

Comparison of Large Language Models in Answering Immuno-Oncology Questions: A Cross-Sectional Study

Giovanni Maria Iannantuono, MD^{1†}, Dara Bracken-Clarke, MD^{2†}, Fatima Karzai, MD¹,
Hyoyoung Choo-Wosoba, PhD³, James L. Gulley, MD, PhD²,
Charalampos S. Floudas, MD, DMSc, MS².

¹Genitourinary Malignancies Branch, Center for Cancer Research, National Cancer Institute, National Institutes of Health, Bethesda, MD, United States.

²Center for Immuno-Oncology, Center for Cancer Research, National Cancer Institute, National Institutes of Health, Bethesda, MD, United States.

³Biostatistics and Data Management Section, Center for Cancer Research, National Cancer Institute, National Institutes of Health, Bethesda, MD, United States.

[†]These authors contributed equally to this work

Corresponding Author:

Charalampos S. Floudas, MD, DMSc, MS

Center for Immuno-Oncology, Center for Cancer Research, National Cancer Institute

10 Center Drive, Bethesda, MD, 20892.

Building 10, Room 7N240A

Tel: 240-858-3032 - Email: charalampos.floudas@nih.gov

Running Head Title: LARGE LANGUAGE MODELS IN IMMUNO-ONCOLOGY

Keywords: Large language models; Artificial intelligence; Immuno-oncology; ChatGPT; Google Bard.

Author Contributions

Iannantuono GM: Conception/Design; Provision of study material; Collection and/or assembly of data; Data analysis and interpretation; Manuscript writing; Final approval of manuscript.

Bracken-Clarke D: Conception/Design; Provision of study material; Collection and/or assembly of data; Data analysis and interpretation; Manuscript writing; Final approval of manuscript.

Karzai F: Manuscript writing; Final approval of manuscript.

Choo-Wosoba H: Data analysis and interpretation; Manuscript writing; Final approval of manuscript.

Gulley JL: Manuscript writing; Final approval of manuscript.

Floudas CS: Conception/Design; Data analysis and interpretation; Manuscript writing; Final approval of manuscript.

47 **ABSTRACT**

48

49 **Background:** The capability of large language models (LLMs) to understand and generate human-
50 readable text has prompted the investigation of their potential as educational and management
51 tools for cancer patients and healthcare providers.

52 **Materials and Methods:** We conducted a cross-sectional study aimed at evaluating the ability of
53 ChatGPT-4, ChatGPT-3.5, and Google Bard to answer questions related to four domains of
54 immuno-oncology (Mechanisms, Indications, Toxicities, and Prognosis). We generated 60 open-
55 ended questions (15 for each section). Questions were manually submitted to LLMs, and responses
56 were collected on June 30th, 2023. Two reviewers evaluated the answers independently.

57 **Results:** ChatGPT-4 and ChatGPT-3.5 answered all questions, whereas Google Bard answered
58 only 53.3% ($p < 0.0001$). The number of questions with reproducible answers was higher for
59 ChatGPT-4 (95%) and ChatGPT3.5 (88.3%) than for Google Bard (50%) ($p < 0.0001$). In terms of
60 accuracy, the number of answers deemed fully correct were 75.4%, 58.5%, and 43.8% for
61 ChatGPT-4, ChatGPT-3.5, and Google Bard, respectively ($p = 0.03$). Furthermore, the number of
62 responses deemed highly relevant was 71.9%, 77.4%, and 43.8% for ChatGPT-4, ChatGPT-3.5,
63 and Google Bard, respectively ($p = 0.04$). Regarding readability, the number of highly readable
64 was higher for ChatGPT-4 and ChatGPT-3.5 (98.1%) and (100%) compared to Google Bard
65 (87.5%) ($p = 0.02$).

66 **Conclusion:** ChatGPT-4 and ChatGPT-3.5 are potentially powerful tools in immuno-oncology,
67 whereas Google Bard demonstrated relatively poorer performance. However, the risk of
68 inaccuracy or incompleteness in the responses was evident in all three LLMs, highlighting the
69 importance of expert-driven verification of the outputs returned by these technologies.

70 **IMPLICATIONS FOR PRACTICE**

71 Several studies have recently evaluated whether large language models may be feasible tools for
72 providing educational and management information for cancer patients and healthcare providers.
73 In this cross-sectional study, we assessed the ability of ChatGPT-4, ChatGPT-3.5, and Google
74 Bard to answer questions related to immuno-oncology. ChatGPT-4 and ChatGPT-3.5 returned a
75 higher proportion of responses, which were more accurate and comprehensive, than those returned
76 by Google Bard, yielding highly reproducible and readable outputs. These data support ChatGPT-
77 4 and ChatGPT-3.5 as powerful tools in providing information on immuno-oncology; however,
78 accuracy remains a concern, with expert assessment of the output still indicated.

