

1 **Is ChatGPT smarter than Otolaryngology trainees?**

2 **A comparison study of board style exam questions**

3 J Patel<sup>1</sup>; PZ Robinson<sup>2</sup>; EA Illing<sup>1</sup>; BP Anthony<sup>1</sup>

4 Indiana University School of Medicine, Department of Otolaryngology – Head and Neck

5 Surgery<sup>1</sup>, Indiana University School of Medicine<sup>2</sup>

6

7 **Running Title:** ChatGPT verses Otolaryngology trainees

8 **Financial Support:** None

9 **Conflict of interest:** None

10 **Corresponding author**

11 Benjamin Anthony, MD

12 Indiana University School of Medicine | Indiana University Health Physicians

13 1130 W. Michigan Street | FH 400

14 Indianapolis, IN 46202

15 Phone | 317.963.7082

16 Email | [bpanthon@iu.edu](mailto:bpanthon@iu.edu)

17

18 **Abstract**

19 **Objectives:** This study compares the performance of the artificial intelligence (AI) platform  
20 Chat Generative Pre-Trained Transformer (ChatGPT) to Otolaryngology trainees on board style  
21 exam questions.

22 **Methods:** We administered a set of 30 Otolaryngology board style questions to medical students  
23 (MS) and Otolaryngology residents (OR). 31 MSs and 17 ORs completed the questionnaire. The  
24 same test was administered to ChatGPT version 3.5, five times. Comparisons of performance  
25 were achieved using a one-way ANOVA with Tukey Post Hoc test, along with a regression  
26 analysis to explore the relationship between education level and performance.

27 **Results:** The average scores increased each year from MS1 to PGY5. A one-way ANOVA  
28 revealed that ChatGPT outperformed trainee years MS1, MS2, and MS3 ( $p = <0.001, 0.003,$  and  
29  $0.019,$  respectively). PGY4 and PGY5 otolaryngology residents outperformed ChatGPT ( $p =$   
30  $0.033$  and  $0.002,$  respectively). For years MS4, PGY1, PGY2, and PGY3 there was no statistical  
31 difference between trainee scores and ChatGPT ( $p = .104, .996,$  and  $1.000,$  respectively).

32 **Conclusion:** ChatGPT can outperform lower-level medical trainees on Otolaryngology board-  
33 style exam but still lacks the ability to outperform higher-level trainees. These questions  
34 primarily test rote memorization of medical facts; in contrast, the art of practicing medicine is  
35 predicated on the synthesis of complex presentations of disease and multilayered application of  
36 knowledge of the healing process. Given that upper-level trainees outperform ChatGPT, it is  
37 unlikely that ChatGPT, in its current form will provide significant clinical utility over an  
38 Otolaryngologist.

39 **Keywords:** Artificial Intelligence, Medical Education, Comprehensive Otolaryngology

40 **Level of Evidence:** Level 2

## 41 **Introduction**

42           Current developments in artificial intelligence (AI) technology using advanced language  
43 models have generated a significant amount of public interest. Chat Generative Pre-Trained  
44 Transformer (ChatGPT), an AI-based language model developed by OpenAI, stands out for its  
45 ability to generate human-like responses in written format. Recent improvements to ChatGPT  
46 have garnered significant attention as this sophisticated AI platform finds its place in modern  
47 society. Fueled by vast databases, ChatGPT provides precise, personalized answers, a testament  
48 to its prowess in understanding the intricacies of human language. Based on this repository of  
49 knowledge, this language model effortlessly mirrors real-life conversations and boasts profound  
50 knowledge across diverse subjects(1).

51           The role of AI in medicine has been met with both hopeful intrigue as well as skepticism.  
52 AI-powered systems like ChatGPT can provide immediate access to information for patients and  
53 healthcare providers to augment healthcare decisions. ChatGPT seems to have an obvious role in  
54 patient education and medical education due to its ability to generate knowledgeable responses to  
55 fact-based questions with categorical answers. ChatGPT could possibly even play a direct role in  
56 augmenting patient care decisions and treatment. However, the accuracy and reliability of AI  
57 systems like ChatGPT has not yet been firmly established in medicine. Nevertheless, efforts  
58 continue to further develop this technology to determine if it holds value for patient care.

