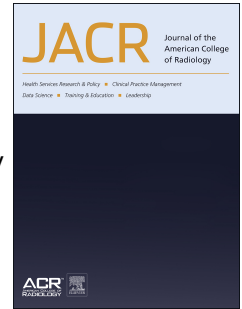


# Journal Pre-proof



Enhancing patient communication with Chat-GPT in radiology: evaluating the efficacy and readability of answers to common imaging-related questions

Emile B. Gordon, M.D., Alexander J. Towbin, M.D., Peter Wingrove, M.D., Umber Shafique, M.D., Brian Haas, M.D., Andrea B. Kitts, M.S., M.P.H., Jill Feldman, M.A., Alessandro Furlan, M.D.

PII: S1546-1440(23)00775-5

DOI: <https://doi.org/10.1016/j.jacr.2023.09.011>

Reference: JACR 6357

To appear in: *Journal of the American College of Radiology*

Received Date: 24 July 2023

Revised Date: 12 September 2023

Accepted Date: 20 September 2023

Please cite this article as: Gordon EB, Towbin AJ, Wingrove P, Shafique U, Haas B, Kitts AB, Feldman J, Furlan A, Enhancing patient communication with Chat-GPT in radiology: evaluating the efficacy and readability of answers to common imaging-related questions, *Journal of the American College of Radiology* (2023), doi: <https://doi.org/10.1016/j.jacr.2023.09.011>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier Inc. on behalf of American College of Radiology

# Enhancing patient communication with Chat-GPT in radiology: evaluating the efficacy and readability of answers to common imaging-related questions

Emile B. Gordon, M.D.<sup>a</sup>, Alexander J. Towbin, M.D.<sup>b</sup>, Peter Wingrove, M.D.<sup>a</sup>, Umber Shafique, M.D.<sup>c</sup>, Brian Haas, M.D.<sup>d</sup>, Andrea B. Kitts, M.S., M.P.H.<sup>e</sup>, Jill Feldman, M.A.<sup>f</sup>, Alessandro Furlan, M.D.<sup>a\*</sup>

<sup>a</sup> University of Pittsburgh Medical Center, Department of Radiology

<sup>b</sup> University of Cincinnati, Department of Radiology

<sup>c</sup> Indiana University School of Medicine, Department of Radiology

<sup>d</sup> University of California San Francisco, Department of Radiology and Biomedical Imaging,

<sup>e</sup> Rescue Lung Society

<sup>f</sup> EGFR Resisters

\***Corresponding author:** Emile B. Gordon, emile.gordon@duke.edu

**Data statement:** The author(s) declare(s) that they had full access to all of the data in this study, and the author(s) take(s) complete responsibility for the integrity of the data and the accuracy of the data analysis.

**Conflict of interest:** The authors declare no financial conflict of interest. Andrea B. Kitts is an Associate Editor at JACR.

**Funding:** No funding was provided.

**Disclosure:** During the preparation of this work the author(s) used GPT-4 in order enhance the clarity of select portions of the text. After using this tool/service, the author(s) reviewed and edited the content as needed and take(s) full responsibility for the content of the publication.

**Author contributions:** All authors were responsible for project design, analyzing the data, and manuscript editing. In addition, P.W. conducted most statistical analyses; E.B.G collected the data, analyzed the data, wrote the first draft of the manuscript, and finalized submitted draft.

ChatGPT shows promise in providing accurate, consistent, and relevant responses to patients' imaging-related queries, with prompts further improving responses. However, accuracy was imperfect, and all responses exceeded the recommended reading levels for patient educational materials.

Journal Pre-proof

# 1 Enhancing patient communication with Chat-GPT in 2 radiology: evaluating the efficacy and readability of 3 answers to common imaging-related questions 4

## 5 Abstract

### 6 **Purpose:**

7 To assess ChatGPT's accuracy, relevance, and readability in answering patients' common  
8 imaging-related questions and examine the effect of a simple prompt.

### 9 **Methods:**

10 22 imaging-related questions were developed from categories previously described as important  
11 to patients: safety, the radiology report, the procedure, preparation before imaging, meaning of  
12 terms, and medical staff. These questions were posed to ChatGPT with and without a short  
13 prompt instructing the model to provide an accurate and easy-to-understand response for the  
14 average person. Four board-certified radiologists evaluated the answers for accuracy,  
15 consistency, and relevance. Two patient advocates also reviewed responses for their utility for  
16 patients. Readability was assessed by Flesch Kincaid Grade Level (FKGL). Statistical  
17 comparisons were performed using chi-square and paired t-tests.

