

Evaluation of Information Quality and Readability of Artificial Intelligence–Powered Chatbots in Systemic Isotretinoin Use

© Huriye Aybüke Koç¹, © Elif Özenir¹, © Cansu Altınöz Güney²

¹Department of Dermatology and Venereology, Giresun University Faculty of Medicine, Giresun, Türkiye

²Clinic of Dermatology and Venereology, Dinar State Hospital, Afyonkarahisar, Türkiye

Abstract

Aim: This study aimed to evaluate and compare the readability and quality of information in responses generated by artificial intelligence (AI) models to patients' frequently asked questions about systemic isotretinoin, a medication commonly prescribed in dermatology.

Materials and Methods: Thirty-four frequently asked questions from patients using isotretinoin were prepared by a team of dermatology specialists. These questions were posed to three AI-based text-generation tools (ChatGPT, Gemini 2.0, and Copilot), and the responses were analyzed. The resulting texts were compared in terms of readability levels [Flesch Reading Ease score (FRES), Flesch-Kincaid Grade Level (FKGL), Simple Measure of Gobbledygook (SMOG), Gunning Fog index (GFOG), Coleman-Liau index (CLI), and Automated Readability index (ARI)], sentence lengths, and content quality, which was evaluated by dermatologists.

Results: None of the AI models achieved the optimal readability threshold ($FRES \geq 60$). Readability metrics differed significantly among models. Gemini produced responses that were significantly less readable and more complex than those produced by ChatGPT and Copilot across all readability indices, including FRES, FKGL, SMOG, GFOG, CLI, and ARI; post-hoc analyses confirmed differences between Gemini and the other models. Sentence counts also differed significantly, with Gemini generating longer responses than Copilot. In contrast, Likert-based quality scores and response appropriateness were comparable across models, with no statistically significant differences observed.

Conclusion: This study demonstrates that AI models produce academic responses that are difficult for those unfamiliar with medical terminology to understand, and can generate outputs with variable readability in health-related content. These findings highlight the need for careful evaluation of AI-based content for use in healthcare.

Keywords: Artificial intelligence, isotretinoin, readability, quality

INTRODUCTION

Chatbots are computer programs that can understand and respond to speech and text in a human-like manner using various algorithms. Large language models are used to simulate human conversation. Many companies are using this technology to develop their own chatbots.^{1,2} A major milestone in this field was the introduction of ChatGPT in 2022. These chatbots have various applications, including serving as dialogue systems, providing language translation,

and generating content. Alongside ChatGPT, other chatbots utilizing large language models have been launched.^{1,3} Microsoft Copilot is another chatbot with different functionalities. Unlike ChatGPT, it can search the internet and update its knowledge base.⁴ Google Gemini, developed in collaboration with Google teams, integrates different types of information: text, code, audio, images, and video to serve as a writing, planning, and learning assistant.¹ Besides these,

Address for correspondence: Elif Özenir,

Department of Dermatology and Venereology, Giresun University Faculty of Medicine, Giresun, Türkiye

Email: elfozenir@gmail.com

ORCID ID: 0009-0002-0000-6384

Submission: 17-Dec-2025

Epub: 06-Mar-2026

Acceptance: 18-Feb-2026

Access this article online

Quick Response Code:



Website:

www.turkjdermatol.com

DOI:

10.4274/tjd.galenos.2026.53244



This is an open access journal, and articles are distributed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License, which allows others to remix, tweak, and build upon the work non-commercially, as long as appropriate credit is given.

How to cite this article: Koç HA, Özenir E, Altınöz Güney C. Evaluation of information quality and readability of artificial intelligence–powered chatbots in systemic isotretinoin use. *Turk J Dermatol.* [Epub Ahead of Print]

there are currently over 5,100 chatbots, approximately 100 of which are used in healthcare applications.⁵ Chatbots offer advantages such as improving diagnostic accuracy, supporting personalized treatment plans, facilitating the translation of the latest medical literature into clinical practice, and contributing to patient education, thereby enhancing healthcare services.^{5,6}

Many dermatological diseases have a chronic course and require long-term follow-up and treatment. However, patients do not always have easy access to a dermatologist. Therefore, patients increasingly turn to online platforms, including social media and artificial intelligence (AI) chatbots, to obtain information about their diseases and the medications they use. Although these assistive tools are well-designed, concerns remain regarding the accuracy, currency, and reliability of the medical information they provide, and regarding the transparency with which user data are handled. To address these concerns, various studies on AI chatbots have been conducted in different specialties.⁶⁻⁹ In this study, we investigated the readability and reliability of chatbot responses to the most frequently asked questions about isotretinoin, a medication commonly prescribed in dermatology outpatient clinics.