79

80

81

82

83

84

85

86

87

88

89

90 **1. INTRODUCTION**

91 Large language models (LLMs) are a recent breakthrough in the domain of generative artificial
92 intelligence (AI) (1). Generative AI includes technologies based on “natural language processing”
93 (NLP) which uses computational linguistics and deep learning (DL) algorithms to enable
94 computers to interpret and generate human-like text (2). Large language models are complex
95 systems trained on large quantities of text data which are able to create new content in response to
96 prompts such as text, images, or other media (3). This versatility has led to the investigation of
97 their potential applications in the field of medicine and healthcare in light of its self-evident
98 potential benefits in these domains (4). Indeed, the availability of user-friendly tools able to
99 provide detailed, accurate and current information would be crucial in promoting patient and
100 healthcare providers’ education and awareness, particularly in the case of complex health
101 conditions like cancer (5).

102
103 Thus far, many studies have assessed the potential of ChatGPT, an advanced LLM based on a
104 generative pre-trained transformer (GPT) architecture, for providing screening and/or management
105 information in solid tumors (6). Following the rollout of ChatGPT, more LLMs trained on different
106 data were released, expanding the selection of these new AI-based tools. Consequently, an
107 increasing number of studies are investigating and comparing the potential ability of ChatGPT
108 with other LLMs as easy-to-use interfaces to gather information related to a specific cancer-related
109 topic (7). So far, initial evidence suggests a possible role of these technologies as “virtual
110 assistants” for healthcare professionals and patients in providing information about cancer,
111 unfortunately counterbalanced by a significant error rate. Therefore, further studies are needed to
112 investigate the potential applicability of these tools in other fields (7).

113 The past several years have seen profound changes in the field of immuno-oncology (IO). The
114 advent of immune-checkpoint inhibitors (ICIs) has paved the way towards a new era in cancer
115 treatment, enhancing the chance of long-term survival in patients with metastatic disease, and
116 providing new treatment options in earlier-stage settings (8). Presently, an increasing number of
117 cancer patients are either candidates for or already receiving ICIs or other immunotherapies,
118 subject to both the enormous potential benefits but also the immune-related adverse events that
119 may be caused by these treatments (9). In this context, LLMs may represent a valid tool for
120 healthcare professionals and patients (and their caregivers) receiving these treatments. Therefore,
121 we sought to assess and compare the ability of three prominent LLMs to provide educational and
122 management information in the IO field.

123
124
125
126
127
128
129
130
131
132
133
134
135

136 **2. MATERIALS AND METHODS**

137 **2.1 Large language models**

138 In this cross-sectional study we compared the performance of three LLMs: ChatGPT-3.5 (10),
139 ChatGPT-4 (10), and Google Bard (11). ChatGPT is an LLM based on the GPT architecture and
140 developed by OpenAI, a company based in San Francisco (USA). ChatGPT is built upon either
141 GPT-3.5 and GPT-4; the former is freely available to all the users, whereas the latter is an advanced
142 version with additional features and provided under the name “ChatGPT Plus” to paid subscribers
143 (10). Google Bard is based on the Pathways Language Model (PaLM) family of LLMs, developed
144 by GoogleAI (11).

145

146 **2.2 Questions and responses’ generation**

147 We generated 60 open-ended questions based on our clinical experience covering four different
148 domains of IO including “mechanisms” (of action), “indications” (for use), “toxicities”, and
149 “prognosis” (Suppl. Mat. A). In order to standardize assessment, particularly of “relevance” and
150 “accuracy”, and to reduce bias, a sample answer for each question was generated *a priori* prior to
151 question submission. Questions were manually and directly submitted to the web chat interfaces
152 of the three above-mentioned LLMs on June 30th 2023 and responses were collected (Suppl. Mat.
153 B). We assessed the reproducibility, accuracy, relevance, and readability (Table 1) of responses
154 provided by each LLM. Two reviewers (GMI and DBC) rated the answers independently. During
155 the rating process, reviewers were blinded to the LLM being assessed. Inconsistencies between the
156 reviewers were discussed with an additional reviewer (CSF) and resolved by consensus. Cohen's
157 kappa coefficient was calculated to evaluate inter-rater reliability during the rating process (12).

158

159 First, we assessed the ability of each LLM to provide reproducible responses. Therefore, each
160 individual question was submitted three times on each LLM. In the case of non-reproducible
161 answers, questions were not considered for further analysis. Subsequently, the accuracy, relevance,
162 and readability of responses deemed reproducible were assessed using a 3-point scale (Table 2)
163 (Figure 1). Reviewers graded the accuracy of answers according to available information as of
164 2021, as the training datasets of ChatGPT are updated to September 2021. Finally, word- and
165 character-counts were calculated for each answer.

166

167 **2.3 Statistical analyses**

168 Categorical variables were presented with proportions and numeric variables as measures of
169 central tendency. Comparisons between categorical variables were performed with two-sided
170 generalized Fisher's exact tests for testing any potential differences in these three LLMs. In the
171 case of numeric continuous variables, a Kruskal-Wallis test was utilized. Statistical tests were not
172 performed within each of the four domains, but rather were performed only to evaluate overall
173 performance by combining those four domains, due to insufficient sample sizes within each
174 domain (i.e., only up to 15 available observations). All statistical results should be interpreted as
175 exploratory; all statistical analyses were performed and all plots generated using R version 4.2.2
176 (The R Foundation for Statistical Computing, 2022). This study was conducted in accordance with
177 Strengthening the Reporting of Observational Studies in Epidemiology (STROBE) reporting
178 guidelines (13).