### 18 **Results:**

19 264 answers were assessed for both unprompted and prompted questions. Unprompted responses  
20 were accurate 83% (218/264) of the time, which did not significantly change for prompted  
21 responses (87% [229/264];  $P=0.2$ ). The consistency of the responses increased from 72%  
22 (63/88) to 86% (76/88) when prompted ( $P=0.02$ ). Nearly all responses (99% [261/264]) were at  
23 least partially relevant for both question types. Fewer unprompted responses were considered

24 fully relevant at 67% (176/264), though this increased significantly to 80% when prompted  
25 (210/264) ( $P=0.001$ ). The average FKGL was high at 13.6 [12.9-14.2], unchanged with the  
26 prompt (13.0 [12.41-13.60],  $P=0.2$ ). None of the responses reached the eighth-grade readability  
27 recommended for patient-facing materials.

### 28 **Conclusions:**

29 ChatGPT demonstrates the potential to respond accurately, consistently, and relevantly to  
30 patients' imaging-related questions. However, imperfect accuracy and high complexity  
31 necessitate oversight before implementation. Prompts reduced response variability and yielded  
32 more targeted information but did not improve readability.

### 33 **Relevance and Application:**

34 ChatGPT has the potential to increase accessibility to health information and to streamline the  
35 production of patient-facing educational materials, though its current limitations require cautious  
36 implementation and further research.

## 37 **Introduction**

38 The American College of Radiology (ACR) has prioritized effective patient communications in  
39 radiology, encouraging its improvement to lower the barrier of healthcare accessibility (1).

40 Patient concerns are often addressed by direct communication with healthcare workers or by  
41 curated educational information. They can turn to a range of sources, including websites  
42 maintained by radiology departments or by representative bodies such as  
43 [www.radiologyinfo.org](http://www.radiologyinfo.org), which hosts an extensive online repository of descriptions of  
44 procedures, common radiology language, and the significance of common imaging findings.

45 While patients can turn to these reliable information sources and knowledgeable

46 healthcare team members, accessibility is limited by resources such as time, personnel, and  
47 funding.

48 The recent emergence of ChatGPT, a conversational software developed by OpenAI (San  
49 Francisco, California), has garnered attention for its ability to generate human-like text responses  
50 to natural language inquiries, which shows promise as an innovative tool to augment patient  
51 communication in radiology. This generative pre-trained language model (GPLM) employs deep  
52 learning algorithms to analyze text and respond to user inquiries or prompts. ChatGPT could be  
53 used to create reference materials or serve to answer patients' questions prior to a radiology  
54 examination. It has already been tested in several patient-facing medical contexts, including  
55 answering questions related to cardiovascular disease prevention (2) and breast cancer screening  
56 recommendations (3). Yet, the role of this language model in addressing patients' imaging-  
57 related questions remains unexplored.

58 The aim of this paper is to assess the accuracy, relevance, and readability of answers  
59 provided by ChatGPT to common imaging-related questions. A secondary aim of the paper is to  
60 evaluate the influence a patient-directed prompt may have on these parameters.

## 61 **Methods**

62 The study was conducted in March 2023. A list of 22 questions was developed to simulate  
63 common patient concerns in radiology. These questions were tailored based on the professional  
64 expertise of four board-certified radiologists and existing literature (Table 1), using categories  
65 identified by Alarifi et al., who conducted a thematic analysis of online questions from platforms  
66 like Yahoo Answers, Reddit, Quora, and Wiki Answers. Our questions focused on key areas  
67 identified in the study: safety, radiology reports, procedures, preparation before imaging,

68 meaning of terms, and medical staff interactions (4). Questions about the cost of imaging were  
69 excluded.

70 To assess for consistency and account for variability in responses, every question was  
71 posed to ChatGPT (version 3.5) (<https://chat.openai.com>) independently three times, within  
72 unique chat instances. The process was repeated and added with a structured prompt to assess  
73 the modifying effect of a short patient-directed prompt: "Provide an accurate and easy-to-  
74 understand response that is suited for an average person."

75 The output of both unprompted and prompted questions was independently assessed for  
76 accuracy and relevance by four board-certified radiologists who reviewed the responses  
77 sequentially. Accuracy was assessed based on the readers' knowledge and available resources,  
78 including RadiologyInfo<sup>TM</sup>, the patient-facing public information website operated by the  
79 Radiological Society of North America (RSNA), and the ACR (5). Relevance was assessed by  
80 addressing the question, "Is the response relevant and helpful to the average patient?" Responses  
81 were subsequently graded as 1) not relevant or helpful, 2) partially relevant or helpful, or 3) fully  
82 relevant or helpful. Each radiologist also estimated the consistency of the three responses  
83 question, labeling the set of answers as either "consistent" or "inconsistent."

84 The readability of the answers was assessed with the Flesch Kincaid grade level (FKGL),  
85 calculated using an online tool (Readability analyzer: <https://datayze.com/readability-analyzer>).