MATERIALS AND METHODS

Chatbots

Chatbots were selected considering factors such as user fees, login requirements, and inspiration from previous similar studies. The chatbots selected for the study and their versions were ChatGPT 4, Google Gemini 2.0, and Microsoft 365 Copilot; the readability and quality of the responses from these chatbots were evaluated. In the remainder of this paper, these chatbots are referred to as ChatGPT, Gemini, and Copilot. The chatbots were accessed using a personal computer (MacBook Air M2) connected to a home broadband connection. Data were collected between July 1 and July 5, 2025.

Questions

The questions about isotretinoin, the most frequently prescribed active ingredient in dermatologic practice, were prepared by expert dermatologists. Of the prepared questions, 34 were selected. Each question was asked individually to the chatbots, and the responses were recorded in separate documents for review and analysis. Responses were examined by two dermatologists specializing in the field. All items were evaluated jointly by two dermatologists. Both the 5-point Likert scale ratings and the appropriateness classifications (appropriate, incomplete, and inappropriate) were determined by discussion and assigned by consensus. Since the ratings were not performed independently, interrater reliability

analysis was not applicable. The British Association of Dermatologists' guidelines were taken as the criterion for response accuracy (Supplementary).¹⁰ The Likert scale developed by Kumari et al.,¹¹ shown in Table 1, was used to evaluate the accuracy of responses. Furthermore, chatbot responses were classified into three categories: "appropriate", "incomplete", and "inappropriate". An appropriate response was defined as accurate, complete, and consistent with what an expert would advise a patient in the same situation; an inappropriate response was defined as inconsistent with expert opinion or containing incorrect information; and an incomplete response was defined as correct and relevant but lacking sufficient detail. Prior to each question, the chatbot sessions were reset.

Readability Analysis

After the accuracy of the responses had been verified by two independent dermatologists, a readability analysis was performed for each response. The following measures were used: Flesch Reading Ease score (FRES), Flesch-Kincaid Grade level (FKGL), Simple Measure of Gobbledygook (SMOG), Gunning Fog index (GFOG), Coleman-Liau index (CLI), and Automated Readability index (ARI). In addition, sentence lengths for each response were compared.

Flesch Reading Ease Score: It is a measure of text readability calculated based on the average number of words per sentence and the average number of syllables per word. Scores range from 0 to 100, with higher scores indicating easier readability. Scores between 70 and 80 correspond to approximately an eighth-grade reading level.¹²

Flesch-Kincaid Grade Level: A readability measure that determines reading difficulty according to the United States school grade level. Scores range from 0 to 18, with higher scores indicating greater difficulty. Scores above 12 indicate that the text is written in an academic style.¹³

Simple Measure of Gobbledygook: Designed to assess the appropriateness of a text for the reader's age. It counts ten sentences from the beginning, middle, and end of the text to determine the level. It counts words with three or more syllables across 30 sentences. The syllable counts are then converted to a corresponding reading-grade level.¹⁴

Gunning Fog Index: This metric determines how difficult a text is to read based on sentence and word length. Scores range from 1 to 18, with each score corresponding to the number of years of education needed to understand the text. To facilitate comprehension by the general public, an average score of 8 is recommended. Scores of 17 and above are considered primarily understandable to individuals with postgraduate education.¹⁵

Coleman-Liau Index: CLI is a readability assessment that measures how challenging a text is and helps determine its appropriate grade level. It is commonly used in the USA and several other countries. Unlike many other grade-level estimators, CLI relies on the number of characters per word rather than on the number of syllables per word.¹⁶

Automated Readability Index: This measure estimates the number of years of education required to understand a text on first reading. It takes into account the average number of characters per word and the average number of words per sentence. ARI uses a specific formula to determine the grade level of the text.¹⁷

Statistical Analysis

Data were analyzed using SPSS for Windows, version 26.0. Descriptive statistics—including mean, standard deviation, median, and percentage distributions—were used to summarize the data. The Kruskal-Wallis test was applied to compare continuous and ordinal variables, such as readability scores and Likert-scale ratings, across the three AI models. For categorical variables, such as response appropriateness, the chi-square test was employed. A *P*-value of less than 0.05 was considered statistically significant.

