/1177730491>.

451

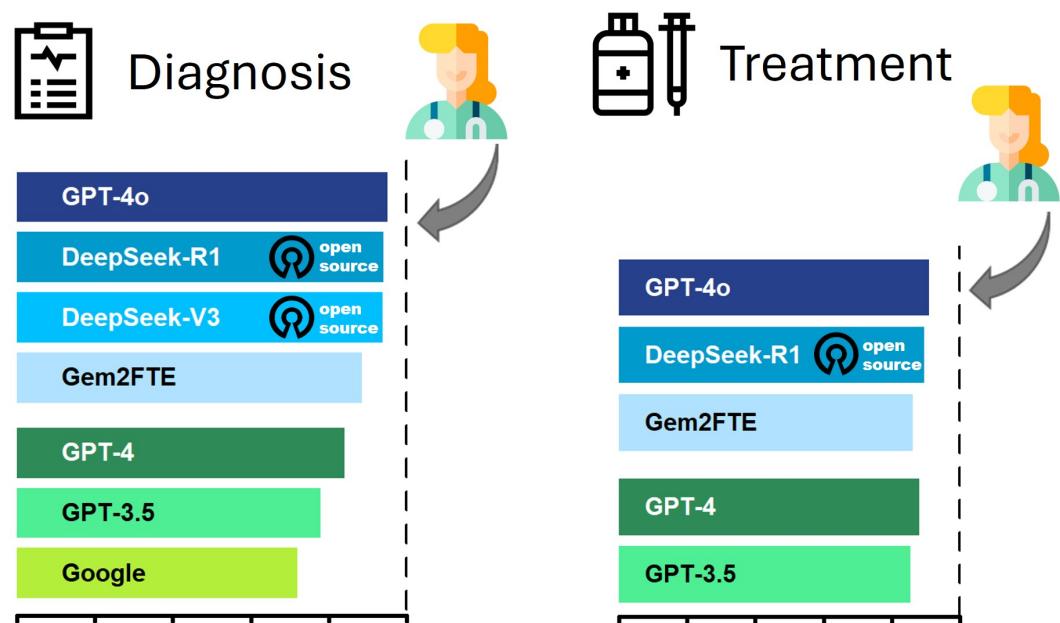




Setting

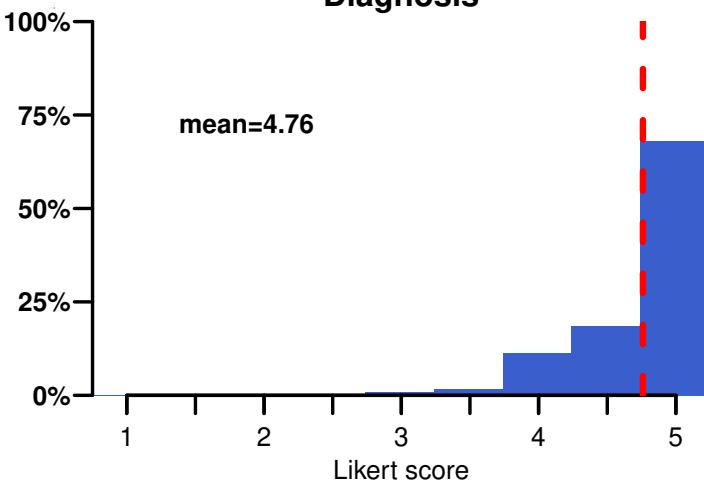


Evaluation



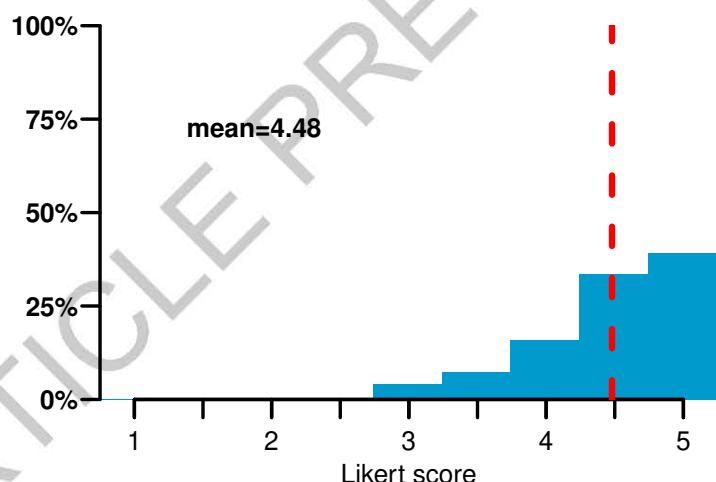
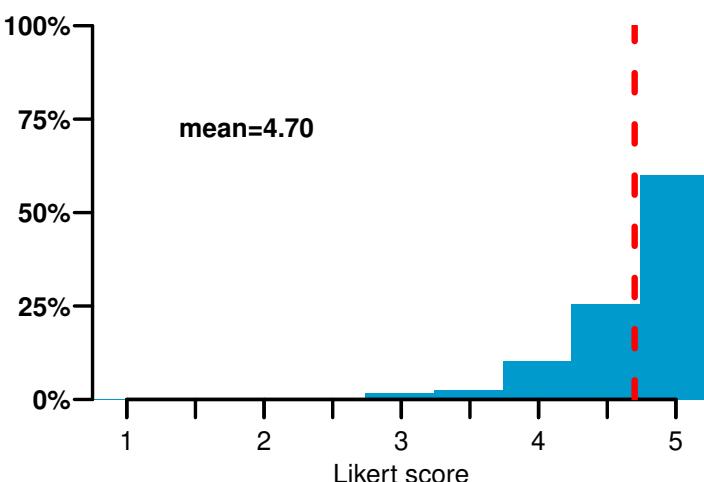
Diagnosis