86 Results for accuracy, relevance, and consistency were considered independent and  
87 pooled. Statistical comparisons between unprompted and prompted responses for these were  
88 performed with chi-square tests. Given the study's exploratory focus, which was not intended for  
89 immediate clinical implementation, we assessed the risk of Type I errors as relatively modest but  
90 Type II errors as more problematic in the context of a rapidly evolving field. As such, we opted

91 against employing a Bonferroni correction. Unpaired Student *t* tests were used to compare  
92 readability and the average response character count. Response character length was calculated  
93 with the Excel LEN function. Statistical analyses were conducted in Stata 14.2. Graphs were  
94 created in Graphpad Prism 9.0. Heatmaps were created with *R* (*R* Core Team, 2023) using the  
95 heatmap package (v1.0.12; Kolde, 2019).

96 For an exploratory analysis, we engaged two patient advocates to provide a lay  
97 perspective on the utility of both prompted and unprompted responses and to identify if there  
98 was a preference between the two. The utility was measured on a 1-3 scale where 1 is not  
99 relevant or helpful, and 3 is fully relevant and helpful.

## 100 Results

101 22 questions were posed to ChatGPT three times each, both with and without an accompanying  
102 prompt, resulting in a total of 132 independent responses. These responses were independently  
103 evaluated by four-board certified radiologists leading to a total of 528 evaluations (264 each for  
104 unprompted and prompted questions). The variability in the grading of accuracy and consistency  
105 among these radiologists is illustrated in Figure S1. The radiologists also performed 176  
106 additional evaluations for consistency across the three outputs for each question, resulting in 88  
107 evaluations each for unprompted and prompted questions.

108 The consistency and accuracy of responses are summarized in Table 2. 83% (218/264) of  
109 the unprompted responses were assessed as accurate. There was no statistically significant  
110 difference in the percentage of accurate answers between unprompted and prompted questions  
111 (83% vs. 87%,  $P=0.2$ ). ChatGPT answered questions for most topics accurately (range of 79% to  
112 96%), except for the "safety" topic, which had a comparably low accuracy rating at 71%



113 (34/48). ChatGPT answered 72% (63/88) of the three repeated questions consistently, which  
114 significantly increased to 86% (76/88) after adding the patient-directed prompt ( $P=0.02$ ).

115 The radiologist reviewers rated very few responses as irrelevant to the question asked.  
116 98.5 % (260/264) and 98.8% (261/264) of answers were determined to be at least partially  
117 relevant for unprompted and prompted questions, respectively (Table 2). Only 67% (176/264) of  
118 unprompted responses were considered fully relevant. When prompted, the number of fully  
119 relevant responses increased significantly to 80% (210/264) ( $P=0.001$ ).

120 The readability of the responses was assessed by Flesch Kincaid grade level (FKGL),  
121 summarized in Figure 1. The average FKGL for unprompted and prompted responses were not  
122 significantly different, despite the request to tailor the response to an average person, measuring  
123 13.6 [12.9-14.2], and 13.0 [12.41-13.60], respectively ( $P=0.2$ ). None of the responses to either  
124 unprompted or prompted questions were at or below the eighth-grade reading level generally  
125 recommended for patient-facing health educational materials (Fig 1B). A minority of all  
126 responses reached below a twelfth-grade level for both unprompted (30% [20/66]) and prompted  
127 (41% [27/66]) questions.

128 The average question length was 1145 [CI 1006-1284] total characters, which was  
129 unchanged when questions were prompted (1147 [CI 1035-1258] total characters;  $P=0.98$ ).

130 Our initial findings of the lay perspective on the utility of responses are divergent from  
131 that of our experts (Table 4). Although the two patient advocates deemed most responses to be  
132 at least partially relevant and helpful (92-97% [122-128/132]), only a minority of unprompted  
133 (42%) and a slight majority of prompted responses (57%) were considered fully relevant. When  
134 given a choice between unprompted or prompted responses, the patient advocates showed a  
135 preference for the latter 71% [31/44] (versus 30% [13/44]).

## 136 Discussion

137 Our results demonstrate that the popular online tool, ChatGPT has the potential to provide  
138 accurate and relevant answers to patient-directed imaging-related questions. These findings  
139 contribute to the growing body of literature that highlights the potential of ChatGPT to automate  
140 time-consuming tasks within the healthcare setting. Automating the development of patient  
141 health educational materials and providing on-demand access to medical questions holds great  
142 promise to improve patient access to health information.

143 While the accuracy, consistency, and relevance of the ChatGPT responses to imaging-  
144 related questions are impressive for a GPLM, they are imperfect; by clinical standards, the  
145 frequency of inaccurate statements that we observed precludes its use without careful human  
146 supervision or review.