RESULTS

In our study, 34 questions regarding systemic isotretinoin use were asked of each of three AI chatbots. While a FRES of ≥ 60 is required for optimal readability, none of the examined models reached this threshold. With the exception of FRES, none of the models achieved an acceptable readability level on the other five measures; all models scored well above the

thresholds, indicating generally low readability of the content (Table 1).

Readability scores differed significantly among the three AI models across all objective indices (FRES: *P* = 0.013; FKGL: *P* = 0.002; SMOG: *P* < 0.001; GFOG: *P* = 0.006; CLI: *P* = 0.003; ARI: *P* < 0.001). Dunn-Bonferroni post-hoc analyses indicated that these differences were consistently driven by Gemini. Compared with both ChatGPT and Copilot, Gemini produced responses that were significantly less readable, as reflected by lower FRES scores (vs. ChatGPT: *P* = 0.011; vs. Copilot: *P* = 0.010) and higher grade-level indices (FKGL: *P* = 0.002 for both comparisons; SMOG: *P* = 0.001 and *P* < 0.001; GFOG: *P* = 0.005 and *P* = 0.007; CLI: *P* = 0.002 and *P* = 0.006; ARI: *P* < 0.001 for both). No significant differences were observed between ChatGPT and Copilot on any readability metric (Table 1).

Sentence counts also differed significantly (*P* = 0.009); Gemini generated longer responses than Copilot (*P* = 0.003); other pairwise comparisons were not significant.

Likert scale ratings evaluating the quality of responses were similar across models (median = 4), with no statistically significant difference observed (*P* = 0.259). Similarly, the distribution of response appropriateness did not differ significantly among models (*P* = 0.701), and most responses were rated as appropriate (Table 1).

DISCUSSION

Today, increasing and unmet healthcare demands are leading individuals to seek information from alternative sources. Among these, online tools are the most frequently used due

Table 1. Comparison of the relevance and readability of responses from three AI models to questions asked by isotretinoin users

	ChatGPT (Mean ± SD)	Gemini (Mean ± SD)	Copilot (Mean ± SD)	<i>P</i> -value*
Flesch Reading Ease score	-2.66 ± 15.08	-11 ± 12.76	-2.71 ± 13.26	0.013
Flesch-Kincaid Grade level	16.91 ± 2.33	18.67 ± 1.37	16.97 ± 2.16	0.002
Simple Measure of Gobbledygook	12.7 ± 1.58	13.97 ± 1	12.62 ± 1.85	< 0.001
Gunning Fog index	18.85 ± 3.17	20.77 ± 2.09	18.95 ± 2.72	0.006
Coleman-Liau index	19.12 ± 2.49	20.8 ± 1.75	19.44 ± 1.89	0.003
Automated Readability index	12.73 ± 2.51	14.85 ± 1.41	12.8 ± 2.12	< 0.001
Sentence count	13.79 ± 5	16.88 ± 6.71	12.18 ± 3.04	0.009
Likert scale (1–5)	4.35 (median:4)	4.18 (median:4)	4.09 (median:4)	0.259
Appropriateness**	26 (76.5%)	24 (70.6%)	24 (70.6%)	0.701
Appropriate	8 (23.5%)	8 (23.5%)	9 (26.5%)	
Incomplete inappropriate	0 (0.0%)	2 (5.9%)	1 (2.9%)	
*Kruskal-Wallis test **chi-square test SD: Standard deviation, AI: Artificial intelligence				

to their accessibility. Previous studies have shown that a significant proportion of patients turn to online resources for information regarding their diseases and treatments.^{18,19}

Our current study provides insight into the performance and reliability of chatbot responses in the medical context. Chatbots are increasingly used across many fields, including medicine. Due to growing healthcare needs and various limitations, unmet demands are increasingly being addressed through tools such as online chatbots. Our study found that chatbots provided answers that varied in length and readability to the same questions. Even though we selected the questions and responses from a publicly accessible online guide, each chatbot generated content and response lengths that differed. While a FRES of ≥ 60 is required for optimal readability, none of the examined models reached this threshold.