59           ChatGPT has been tested with a diverse list of standardized examinations, such as the  
60 uniform Bar Examination, the Scholastic Assessment test (SAT), the Graduate Record  
61 Examination (GRE), high school advanced placement exams and more(2). Despite medicine  
62 being filled with niche terminology, acronyms, and multidisciplinary topics, ChatGPT has been  
63 able to exhibit a broad knowledge of medicine. Indeed, ChatGPT was found to likely be able to

64 pass the USMLE Step 1 examination(3). With regards to subspecialty fields, the literature has  
65 shown that ChatGPT is passable or near passable in board exams for Ophthalmology, Pathology,  
66 Neurosurgery, Cardiology, and Otolaryngology(3-9); however, ChatGPT did quite poorly on the  
67 multiple-choice Orthopedic board exam(10). As a repository of advanced medical knowledge,  
68 ChatGPT underperformed in comparison to the widely used UpToDate medical reference(11).  
69 AI based language models could be a great tool when patients desire reliable information on  
70 upcoming procedures, information on prescriptions, and other aspects of their care that carry  
71 significant weight to the patient(12), but their utility in advanced medical decision making  
72 remains to be investigated.

73 This current project compares the performance of ChatGPT version 3.5 to medical  
74 trainees at a US medical school and residency on board style questions for the Otolaryngology –  
75 Head and Neck Surgery board exam. The spectrum of questions ranged from fundamental  
76 concepts learned during the infancy of medical school to the complexities of advanced medical  
77 and surgical patient management derived by the end of resident training. Our primary aim is to  
78 assess if and when ChatGPT can outperform human learners on Otolaryngology board style  
79 questions.

## 80 **Materials and Methods**

81 This study was exempt from requiring approval by the institutional review board at  
82 Indiana University. The study started collecting data on October 2nd, 2023, through January 5th,  
83 2024. 30 multiple choice Otolaryngology board-style questions were asked to all years of  
84 medical students and Otolaryngology residents. The same questions were also asked to  
85 ChatGPT. Given that ChatGPT is a reiterative, learning-based model with a potential for  
86 different answers each time a question is asked, the test was administered to ChatGPT five times.

87           Questions were dispersed by using Google Forms to all medical students, years 1-4,  
88 (MS1-MS4) and Otolaryngology residents, years 1-5, (PGY1-PGY5) at Indiana University  
89 School of Medicine. Participants were blinded to the purpose of this exam to avoid bias, thus  
90 they were not provided informed consent on underlying purpose of the study. They were simply  
91 asked to answer questions to test the quality of the questions written. No compensation or  
92 incentives were provided for the completion of this questionnaire. The only identifying data  
93 collected was the education level of each participant (MS1-PGY5). At the beginning of the  
94 study, the participants were given clear instructions: “Thanks so much for taking the time to  
95 answer this 30-question quiz that covers topics within Otolaryngology. We ask that you take this  
96 quiz in one sitting and do not use outside resources. This will allow us to accurately evaluate the  
97 questions written.”

98           For ChatGPT, the model was prompted with the following: “You are a medical  
99 professional and I want you to pick an answer from the multiple-choice question I provide.” For  
100 example, in one administration, ChatGPT responded with: “Of course, I would be happy to help  
101 you with multiple choice questions related to medical topics. Please provide the question and its  
102 options, and I'll do my best to provide you with the correct answer and explanation.” Following  
103 this prompt, each of the 30 questions were provided one at a time. The answer and reasoning  
104 were recorded. The test was administered five times, once each day on five different days. This  
105 methodology was utilized to help capture the variability that language models can exhibit. We  
106 believe this allowed ChatGPT additional chances to retrieve the correct information within the  
107 vast databases it utilizes.