147 Our results also show that the readability of the ChatGPT responses far exceeded the  
148 recommended eighth-grade level for patient health educational materials (6, 7). The ability to  
149 understand health information presented to patients is crucial for their capacity to make informed  
150 medical decisions. As it currently stands, the high complexity of the responses clouds the  
151 promise of true patient access to health information. However, the higher levels of complexity  
152 exhibited by ChatGPT's output are also observed for most patient-facing resources that are  
153 currently available across medical fields, including within radiology, which are the likely sources  
154 for the model's pre-training (8, 9).

155 The readability of the responses did not change significantly when accompanied by a  
156 prompt to provide a response to an average person. The lack of improvement may indicate a  
157 resistance to the model's ability to tailor its readability to recommended levels, which would be a  
158 deviation from its pre-trained knowledge, or could represent a deficiency of our prompt.

159 Understanding and addressing these limitations is essential to effectively increase the  
160 accessibility of ChatGPT and could potentially involve the use of more detailed prompting.

161 Prompts play a critical role in affecting the output of GPLMs by providing context and  
162 serving as a source of external knowledge that the model can effectively leverage (10). They can  
163 profoundly affect the output, possibly even dissuading the model from reaching an accurate  
164 response that it would have otherwise concluded (11). While our prompt did not significantly  
165 change the percentage of accurate answers, it elicited a significant increase in fully relevant  
166 responses, potentially highlighting the value of a carefully crafted prompt to provide more  
167 targeted information. Additionally, our prompt may have introduced unmeasured differences that  
168 were noticed by our patient advocates, who preferred the responses with additional prompting.  
169 The efficacy of proper prompt engineering should be evaluated in the medical context, which  
170 includes patient involvement for patient-facing materials, to optimize the output of accurate and  
171 relevant medical information.

172 The layperson's assessment of response utility was unexpectedly critical of responses  
173 both with and without a prompt, assessing fewer than half as fully relevant/helpful. This  
174 contrasts with physicians' evaluations, although we did not tailor this secondary aim for  
175 statistical comparison. This discrepancy underscores the need for more inclusive evaluations of  
176 patient-facing materials, challenging the approach of using only expert reviewers.

177 There are several limitations to our study. The rapidly evolving nature of technology  
178 introduces unpredictability in evaluating the platform's performance, which may change in the  
179 future. Additionally, radiologists —not patients— wrote the questions and evaluated the answers,  
180 which does not fully capture the unpredictable variety of ways patients seek, consume, or  
181 understand medical information. Variability in question context and phrase is expected to affect

182 the output just as prompts do; a straightforward example is that the model responds in the same  
183 language that the question was asked. Additionally, it is unknown how ChatGPT would perform  
184 when questions are posed and responded to in a language other than English. Furthermore, some  
185 areas of knowledge may have been less accessible or represented during the model's pre-training,  
186 which may lead to less accurate responses; for example, while overall accurate, our study  
187 showed that certain topics such as "safety" could potentially be less accurate or relevant.  
188 Additionally, as an exploratory study, extrapolation of the findings is limited; the inherent  
189 stochastic output of GPLMs and the tendency to "confabulate", or confidently produce factually  
190 incorrect responses, may become more clinically apparent in larger cohorts.

191 In conclusion, our exploratory analysis underscores the potential of ChatGPT to  
192 streamline time-consuming tasks in patient health education. However, the immediate  
193 application of this system necessitates a cautious approach and further investigation. It is vital to  
194 tackle challenges related to readability and mitigate the risks of presenting potentially misleading  
195 information to patients. Exploring strategies such as effective, prompt engineering will  
196 contribute to optimizing ChatGPT's output, ensuring its safety and effectiveness for patient use.

## 197 **Summary statement:**

198 ChatGPT shows promise in providing accurate, consistent, and relevant responses to patients'  
199 imaging-related queries, with prompts further improving responses. However, accuracy was  
200 imperfect, and all responses exceeded the recommended reading levels for patient educational  
201 materials.

## 202 Take Home Points

- 203 • ChatGPT provided an accurate response to patients' common imaging-related questions  
204 83% of the time and 87% when asked with a simple prompt, although this difference was  
205 not statistically significant ( $P=0.2$ ).
- 206 • Although almost always partially relevant (99%), the proportion of ChatGPT responses  
207 considered fully relevant significantly rose from 67% to 80% ( $P=0.001$ ) when a simple  
208 prompt accompanied the questions. Prompting also improved the response consistency  
209 from 72% to 86% ( $P=0.02$ ).
- 210 • ChatGPT responses were uniformly complex, with no response reaching the  
211 recommended eighth-grade level for patient-facing materials (average Flesch Kincaid  
212 grade level of 13.6 unprompted and 13.0 prompted [ $P=0.2$ ]).  
213