When examining accuracy and appropriateness scores, we found that the three chatbots demonstrated similar performance. Similar to our findings, a previous study comparing ChatGPT and Google Bard on educational questions posed by patients with obstructive sleep apnea found the responses from both chatbots to be appropriate and accurate.²⁰ However, other studies in dermatology, hematology, neurosurgery, lung cancer, and urology have shown differing accuracy rankings among ChatGPT, Gemini, and Copilot.²¹ These varying results may be related to differences in the algorithms used by chatbots, the training data, which can vary by country, and updates made to chatbots over time. While accuracy rates differ across studies, one common finding in our study is that no chatbot achieved 100% accuracy.

Statistically significant differences in readability scores were most pronounced between ChatGPT and Gemini and between Gemini and Copilot. Across readability scales, Gemini consistently received significantly higher scores than the other two chatbots, indicating lower readability and increased response complexity. Examination of sentence lengths across all three AI chatbots revealed a significant difference between Gemini and Copilot. Gemini provided longer and more detailed answers with the highest number of sentences, while Copilot offered shorter answers with simpler sentence structures.

Across all readability indices (FRES, FKGL, GFOG, ARI, SMOG, and CLI), the responses generally corresponded to a university-level or higher reading difficulty. Negative FRES values, particularly pronounced in Gemini, were attributable to very long sentences and the use of multisyllabic terms, as reflected in the formula: $FRES = 206.835 - (1.015 \times \text{average words per sentence}) - (84.6 \times \text{average syllables per word})$. When the subtraction components exceed 206.835, negative scores occur, indicating exceptionally high textual complexity.¹²

These findings align with previous studies evaluating chatbot readability in lung cancer, radiology, urology, and chronic kidney disease contexts, all of which reported low readability levels.^{17,21-26} However, a discipline-specific study in urology has reported different readability outcomes.²⁷

Likert-scale ratings of response quality were similar across models, and no statistically significant differences were found. The quality of chatbot responses ranged from 81.8% to 87; none achieved the perfect score of 5 and therefore cannot be considered 100% reliable. This finding is consistent with previous studies comparing AI chatbots.^{1,21-23,28-30} The quality and accuracy observed in our study suggest that chatbots may be useful for providing relatively accurate information about diseases. Consequently, they may provide valuable assistance to individuals concerning systemic isotretinoin, one of dermatology's fundamental treatment options. Chatbots could educate patients about the use, side effects, and treatment process of systemic isotretinoin, and about when to seek professional medical support. Additionally, by answering simple questions, chatbots can support patients in managing their treatment, thereby reducing the workload on the healthcare system and enabling dermatologists to devote more time to complex and serious cases.

Study Limitations

This study has several limitations. First, only three chatbots were selected, excluding other online platforms accessible to patients. Although the accuracy of responses was evaluated jointly by two dermatologists, inter-rater agreement statistics were not reported. This limitation is acknowledged and could be addressed in future studies by including measures such as Cohen's kappa. The question list was limited and prepared based on the most frequently asked questions in dermatology outpatient clinics. In real life, patient queries may be more diverse and multifaceted.

CONCLUSION

This study demonstrated that responses generated by chatbots about systemic isotretinoin, which is one of the most frequently prescribed treatments in dermatology, have readability levels corresponding to university education and above, making them relatively difficult to read. The highest response quality reached 87%, and no chatbot provided answers with 100% quality. Although the selected questions focused on a specific drug, the responses were based on publicly available online guidelines. Moreover, individuals of diverse ages and educational backgrounds who initiate systemic isotretinoin treatment in dermatology clinics may consult AI chatbots for information.

The high readability of these responses could lead to misinterpretation of information. This indicates that people seeking information through AI chatbots might be misled, potentially resulting in unnecessary anxiety and increased demand for consultations with doctors, thereby placing an excessive burden on the healthcare system. Conversely, if readability is low for matters requiring urgent intervention, this may delay necessary treatment and negatively impact patients' health.

For healthcare professionals, higher readability levels may be advantageous, providing more detailed and informative content. In the future, it would be beneficial to program chatbots to generate responses tailored to different age groups and educational levels. This approach could make AI chatbots more accessible and effective for diverse populations. More comprehensive and large-scale studies are needed to explore this further.

Ethics

Ethics Committee Approval: Ethical approval was not obtained for the study as it involved publicly accessible data and did not include any patient-specific information.

Informed Consent: Not applicable.

Authorship Contributions

Surgical and Medical Practices: H.A.K., E.Ö., Concept: H.A.K., C.A.G., Design: H.A.K., C.A.G., Data Collection or Processing: E.Ö., Analysis or Interpretation: C.A.G., Literature Search: H.A.K., Writing: H.A.K., E.Ö.