108           **Participants:** The 30-question survey was completed by medical students and  
109 Otolaryngology residents at Indiana University (n = 48) and ChatGPT model 3.5 (n = 5). There

110 were 9 education level groups across the human participants, MS1 (n = 8), MS2 (n = 7), MS3 (n  
111 = 10), MS4 (n = 6), PGY1 (n = 4), PGY2 (n = 4), PGY3 (n = 4), PGY4 (n = 2), and PGY5 (n =  
112 3). See Table 1.

113 Table 1 - Demographics of participants

Level of Education	Number of participants
MS 1	8
MS 2	7
MS 3	10
MS 4	6
PGY – 1	4
PGY – 2	4
PGY – 3	4
PGY – 4	2
PGY – 5	3
MS – Medical Student Year, PGY – Post graduate year	

114  
115 **Statistical analysis:** Statistical analysis was conducted using Statistical Package for the  
116 Social Sciences (SPSS). A one-way ANOVA was conducted to compare Otolaryngology Board  
117 Exam Scores between human participants at each medical education level and ChatGPT. The  
118 ANOVA was implemented to identify if group differences were present between the 9 education  
119 levels (MS1-PGY5) and ChatGPT. Tukey’s Honest Significant Difference Test (HSD) post hoc  
120 test was utilized to identify which of the 9 education levels (MS1-PGY5) differed to ChatGPT. A

121 regression analysis was conducted to explore the relationship between education level and score,  
 122 specifically to explore whether education level predicted score.

123

124 **Results**

125 A regression revealed that the education level significantly predicted score  $R^2 =$   
 126 .765,  $F(1, 46) = 150.003$ ,  $p < .001$ . The average score of human participants increased linearly as  
 127 education level increased by years (MS1-PGY5) (MS1 = 28.75%; MS2 = 31.44%; MS3 = 36%;  
 128 MS4 = 37.77%; PGY1 = 49.18%; PGY2 = 56.68%; PGY3 = 70.83%; PGY4 = 81.65%; PGY5 =  
 129 84.47%,). See table 2.

130 Table 2 – Percent correct and mean difference between ChatGPT and Medical Trainees.

Group A	Average % Correct	Group B	Average % Correct	Mean Difference (A-B)	Sig.	95% Confidence Interval	
						Lower Bound	Upper Bound
ChatGPT	54.66	MS1	28.75	25.91*	<.001	8.40	43.43
		MS2	31.44	23.22*	.003	5.22	41.21
		MS3	36.00	18.66*	.019	1.83	35.49
		MS4	37.77	16.89	.104	-1.72	35.50
		PGY-1	49.18	5.49	.996	-15.13	26.10
		PGY-2	56.68	-2.01	1.000	-22.63	18.60
		PGY-3	70.83	-16.17	.242	-36.78	4.45
		PGY-4	81.65	-26.99*	.033	-52.70	-1.28
		PGY-5	84.47	-29.81*	.002	-52.25	-7.36
MS – Medical Student Year, PGY – Post graduate year							

131

132 The average score of ChatGPT was 54.66% across the 5 administrations. At times,  
 133 ChatGPT did provide different answers to questions with different explanations. However, there  
 134 was not a consistent increase in percent correct overtime. By mean, ChatGPT out-performed  
 135 human participants from education level MS1-PGY1 but underperformed in comparison to  
 136 PGY2-PGY5. See Fig 1.

137 Figure 1 title: Board Exam Scores between Medical Trainees and ChatGPT.

138

139 A one-way ANOVA revealed that there were statistically significant differences in the  
140 average score between at least two of the 10 groups ( $F(9, 43) = [20.393]$ ,  $p < .001$ ).

141 Tukey's HSD test for multiple comparisons were implemented to identify which groups  
142 differed significantly from each other, particularly from ChatGPT. Results revealed that the score  
143 significantly differed between ChatGPT and MS1 ( $p < .001$ , 95% C.I. = 8.3905, 43.4295), MS2  
144 ( $p = .003$ , 95% C.I. = 5.2228, 41.2115), MS3 ( $p = .019$ , 95% C.I. = 1.8278, 35.4922), PGY-4 ( $p$   
145 = .033, 95% C.I. = -52.7016, -1.2784), PGY-5 ( $p = .002$ , 95% C.I. = -52.2496, -7.3637).