## References

1. Siewert B, Bruno MA, Fleishon HB, et al. Summary of the 2022 ACR Intersociety Meeting. *Journal of the American College of Radiology* 2023;20(5):479–486.
2. Sarraju A, Bruemmer D, Van Iterson E, Cho L, Rodriguez F, Laffin L. Appropriateness of Cardiovascular Disease Prevention Recommendations Obtained From a Popular Online Chat-Based Artificial Intelligence Model. *JAMA [Internet]* 2023 [cited 2023 Feb 24]; Available from: <https://jamanetwork.com/journals/jama/fullarticle/2801244>
3. Haver HL, Ambinder EB, Bahl M, Oluyemi ET, Jeudy J, Yi PH. Appropriateness of Breast Cancer Prevention and Screening Recommendations Provided by ChatGPT. *Radiology*

- [Internet] Radiological Society of North America; 2023 [cited 2023 Apr 5]; Available from: <https://pubs.rsna.org/doi/10.1148/radiol.230424>
4. Alarifi M, Patrick T, Jabour A, Wu M, Luo J. Understanding patient needs and gaps in radiology reports through online discussion forum analysis. *Insights Imaging* 2021;12(1):50.
  5. RadiologyInfo [Internet]. Radiologyinfo.org 2023 [cited 2023 Mar 20]. Available from: <https://www.radiologyinfo.org/en>
  6. Doak CC, Doak LG, Root JH. Teaching Patients with Low Literacy Skills. *AJN The American Journal of Nursing* 1996;96(12):16M.
  7. Agency for Healthcare Research and Quality. Assess, Select, and Create Easy-to-Understand Materials: Tool #11 [Internet]. 2020. Available from: <https://www.ahrq.gov/health-literacy/improve/precautions/tool11.html>
  8. Hansberry DR, Agarwal N, Baker SR. Health Literacy and Online Educational Resources: An Opportunity to Educate Patients. *American Journal of Roentgenology American Roentgen Ray Society*; 2015;204(1):111–116.
  9. Rooney MK, Santiago G, Perni S, et al. Readability of Patient Education Materials From High-Impact Medical Journals: A 20-Year Analysis. *Journal of Patient Experience SAGE Publications Inc*; 2021;8:2374373521998847.
  10. Lee P, Bubeck S, Petro J. Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine. Drazen JM, Kohane IS, Leong T-Y, editors. *N Engl J Med* 2023;388(13):1233–1239.

11. Zuccon G, Koopman B. Dr ChatGPT, tell me what I want to hear: How prompt knowledge impacts health answer correctness [Internet]. arXiv; 2023 [cited 2023 Apr 18]. Available from: <http://arxiv.org/abs/2302.13793>

## 214 Table and Figure legends:

**Table 1:** Common imaging-related topics and patient questions

**Table 2:** Accuracy and consistency of responses

**Table 3:** Responses assessed as at least partially or fully relevant.

**Table 4:** Utility of responses, as assessed by patient advocates.

**Figure 1:** Readability of ChatGPT responses to questions accompanied without (–) or with (+) a prompt. A) Violin plots of Flesch Kincaid grade level (FKGL) show the distribution of responses with median and IQR; striped lines at FKGL=12 and FKGL=8 for reference, the latter of which is the upper limit of the recommended readability of patient-facing educational materials B) Heatmap of the FKGL for ChatGPT responses by topic and prompt status. Each row represents a question, and each cell represents an individual response to questions posed in triplicate.

**Figure S1:** Individual reviewer scores of ChatGPT responses to questions accompanied without or with a prompt. A) Heatmap of accuracy scores assigned by reviewers for ChatGPT responses by topic and prompt status. Each row represents a question, and each cell represents an individual response to questions posed in triplicate. B) Heatmap of the relevance scores assigned by reviewers for ChatGPT responses by topic and prompt status. Each row represents a question, and each cell represents an individual response to questions posed in triplicate.