179

180

181

182 3. RESULTS

183 Assessment of inter-rater reliability with Cohen’s kappa during the rating process demonstrated
184 “strong” to “near perfect” agreement between reviewers (Suppl. Mat. C). ChatGPT-3.5 and
185 ChatGPT-4 provided at least one response to all questions (60 [100%]), while Google Bard
186 responded only to 32 (53.3%) queries ($p < 0.0001$). Specifically, the percentages of responses
187 provided by Google Bard were different across the four domains, with better performances in the
188 “mechanisms” (14 [93.3%]) and “prognosis” domains (13 [86.7%]) compared to the “indications”
189 (5 [33.3%]), and “toxicities” (0 [0%]) domains. Regarding reproducibility, the numbers of
190 questions with reproducible answers were similar between ChatGPT-3.5 and ChatGPT-4 (53
191 [88.3%] and 57 [95%], respectively), while it was lower (16 [50%]) for Google Bard ($p < 0.0001$).
192 Although ChatGPT-3.5 and ChatGPT-4 performed similarly across all domains, ChatGPT-4
193 achieved 100% reproducible responses in two domains (“mechanisms” and “indications”) in which
194 ChatGPT-3.5 achieved only 86.7%. Google Bard was variably capable and accurate across the
195 different sections. Despite a significant number of answers deemed reproducible in the
196 “mechanisms” (6 [40%]) and “prognosis” (9 [60%]) sections, a poor performance was observed
197 in the “indications” (1 [6.7%]) and “toxicities” (0 [0%]) domains (Figure 2). In terms of accuracy,
198 the numbers of answers deemed fully correct were 31 (58.5%), 43 (75.4%), and 7 (43.8%) for
199 ChatGPT-3.5, ChatGPT-4 and Google Bard, respectively ($p = 0.03$). Furthermore, regarding
200 relevancy, the numbers of responses deemed highly relevant were 41 (77.4%), 41 (71.9%), and 7
201 (43.8%) for ChatGPT-3.5, ChatGPT-4 and Google Bard, respectively ($p = 0.04$). Readability was
202 deemed optimal across all three LLMs. However, the numbers of highly readable answers were
203 greater for ChatGPT-3.5 and ChatGPT-4 (52 [98.1%] and 57 [100%]) compared to Google Bard
204 (14 [87.5%]) ($p = 0.02$) (Figure 3). The median numbers of words and their corresponding ranges

205 for the responses provided by ChatGPT-3.5, ChatGPT-4, and Google Bard were 297 (197 - 404),
206 276 (139 - 395), and 290.5 (12 - 424), respectively ($p = 0.06$). Finally, the median numbers of
207 characters and their corresponding ranges were 1829 (1119 - 2470), 1589 (854 - 2233), and 1532
208 (75 - 2070), respectively ($p < 0.0001$).

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224 **4. DISCUSSION**

225 In recent decades, significant effort has been made to harness the potential of AI in medicine and
226 healthcare (14). Artificial intelligence can be defined as “the science and engineering of making
227 intelligent machines, especially intelligent computer programs” (15). It is composed of multiple
228 subfields, based on different algorithms and principles, including knowledge representation,
229 machine learning (ML), DL, and NLP (2,16). Specifically, NLP uses computational language and
230 DL to enable computers to understand text in the same way as humans (2). Recent progress in NLP
231 has led to major breakthroughs in the field of generative AI, as evidenced by the advent of LLMs
232 (3). These can recognize, summarize and generate novel content using statistical connections
233 between letters and words. Indeed, LLMs can also be considered as “few shot learners” due to
234 their ability to readily adapt to new domains with few information after being trained (17).

235
236 Over the last year, the release of ChatGPT (10) has attracted considerable attention, which only
237 increased following the release of other LLMs such as Google Bard (11), Bing AI (18) and,
238 Perplexity (19). The remarkable adaptability of these AI-based technologies to a broad and
239 extensive range of disciplines was immediately apparent following their introduction (20). This is
240 also evidenced by the rapid publication of large numbers of studies designed to investigate their
241 role in multiple and diffuse fields, including medicine and healthcare. Initial data have
242 demonstrated LLMs to be highly applicable to the field of cancer care, especially in providing
243 information about the screening and/or management of specific solid tumors (7). However, to the
244 authors’ knowledge, their potential role in the field of IO has not yet been investigated, despite the
245 rapidly expanding knowledge in all the aspects of IO (basic, translational, and clinical research)
246 and the large number of cancer patients currently receiving immunotherapy (8,9).