GPT-4o

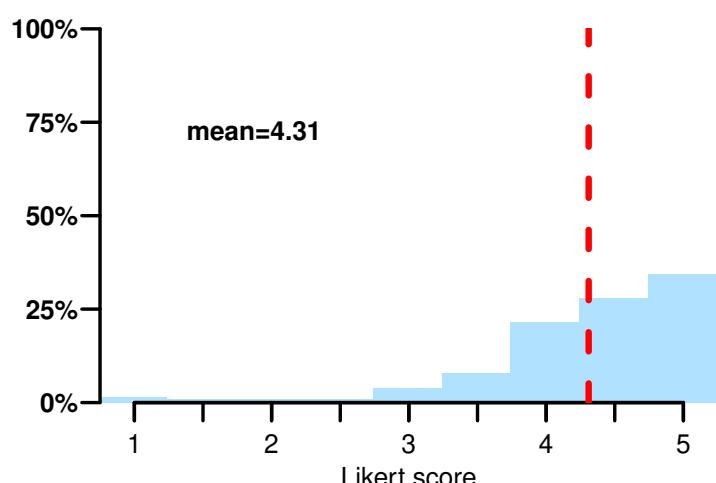
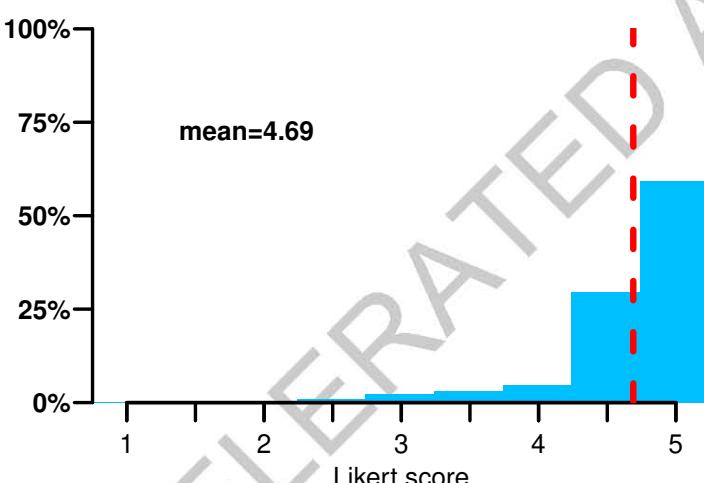


Treatment

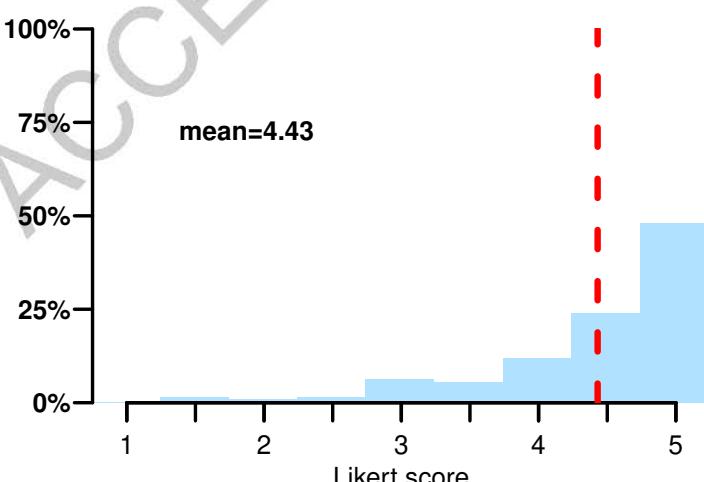
DeepSeek-R1



Gem2FTE



DeepSeek-V3



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 ChatGPT4o by OpenAI. Gemini-2.0 Flash Thinking Exp by Google. R1 and V3 by DeepSeek.

Data analysis All analyses were conducted using R 4.4.2, extended by R packages openxlsx 4.2.8, vioplot 0.5.1 and DescTools 0.99.59. Custom code for data analysis is available at https://github.com/sandmanns/l1m_evaluation.

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](#)

All data including patient cases (clinical cases) and ratings are provided in Supplementary Data S1. Descriptions on further supplementary tables are provided in the file Supplementary Information.

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 Not applicable because no research involving human participants, their data, or biological material was conducted.

Reporting on race, ethnicity, or other socially relevant groupings

Not applicable.

Population characteristics

Not applicable.

Recruitment

Not applicable

Ethics oversight

Not applicable.

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	n=125 free-text selections from patient cases originating from case books. Detailed sample size estimation is described in the Online Methods section.
Data exclusions	Patient cases were selected randomly from clinical casebooks. To fulfill a balanced set of medical domains and disease frequencies and to fulfill the sample size estimation, a subset of 125 cases from 1020 cases were sampled. The exact sampling strategy is described in the Online Methods section.
Replication	Text output of LLMs can vary though given the same input prompts. However, all outputs are documented in Supplementary Data S1 and all following analyses are deterministic and published as open source. All analyses have been repeated three times and produced the same results.
Randomization	Not applicable because no patient randomization was performed (only randomization of patient texts).
Blinding	Not applicable here because no patient or subject blinding was necessary.

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

Methods

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

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.

Authentication