**Table 1:** Common imaging-related topics and patient questions

<b>Topic</b>	<b>Question</b>
<b>Report</b>	<p>1. What does “clinical correlation” mean in a radiology report?</p> <p>2. What does it mean if a mass is “stable’ in the radiology report?</p> <p>3. In the radiology report, my liver is described as “steatotic”. What does it mean?</p> <p>4. What does it mean if my tumor is “hypermetabolic” on my PET-CT scan?</p>
<b>Safety</b>	<p>5. Do X-rays or CT scans cause cancer?</p> <p>6. What are the risk of MRI during pregnancy?</p> <p>7. What are the risks of Gadolinium deposition?</p> <p>8. I have an allergy to shellfish. Can I get iodine CT contrast?</p>
<b>Procedure</b>	<p>9. How long does it take to obtain a brain MRI?</p> <p>10. What does it feel like to get IV contrast?</p> <p>11. What is it like to be inside an MR machine?</p> <p>12. What are the risks of an imaging-guided renal biopsy?</p>
<b>Preparation</b>	<p>13. Why do I need to drink contrast before some CTs but not others?</p> <p>14. Why do I need intravenous contrast for my CT scan?</p> <p>15. What can I do to reduce my anxiety during an MRI scan?</p> <p>16. Why do I need to isolate myself from others after some but not other nuclear medicine exams?</p>
<b>Meaning</b>	<p>17. What is the difference between an open MRI and a closed MRI?</p> <p>18. What is the difference between a CT and an MRI?</p> <p>19. What is the difference between a CT and an ultrasound?</p> <p>20. How is a blocked vessel opened in interventional radiology?</p>
<b>Medical Staff</b>	<p>21. Who is a radiologist?</p> <p>22. Who can answer questions about my radiology report?</p>



**Table 2:** Accuracy and consistency of responses

	Consistency (% , [n/N])			Accuracy (% , [n/N])		
	No prompt	Prompted	<i>P</i>	No prompt	Prompted	<i>P</i>
<b>All questions</b>	71.6 (63/88)	86.4 (76/88)	0.016*	82.6 (218/264)	86.7 (229/264)	0.2
<b>Report</b>	87.5 (14/16)	100 (16/16)	0.1	89.6 (43/48)	85.4 (41/48)	0.5
<b>Safety</b>	62.5 (10/16)	75 (12/16)	0.4	70.8 (34/48)	81.3 (39/48)	0.2
<b>Procedure</b>	62.5 (10/16)	81.3 (13/16)	0.2	87.5 (42/48)	91.7 (44/48)	0.5
<b>Preparation</b>	75 (12/16)	87.5 (14/16)	0.4	79.2 (38/48)	77.1 (37/48)	0.8
<b>Meaning</b>	62.5 (10/16)	87.5 (14/16)	0.1	79.2 (38/48)	91.7 (44/48)	0.08
<b>Medical Staff</b>	87.5 (7/8)	87.5 (7/8)	1.0	95.8 (23/24)	100.0 (24/24)	0.3

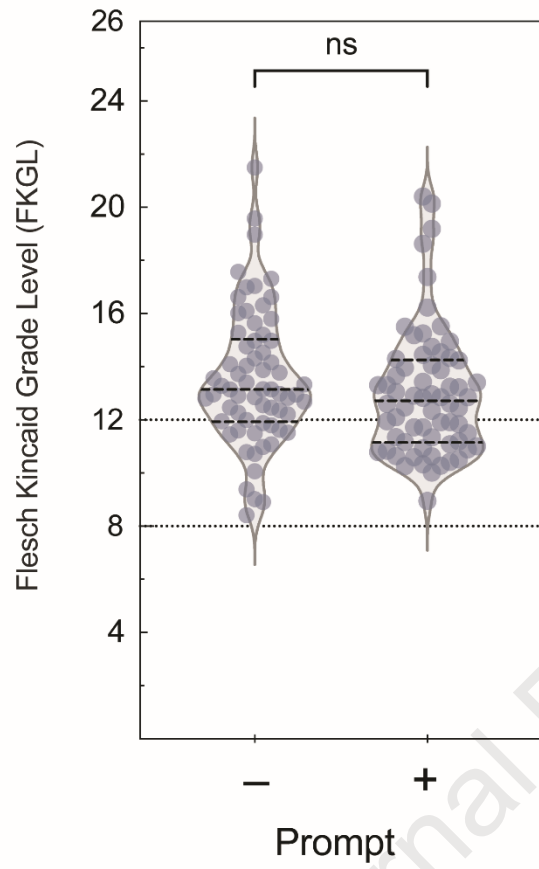
**Table 3:** Responses assessed as at least partially or fully relevant

	At least partially relevant (% , [n/N])			Fully relevant (% , [n/N])		
	No prompt	Prompt	<i>P</i>	No prompt	Prompt	<i>P</i>
<b>All questions</b>	98.5 [260/264]	98.9 [261/264]	0.70	66.7 [176/264]	79.6 [210/264]	0.001**
<b>Report</b>	100.0 [48/48]	100.0 [48/48]		77.1 [37/48]	91.7 [44/48]	0.049*
<b>Safety</b>	100.0 [48/48]	100.0 [48/48]		50.0 [24/48]	64.6 [31/48]	0.15
<b>Procedure</b>	97.9 [47/48]	100.0 [48/48]	0.32	70.8 [34/48]	79.2 [38/48]	0.35
<b>Preparation</b>	95.8 [46/48]	93.8 [45/48]	0.65	64.6 [31/48]	68.8 [33/48]	0.67
<b>Meaning</b>	97.9[47/48]	100 [48/48]	0.32	54.2 [26/48]	91.7[44/48]	0.001***
<b>Medical Staff</b>	100 [24/24]	100 [24/24]		100 [24/24]	83.3 [20/24]	0.037*