247 Therefore, we performed a cross-sectional study aimed for the first time at assessing the potential
248 of three prominent LLMs in answering questions about the field of IO. Our results demonstrated
249 that ChatGPT-4 and ChatGPT-3.5 were able to answer most of the IO-related questions with
250 excellent accuracy and relevance. In contrast, the performance of Google Bard was comparatively
251 poorer, as shown by a lower number of both answered questions and the reproducibility/accuracy
252 of these responses, compared to the other two LLMs. All three LLMs were able to provide highly
253 readable responses, highlighting the power of these generative AI technologies in providing
254 human-readable text. ChatGPT (both v3.5 and v4) clearly demonstrated their potential as a “virtual
255 assistant” for both clinicians and patients or caregivers. ChatGPT (both v3.5 and, especially, v4)
256 has also demonstrated remarkable acumen in both diagnosing and providing management plans
257 for IO toxicities. It has also proved highly effective in suggesting evidence-based and licensed
258 indications for IO therapy, either alone or in combination. Additionally, it has demonstrable
259 efficacy in providing background information on IO drug mechanisms and disease prognoses in
260 generally comprehensible text without excess jargon, albeit often with a lack of sources and broken
261 or inaccurate references.

262
263 However, the results of this study also highlight the differing performance of various LLMs across
264 topics and specific tasks (Table 3), as this demonstrates significant variability. In our study,
265 ChatGPT is demonstrated to be a powerful tool when applied to the field of IO, particularly in
266 comparison to Google Bard. Similar results were also reported in another recently published study
267 assessing these three LLMs in a different cancer-related topic. Specifically, Rahsepar et al.
268 reported the results of a study investigating the ability of ChatGPT-3.5, ChatGPT-4 and Google
269 Bard in answering questions related to lung cancer screening and prevention (21). As in our study,

270 ChatGPT achieved a superior performance to Google Bard. However, the available evidence
271 suggests that the LLM developed by OpenAI is not always accurate, as shown by the results of
272 other studies investigating medical/healthcare topics other than cancer (Table 3). In the studies
273 published by Seth et al., Zúñiga Salazar et al. and Dhanvijay et al., Google Bard performed better
274 in comparison to ChatGPT in non-cancer domains, likely clarifying a potential role for this LLM
275 (22–24). Furthermore, the results of the study by Al-Ashwal et al. showed a better performance
276 for Bing AI in answering questions related to drug-drug interactions in comparison to the other
277 LLMs (25). Therefore, it is essential to compare the performance of different LLMs since their
278 abilities may vary based on both task and domain.

279
280 In addition, despite the promising results of our study and its unequivocal efficacy in synthesizing
281 and evaluating textual data, the potential of ChatGPT for error and hallucination remains (26). The
282 occurrence of “hallucinations” is one of the greatest obstacles to the routine clinical application of
283 LLMs. While potentially tolerable in other domains, this is a critical issue in medicine and the
284 biomedical sciences due to its potential to directly impact patient care. In addition, it must be noted
285 that the datasets on which these models were trained were: (i) confidential and proprietary (thus
286 impossible to assess for data quality or bias), (ii) not specifically selected *ab initio* for addressing
287 biomedical issues and (iii) only valid up to September 2021 (thus lacking up to date information
288 – a major issue in so rapidly evolving a field as medicine in general and IO in particular) (10,27).
289 Therefore, expert assessment of LLMs’ output remains a prerequisite for their clinical use.

290
291 Open-source LLMs trained on specific biomedical datasets in order to accomplish pre-specified
292 tasks offer a potential solution and alternative paradigm. BioGPT, a cutting-edge LLM with a user-

293 friendly interface developed for the biomedical field, represents an excellent example of this (28).
294 BioGPT shares the same architecture as OpenAI's GPT models but was trained on information
295 derived from the biomedical literature. It has demonstrated excellent performance in several tasks,
296 including text generation and categorization, due to its extensive pre-training on massive
297 biomedical datasets (28). Further studies to investigate the utility and performance of LLMs
298 developed on biomedical data, with comparison to those LLMs presently available, are, thus,
299 required.

300

301 **4.1 Limitations**

302 Our study has some limitations that need to be mentioned. Firstly, we have focused only on three
303 prominent LLMs, excluding other LLMs including BingAI and Perplexity. At the time of the
304 design of this study, ChatGPT and Google Bard were the most investigated LLMs and, thus, we
305 elected to focus on them. However, recent evidence has shown the potential of BingAI in the
306 biomedical field. Therefore, our results do not represent the entire spectrum of LLMs available
307 and further assessment of other LLMs in the field of IO is essential. Secondly, the rating process
308 of the answers was made by only two reviewers. However, while a third reviewer was available to
309 resolve any conflicts which arose, this proved unnecessary as a strong to near perfect agreement
310 was demonstrated between the two reviewers Finally, the number of open-ended questions
311 included was relatively small, which may have impacted the analysis, particularly for domain-
312 specific performance.

313

314

315

316 **5. CONCLUSION**

317 ChatGPT-3.5 and ChatGPT-4 have demonstrated significant and clinically meaningful utility as
318 decision- and research-aids in various subfields of IO, while Google Bard demonstrated significant
319 limitations, especially in comparison to ChatGPT. However, the risk of inaccurate or incomplete
320 responses was evident in all LLMs, highlighting the importance of an expert-driven verification of
321 the information provided by these technologies. Finally, despite their potential to positively impact
322 the field of medicine and healthcare, this study reinforced the significance of a human evaluation
323 of LLMs in order to create reliable tools for clinical use.