Conflict of Interest: The authors declared that they have no conflict of interest.

Financial Disclosure: The authors declared that this study received no financial support.

REFERENCES

- Olszewski R, Watros K, Mańczak M, Owoc J, Jeziorski K, Brzeziński J. Assessing the response quality and readability of chatbots in cardiovascular health, oncology, and psoriasis: a comparative study. *Int J Med Inform.* 2024;190:105562.
- Bohr A, Memarzadeh K. The rise of artificial intelligence in healthcare applications. *Artificial Intelligence in Healthcare.* 2020:25–60.
- Nirala KK, Singh NK, Purani VS. A survey on providing customer and public administration based services using AI: chatbot. *Multimed Tools Appl.* 2022;81(16):22215-22246.
- Semeraro F, Gamberini L, Carmona F, Monsieurs KG. Clinical questions on advanced life support answered by artificial intelligence. A comparison between ChatGPT, Google Bard and Microsoft Copilot. *Resuscitation.* 2024;195:110114.
- Diamond C, Rundle CW, Albrecht JM, Nicholas MW. Chatbot utilization in dermatology: a potential amelioration to burnout in dermatology. *Dermatol Online J.* 2022;28(6).
- Yan S, Du D, Liu X, Dai Y, Kim MK, Zhou X, Wang L, Zhang L, Jiang X. Assessment of the reliability and clinical applicability of ChatGPT's responses to patients' common queries about rosacea. *Patient Prefer Adherence.* 2024;18:249-253.
- Musheyev D, Pan A, Loeb S, Kabarriti AE. How well do artificial intelligence Chatbots respond to the top search queries about urological malignancies? *Eur Urol.* 2024;85(1):13-16.
- P Pan A, Musheyev D, Bockelman D, Loeb S, Kabarriti AE. Assessment of artificial intelligence Chatbot responses to top searched queries about cancer. *JAMA Oncol.* 2023;9(10):1437-1440.
- Young JN, Ross O'Hagan, Poplasky D, Levoska MA, Gulati N, Ungar B, Ungar J. The utility of ChatGPT in generating patient-facing and clinical responses for melanoma. *J Am Acad Dermatol.* 2023;89(3):602-604.
- British Association of Dermatologists. Isotretinoin patient guide [Internet]. London: British Association of Dermatologists; [cited 2026 Feb 15]. Available from: <https://www.bad.org.uk/pils/isotretinoin>
- Kumari A, Kumari A, Singh A, Singh SK, Juhi A, Dhanvijay AKD, Pinjar MJ, Mondal H. Large language models in hematology case solving: a comparative study of ChatGPT-3.5, Google Bard, and Microsoft Bing. *Cureus.* 2023;15(8):e43861.
- Bellot P, Tavernier J. Flesch and dale-chall readability measures for INEX 2011 question-answering track. *Lecture Notes in Computer Science.* 2012;2011(7424):235-246.
- Gbedemah ZEE, Fuseini MN, Fordjuor SKEJ, Baisie-Nkrumah EJ, Beecham REM, Amisah-Arthur KN. Readability and quality of online information on sickle cell retinopathy for patients. *Am J Ophthalmol.* 2024;259:45-52.
- Dalillah N, Ismayanti F, Azzahra E, Kusmana S, Rahayu I. SMOG (Simple Measure of Goobledygook) readability index in selecting reading materials and reading literacy skills of primary school student. *Int J Elem Educ.* 2024;13(2):31-38.
- Marshall S, Hanish SJ, Baumann J, Groneck A, DeFroda S. A standardised method for improving patient education material readability for orthopaedic trauma patients. *Musculoskeletal Care.* 2024;22(1):e1869.
- Readable. The Coleman-Liau Readability Index [Internet]. London: Readable; [cited 2026 Feb 15]. Available from: <https://readable.com/readability/coleman-liau-readability-index/>
- Gencer A. Readability analysis of ChatGPT's responses on lung cancer. *Sci Rep.* 2024;14(1):17234.
- Potemkowski A, Broła W, Ratajczak A, Ratajczak M, Zaborski J, Jasińska E, Pokryszko-Dragan A, Gruszka E, Dubik-Jezierzańska M, Podlecka-Piętowska A, Nojszewska M, Gospodarczyk-Szot K, Stępień A, Gocyla-Dudar K, Maciągowska-Terela M, Wencel J, Kazmierski R, Kulakowska A, Kapica-Topczewska K, Pawelczak W, Bartosik-Psujek H. Internet usage by polish patients with multiple sclerosis: a multicenter questionnaire study. *Interact J Med Res.* 2019;8(1):e11146.
- Wong DK, Cheung MK. Online health information seeking and ehealth literacy among patients attending a primary care clinic in Hong Kong: a cross-sectional survey. *J Med Internet Res.* 2019;21(3):e10831.
- Cheong RCT, Unadkat S, Mcneillis V, Williamson A, Joseph J, Randhawa P, Andrews P, Paleri V. Artificial intelligence chatbots as sources of patient education material for obstructive sleep apnoea: ChatGPT versus Google Bard. *Eur Arch Otorhinolaryngol.* 2024;281(2):985-993.
- Aydın FO, Aksoy BK, Ceylan A, Akbaş YB, Duran Güler S, Varan G, Kepez Yıldız B. Kontakt lens kullanıcı desteğinin yapay zeka ile geliştirilmesi: ChatBot'larda doğruluk ve anlaşılabilirliğin değerlendirilmesi. *MN Oftalmoloji.* 2025;32(2):108-114
- Li H, Moon JT, Iyer D, Balthazar P, Krupinski EA, Bercu ZL, Newsome JM, Banerjee I, Gichoya JW, Trivedi HM. Decoding radiology reports: Potential application of OpenAI ChatGPT to enhance patient understanding of diagnostic reports. *Clin Imaging.* 2023;101:137-141.
- Cocci A, Pezzoli M, Lo Re M, Russo GI, Asmundo MG, Fode M, Cacciamani G, Cimino S, Minervini A, Durukan E. Quality of information and appropriateness of ChatGPT outputs for urology patients. *Prostate Cancer Prostatic Dis.* 2024;27(1):103-108.