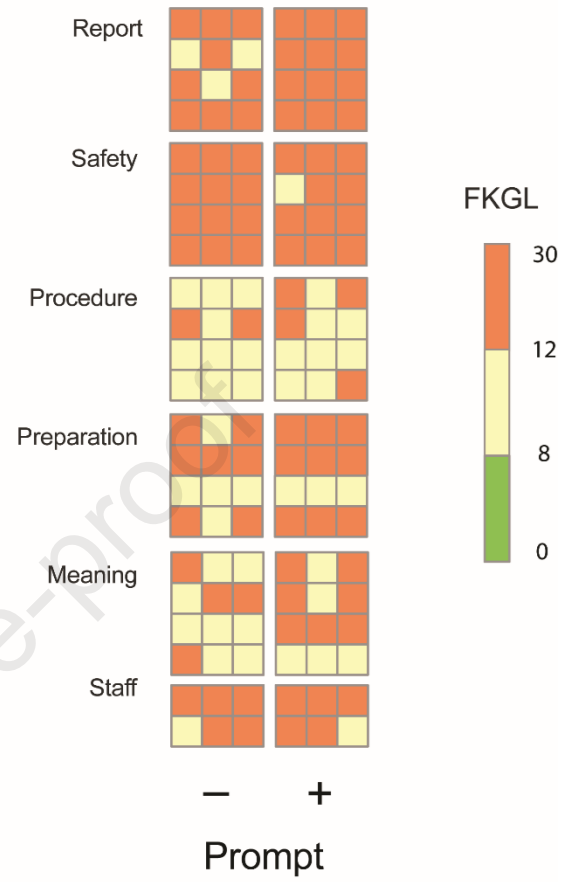
**Table 4:** Utility of responses, as assessed by patient advocates.

	<b>At least partially relevant/helpful</b> (%, [n/N])		<b>Fully relevant/helpful</b> (%, [n/N])	
	<b>No prompt</b>	<b>Prompt</b>	<b>No prompt</b>	<b>Prompt</b>
<b>All questions</b>	92 [122/132]	97 [128/132]	42 [55/132]	57 [75/132]
<b>Report</b>	96 [23/24]	100 [24/24]	54 [13/24]	75 [18/24]
<b>Safety</b>	92 [22/24]	96 [23/24]	38 [9/24]	58 [14]
<b>Procedure</b>	96 [22/24]	100 24/24	50 [12/24]	63 [15/24]
<b>Preparation</b>	83 [20/24]	96 [23/24]	29 [7/24]	38 [9/24]
<b>Meaning</b>	92 [22/24]	100 [24/24]	42 [10/24]	63 [15/24]
<b>Medical Staff</b>	100 [12/12]	83 [10/12]	33 [4/12]	33 [4/12]

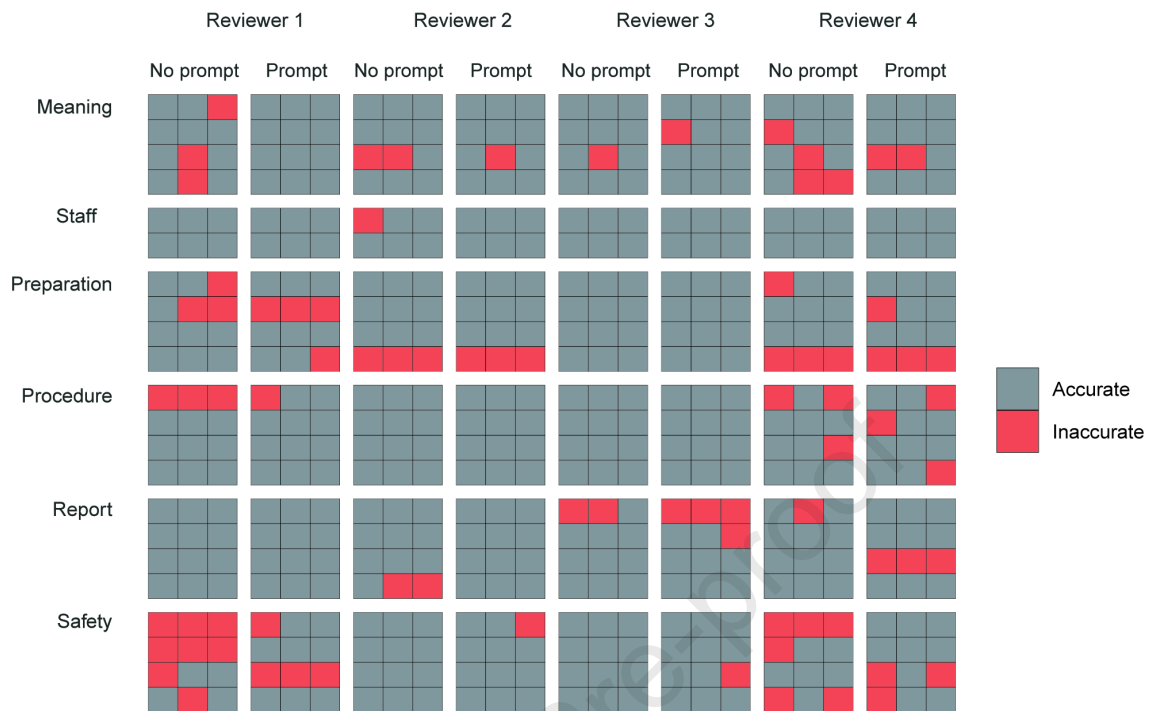
A)