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339 **Conflicts of Interest**

340 Authors declare no conflict of interest.

341

342 **CRedit Roles**

343 Conceptualization (GMI, DBC, CSF); Formal analysis (GMI and HCW); Investigation (GMI and
344 DBC); Methodology (GMI and DBC); Visualization (GMI and CSF); Writing - Original Draft
345 (GMI, DBC, HCW); Writing - Review & Editing (FK, JG, CSF); Supervision (CSF). All authors
346 accepted the final draft of the manuscript.

347

348 **Data Availability Statement**

349 The data underlying this article are available in the article and in its online supplementary material.

350

351 **Funding**

352 This work was supported by the Intramural Research Program, National Institutes of Health,
353 National Cancer Institute, Center for Cancer Research. The interpretation and reporting of these
354 data are the sole responsibility of the authors.

355

356

357

358

359

360

361

362 **REFERENCES**

363

- 364 1. IBM. What is generative AI? [Internet]. 2021 [cited 2023 Oct 13]. Available from:
365 <https://research.ibm.com/blog/what-is-generative-AI>
- 366 2. IBM. What is Natural Language Processing? | IBM [Internet]. [cited 2023 Oct 15]. Available
367 from: <https://www.ibm.com/topics/natural-language-processing>
- 368 3. Birhane A, Kasirzadeh A, Leslie D, Wachter S. Science in the age of large language models.
369 Nat Rev Phys [Internet]. 2023 [cited 2023 Oct 13];5(5). Available from:
370 <https://ora.ox.ac.uk/objects/uuid:9eac0305-0a9a-4e44-95f2-c67ee9eae15c>
- 371 4. Ayers JW, Poliak A, Dredze M, Leas EC, Zhu Z, Kelley JB, et al. Comparing Physician and
372 Artificial Intelligence Chatbot Responses to Patient Questions Posted to a Public Social
373 Media Forum. JAMA Intern Med. 2023 Jun 1;183(6):589–96.
- 374 5. Risk A, Petersen C. Health information on the internet: quality issues and international
375 initiatives. JAMA. 2002 May 22;287(20):2713–5.
- 376 6. Sallam M. ChatGPT Utility in Healthcare Education, Research, and Practice: Systematic
377 Review on the Promising Perspectives and Valid Concerns. Healthc Basel Switz. 2023 Mar
378 19;11(6):887.
- 379 7. Iannantuono GM, Bracken-Clarke D, Floudas CS, Roselli M, Gulley JL, Karzai F.
380 Applications of large language models in cancer care: current evidence and future
381 perspectives. Front Oncol. 2023 Sep 4;13:1268915.
- 382 8. Johnson DB, Nebhan CA, Moslehi JJ, Balko JM. Immune-checkpoint inhibitors: long-term
383 implications of toxicity. Nat Rev Clin Oncol. 2022 Apr;19(4):254–67.
- 384 9. Darvin P, Toor SM, Sasidharan Nair V, Elkord E. Immune checkpoint inhibitors: recent
385 progress and potential biomarkers. Exp Mol Med. 2018 Dec 13;50(12):1–11.
- 386 10. OpenAI. What is ChatGPT? [Internet]. [cited 2023 Oct 13]. Available from:
387 <https://help.openai.com/en/articles/6783457-what-is-chatgpt>
- 388 11. Google. Try Bard, an AI experiment by Google [Internet]. [cited 2023 Oct 13]. Available
389 from: <https://bard.google.com>
- 390 12. McHugh ML. Interrater reliability: the kappa statistic. Biochem Medica. 2012;22(3):276–82.
- 391 13. von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP, et al. The
392 Strengthening the Reporting of Observational Studies in Epidemiology (STROBE)
393 statement: guidelines for reporting observational studies. J Clin Epidemiol. 2008
394 Apr;61(4):344–9.