24. Tepe M, Emekli E. Assessing the responses of large language models (ChatGPT-4, Gemini, and Microsoft Copilot) to frequently asked questions in breast imaging: a study on readability and accuracy. *Cureus*. 2024;16(5):e59960.
25. Acharya PC, Alba R, Krisanapan P, Acharya CM, Suppadungsuk S, Csongradi E, Mao MA, Craici IM, Miao J, Thongprayoon C, Cheungpasitporn W. AI-driven patient education in chronic kidney disease: evaluating chatbot responses against clinical guidelines. *Diseases*. 2024;12(8):185.
26. Eid K, Eid A, Wang D, Raiker RS, Chen S, Nguyen J. Optimizing ophthalmology patient education via ChatBot-generated materials: readability analysis of AI-generated patient education materials and The American Society of Ophthalmic Plastic and Reconstructive Surgery Patient Brochures. *Ophthalmic Plast Reconstr Surg*. 2024;40(2):212-216.
27. Eppler MB, Ganjavi C, Knudsen JE, Davis RJ, Ayo-Ajibola O, Desai A, Storino Ramacciotti L, Chen A, De Castro Abreu A, Desai MM, Gill IS, Cacciamani GE. Bridging the gap between urological research and patient understanding: the role of large language models in automated generation of Layperson's summaries. *Urol Pract*. 2023;10(5):436-443.
28. Podder I, Pipil N, Dhabal A, Mondal S, Pienyii, Mondal H. Evaluation of artificial intelligence-based Chatbot responses to common dermatological queries. *J Med J*. 2024;58(2):271-277.
29. Ali R, Tang OY, Connolly ID, Zadnik Sullivan PL, Shin JH, Fridley JS, Asaad WF, Cielo D, Oyelese AA, Doberstein CE, Gokaslan ZL, Telfeian AE. Performance of ChatGPT and GPT-4 on neurosurgery written board examinations. *Neurosurgery*. 2023;93(6):1353-1365.
30. Rahsepar AA, Tavakoli N, Kim GHJ, Hassani C, Abtin F, Bedayat A. How AI responds to common lung cancer questions: ChatGPT vs Google bard. *Radiology*. 2023;307(5):e230922.

Supplementary Link: <https://d2v96fxpocvxx.cloudfront.net/cf9d60d6-523c-458a-a2e6-78728d3ffbb0/content-images/1a0969cb-47d7-48e7-9bee-9473405a43b3.pdf>