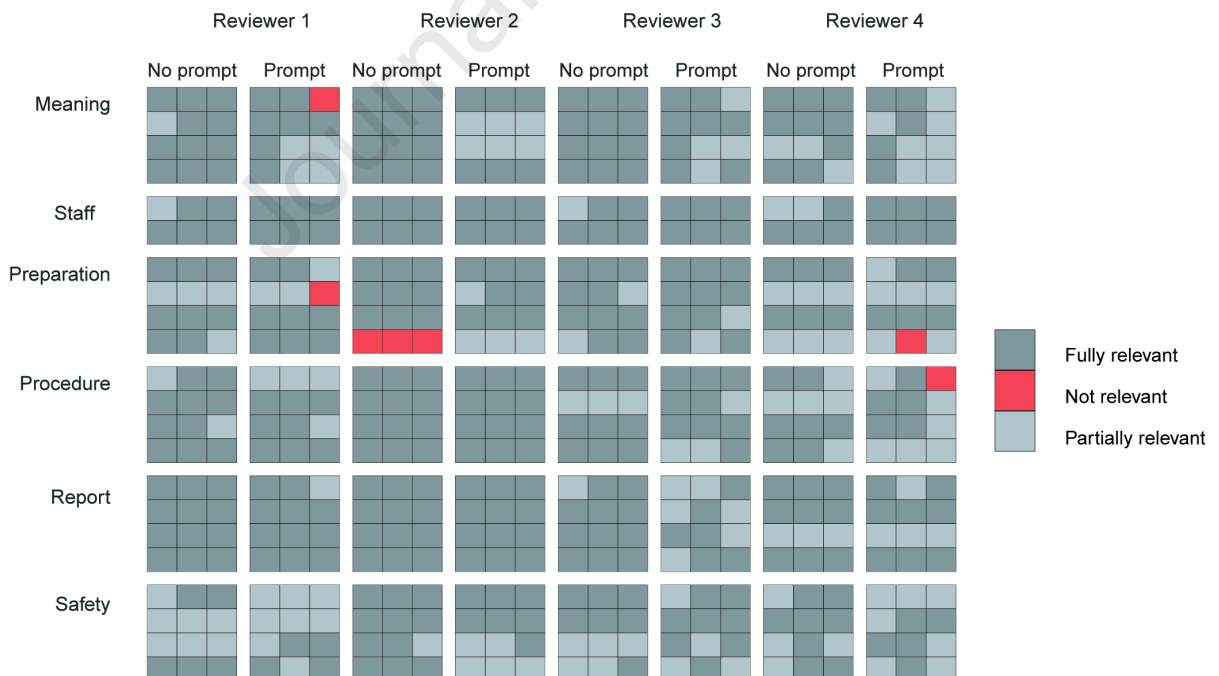
B)



## A) Accuracy scores



## B) Relevance scores



**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

---

One of our co-authors, Andrea Borondy Kitts, is an Associate Editor at JACR.

Journal Pre-proof

- ChatGPT provided an accurate response to patients' common imaging-related questions 83% of the time and 87% when asked with a simple prompt, although this difference was not statistically significant ( $P=0.2$ ).
- Although almost always partially relevant (99%), the proportion of ChatGPT responses considered fully relevant significantly rose from 67% to 80% ( $P=0.001$ ) when a simple prompt accompanied the questions. Prompting also improved the response consistency from 72% to 86% ( $P=0.02$ ).
- ChatGPT responses were uniformly complex, with no response reaching the recommended eighth-grade level for patient-facing materials (average Flesch Kincaid grade level of 13.6 unprompted and 13.0 prompted [ $P=0.2$ ]).