- 395 14. Haug CJ, Drazen JM. Artificial Intelligence and Machine Learning in Clinical Medicine,
396 2023. *N Engl J Med*. 2023 Mar 30;388(13):1201–8.
- 397 15. McCarthy J. What Is Artificial Intelligence?
- 398 16. IBM. AI vs. Machine Learning vs. Deep Learning vs. Neural Networks: What’s the
399 difference? [Internet]. 2023 [cited 2023 Oct 16]. Available from:
400 <https://www.ibm.com/blog/ai-vs-machine-learning-vs-deep-learning-vs-neural-networks/>
- 401 17. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, et al. Language Models are
402 Few-Shot Learners [Internet]. arXiv; 2020 [cited 2023 Oct 16]. Available from:
403 <http://arxiv.org/abs/2005.14165>
- 404 18. Microsoft. Bing AI [Internet]. [cited 2023 Oct 17]. Available from:
405 <https://www.bing.com:9943/search?showconv=1&q=bing AI&sf=codex3p&form=MA13FV>
- 406 19. Perplexity AI. Perplexity [Internet]. [cited 2023 Oct 17]. Available from:
407 <https://www.perplexity.ai/>
- 408 20. Thorp HH. ChatGPT is fun, but not an author. *Science*. 2023 Jan 27;379(6630):313.
- 409 21. Rahsepar AA, Tavakoli N, Kim GHJ, Hassani C, Abtin F, Bedayat A. How AI Responds to
410 Common Lung Cancer Questions: ChatGPT vs Google Bard. *Radiology*. 2023
411 Jun;307(5):e230922.
- 412 22. Seth I, Lim B, Xie Y, Cevik J, Rozen WM, Ross RJ, et al. Comparing the Efficacy of Large
413 Language Models ChatGPT, BARD, and Bing AI in Providing Information on Rhinoplasty:
414 An Observational Study. *Aesthetic Surg J Open Forum*. 2023;5:ojad084.
- 415 23. Zúñiga Salazar G, Zúñiga D, Vindel CL, Yoong AM, Hincapie S, Zúñiga AB, et al. Efficacy
416 of AI Chats to Determine an Emergency: A Comparison Between OpenAI’s ChatGPT,
417 Google Bard, and Microsoft Bing AI Chat. *Cureus*. 2023 Sep;15(9):e45473.
- 418 24. Dhanvijay AKD, Pinjar MJ, Dhokane N, Sorte SR, Kumari A, Mondal H. Performance of
419 Large Language Models (ChatGPT, Bing Search, and Google Bard) in Solving Case
420 Vignettes in Physiology. *Cureus*. 2023 Aug;15(8):e42972.
- 421 25. Al-Ashwal FY, Zawiah M, Gharaibeh L, Abu-Farha R, Bitar AN. Evaluating the Sensitivity,
422 Specificity, and Accuracy of ChatGPT-3.5, ChatGPT-4, Bing AI, and Bard Against
423 Conventional Drug-Drug Interactions Clinical Tools. *Drug Healthc Patient Saf*.
424 2023;15:137–47.
- 425 26. Singhal K, Azizi S, Tu T, Mahdavi SS, Wei J, Chung HW, et al. Large language models
426 encode clinical knowledge. *Nature*. 2023 Aug;620(7972):172–80.
- 427 27. van Dis EAM, Bollen J, Zuidema W, van Rooij R, Bockting CL. ChatGPT: five priorities for
428 research. *Nature*. 2023 Feb;614(7947):224–6.

- 429 28. Luo R, Sun L, Xia Y, Qin T, Zhang S, Poon H, et al. BioGPT: generative pre-trained
430 transformer for biomedical text generation and mining. *Brief Bioinform.* 2022 Nov
431 19;23(6):bbac409.
- 432 29. Koga S, Martin NB, Dickson DW. Evaluating the performance of large language models:
433 ChatGPT and Google Bard in generating differential diagnoses in clinicopathological
434 conferences of neurodegenerative disorders. *Brain Pathol Zurich Switz.* 2023 Aug 8;e13207.
- 435 30. Kumari A, Kumari A, Singh A, Singh SK, Juhi A, Dhanvijay AKD, et al. Large Language
436 Models in Hematology Case Solving: A Comparative Study of ChatGPT-3.5, Google Bard,
437 and Microsoft Bing. *Cureus.* 2023 Aug;15(8):e43861.
- 438 31. Lim ZW, Pushpanathan K, Yew SME, Lai Y, Sun CH, Lam JSH, et al. Benchmarking large
439 language models' performances for myopia care: a comparative analysis of ChatGPT-3.5,
440 ChatGPT-4.0, and Google Bard. *EBioMedicine.* 2023 Sep;95:104770.
- 441 32. Meo SA, Al-Khlaiwi T, AbuKhalaf AA, Meo AS, Klonoff DC. The Scientific Knowledge of
442 Bard and ChatGPT in Endocrinology, Diabetes, and Diabetes Technology: Multiple-Choice
443 Questions Examination-Based Performance. *J Diabetes Sci Technol.* 2023 Oct
444 5;19322968231203987.
- 445 33. Toyama Y, Harigai A, Abe M, Nagano M, Kawabata M, Seki Y, et al. Performance
446 evaluation of ChatGPT, GPT-4, and Bard on the official board examination of the Japan
447 Radiology Society. *Jpn J Radiol.* 2023 Oct 4;
- 448 34. Waisberg E, Ong J, Masalkhi M, Zaman N, Sarker P, Lee AG, et al. Google's AI chatbot
449 "Bard": a side-by-side comparison with ChatGPT and its utilization in ophthalmology. *Eye*
450 *Lond Engl.* 2023 Sep 28;

451

452

453

454

455

456

457

458

459

460

461 **FIGURE LEGENDS**

462 **Figure 1:** Flowchart of the rating process for each triplet of responses.

463

464 **Figure 2:** Spot matrix of the percentages of the answered questions [Blue] and reproducible
465 responses [Orange] for each LLM. Color volume is directly proportional to percentage with the
466 outer black circle representing 100%. Corresponding numeric data are available in Suppl. Mat. D.

467

468 **Figure 3:** Bar plot of the results (accuracy, readability, and relevance) for all three LLMs. This
469 plot was based only on the questions evaluable for accuracy, readability, and relevance.
470 Corresponding numeric data are available in Suppl. Mat. D.

Table 1. Definitions of the outcomes

Outcomes	Definitions	Score
Answer Returned	<i>The ability of LLM to return a meaningful answer to each instance of the question submitted, rather than returning an error or declining to return an answer, independent of the accuracy of this response.</i>	Recorded as Boolean True/False
Reproducibility	<i>The ability of LLM to return a generally similar series of answers across the three separate queries with no fundamental differences or inconsistencies between these three answers.</i>	Recorded as Boolean True/False
Accuracy	<i>The ability of LLM to provide accurate and correct information addressing the question asked and returning all major or critical points required in such an answer. Response <u>not</u> adversely marked for extraneous or irrelevant information here – as long as this information was correct.</i>	Recorded numerically from 1 to 3
Readability	<i>The ability of LLM to return comprehensible and coherent natural language text in English, including appropriate syntax, formatting, and punctuation, independent of the accuracy of this response.</i>	Recorded numerically from 1 to 3
Relevance	<i>The ability of LLM to return information that was relevant and specific to the question asked or immediately adjacent topics without extraneous, unrequested, or tangential information. Accuracy was not specifically assessed here, though the result was adversely marked if the response included immaterial information while neglecting to address the specific question asked.</i>	Recorded numerically from 1 to 3

Note: for scoring of Relevance, the answer returned was not adversely marked for any included disclaimers to the effect that the LLM cannot provide medical advice and any such advice should be sought from a clinician or that anyone with a cancer diagnosis and/or receiving systemic therapy with potential toxicity should contact their treating clinician/s. This was deemed to represent appropriate and medically sound advice and not to be irrelevant or extraneous material.

Table 2. Definitions of the scoring system

	Score		
	1	2	3
Accuracy	<i>Fundamentally inaccurate or incorrect information, including critical errors, omissions and/or entirely incorrect treatment advice.</i>	<i>Partially correct and accurate information, including non-critical errors and/or omitting relevant information or failing to provide specific guideline advice.</i>	<i>Fully accurate and correct information, answering the specific question asked with no significant errors or omissions.</i>
Relevance*	<i>Irrelevant and/or entirely tangential material, not addressing the specific question asked.</i>	<i>Generally relevant material though including significant extraneous and/or tangential information.</i>	<i>Relevant and focused information directly addressing the question asked, including an appropriate expansion on the relevant topic.</i>
Readability	<i>Incoherent, unintelligible and/or garbled text, +/- severely misformatted and/or oxymoronic material resulting in compromised legibility.</i>	<i>Generally coherent and intelligible material with significant formatting and/or parsing errors.</i>	<i>Fully coherent, well-parsed and constructed material, easily and clearly intelligible.</i>

*Note: Inclusion of a disclaimer that the answer was provided by an AI/LLM and cannot be taken as medical advice and/or that any information or questions should also be addressed to a qualified medical practitioner was not scored negatively – as this represents a legitimate and appropriate legal disclaimer.

Table 3. List of studies investigating the utility of ChatGPT and Google Bard across various contexts of medicine and healthcare.

First Author	Year of Publication	LLMs	Domain	Questions (n)	Reviewers (n)
Al-Ashwal FY (25)	2023	ChatGPT - Google Bard - Bing AI	Drug-drug interactions	225 [OE]	NA
Dhanvijay AK (24)	2023	ChatGPT - Google Bard - Bing AI	Physiology	77 [OE]	2
Seth I (22)	2023	ChatGPT - Google Bard - Bing AI	Rhinoplasty	6 [OE]	3
Koga S (29)	2023	ChatGPT - Google Bard	Neurodegenerative disorder	25 [OE]	NA
Kumari A (30)	2023	ChatGPT - Google Bard	Hematology	50 [OE]	3
Lim ZW (31)	2023	ChatGPT - Google Bard	Myopia	31 [OE]	3
Meo SA (32)	2023	ChatGPT - Google Bard	Endocrinology, diabetes, and diabetes technology	100 [MC]	-
Toyama Y (33)	2023	ChatGPT - Google Bard	Radiology	103 [MC]	3
Waisberg E (34)	2023	ChatGPT - Google Bard	Ophthalmology	NA	4
Zuniga Salazar G (23)	2023	ChatGPT - Google Bard - Bing AI	Emergency	176 [OE]	NA

Abbreviations: Multiple choice (MC); Not available (NA); Open-ended (OE).

Figure 1

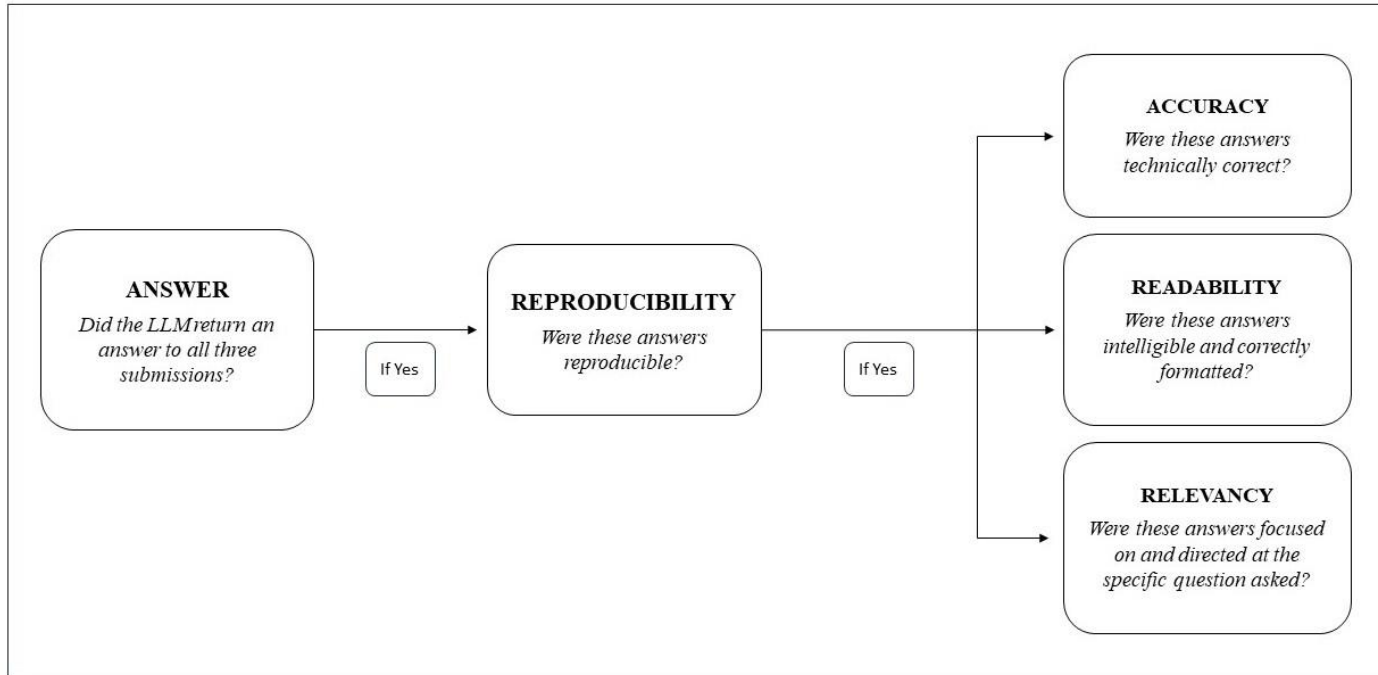


Figure 2

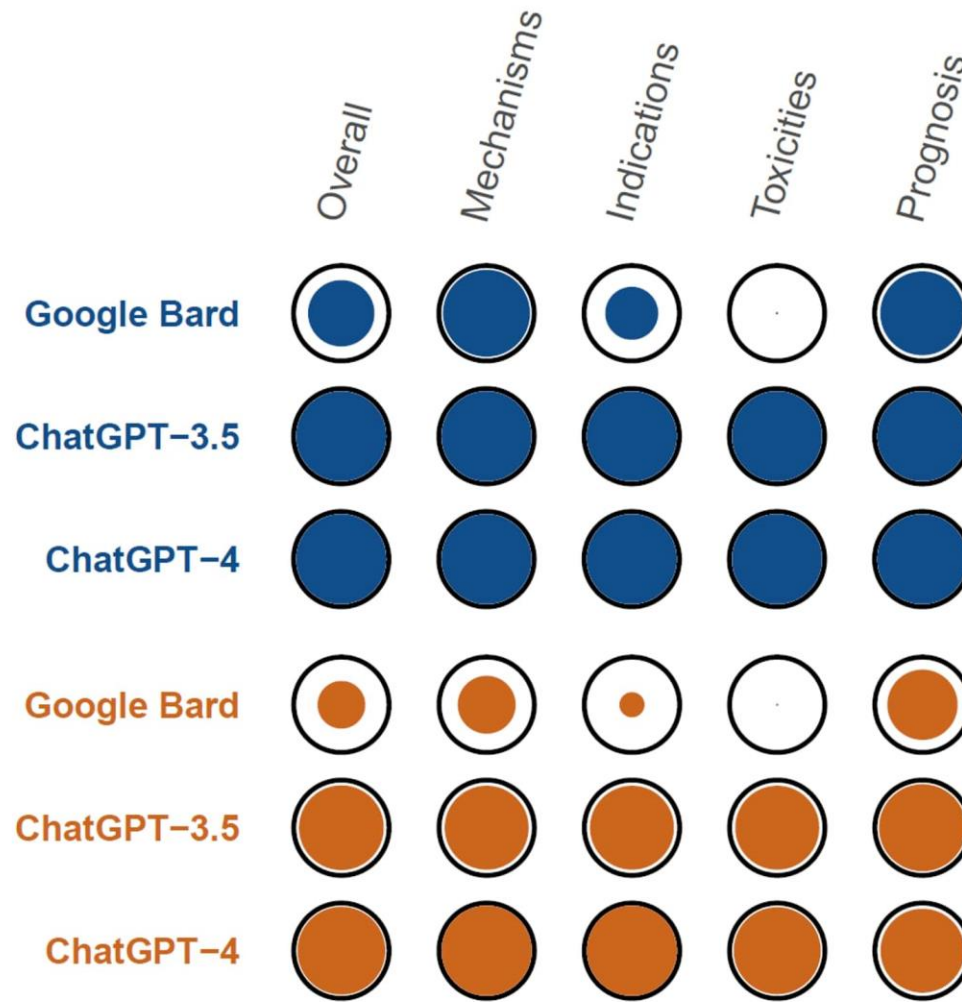


Figure 3