146 Results revealed that the score did not significantly differ between ChatGPT and MS4 ( $p$   
147 = .104, 95% C.I. = -1.7154, 35.5020), nor between ChatGPT and PGY-1 ( $p = .996$ , 95% C.I. = -  
148 15.1302, 26.1002), nor PGY-2 ( $p = 1.000$ , 95% C.I. = -22.6302, 18.6002), nor PGY-3 ( $p = .242$ ,  
149 95% C.I. = -36.7802, 4.4502).

## 150 Discussion

151 Language-centric AI models, exemplified by ChatGPT, are gaining momentum for their  
152 ability to sustain coherent conversations, and demonstrating aptitude on standardized  
153 examinations. Powered by deep machine learning techniques and extensive textual data,  
154 ChatGPT iteratively enhances its abilities via user interactions and reinforcement learning. This  
155 research explicates ChatGPT's deficiency in tackling complex medical multiple-choice  
156 questions, contrasting its performance with that of medical students and Otolaryngology trainees.  
157 Findings reveal ChatGPT's superiority over beginners but eventual inferiority to seasoned  
158 residents on board-style questions targeting Otolaryngology knowledge, indicating a progressive  
159 convergence in performance.



160 One of the key findings that we believe challenged ChatGPT was the nuanced and  
161 context-dependent nature of medical questions. While it provided suitable explanations for its  
162 reasoning on specific queries, there were instances where it seemed to grapple with a lack of  
163 understanding or data support, leading to what appeared as a guess, misinformed, or ill-informed  
164 answer. This was seen through multiple repetitions of the question with either similar answer  
165 choice but different explanation and vice versa. While illustrating the robust power of this  
166 language model, these inconsistencies beg the questions about continued knowledge gaps in  
167 specific queries on AI language models. Thus, while the model demonstrated an impressive  
168 ability to generate human-like responses in natural language, it continues to struggle with the  
169 intricacies and subtleties inherent in otolaryngology, and perhaps medicine generally.

170 Different from the patterns shown by repeated administration to ChatGPT, medical  
171 learners exhibited marked growth in their knowledge base, showcasing a linear progression in  
172 their average correct responses on the exam over years of continued training. This aligns with  
173 our expectations, as their evolving domain-specific knowledge, clinical experiences, and the  
174 ability to interpret complex scenarios increases with seniority.

175 Further examining our findings, the interpretability of responses emerged as a critical  
176 factor in evaluating the performance of ChatGPT. Despite its ability to generate coherent and  
177 grammatically correct answers, deciphering the underlying reasoning process posed a significant  
178 challenge. For example, ChatGPT was able to identify the correct answer without offering the  
179 accurate explanation, and vice versa. Upon multiple assessments of the same question, the  
180 rationale and explanation underwent changes at times, resulting in a different answer choice.  
181 This implies a potential learning process, where continuous exposure to queries builds on the  
182 model's knowledge base, enabling it to generate more accurate responses, indeed, an avenue for

183 future research to investigate. Consequently, ChatGPT remains rudimentary in its ability to  
184 become the gold standard for querying medical questions. This may be in part due to its lack of a  
185 deep understanding of patient-specific factors, consideration of evolving clinical contexts, and  
186 the incorporation of the latest medical research, specifically in Otolaryngology. Future research  
187 should explore how AI language models can be trained to better perform answer medical queries.  
188 Further investigation should continue to be done to test the growth of ChatGPT as the model  
189 advances.

190 Human participants, in contrast, are adept at synthesizing information, applying critical  
191 thinking skills, and adapting responses to the intricacies of each scenario. This foundational skill  
192 is nurtured throughout the educational journey, particularly for individuals in the medical field.  
193 Resultantly, senior Otolaryngology residents demonstrate superior deductive abilities in  
194 answering multiple-choice questions compared to ChatGPT. Nevertheless, medical trainees  
195 historically rely on diverse study aids to cultivate this deductive ability and expand their  
196 knowledge base. As AI continues to advance, it is essential to acknowledge ChatGPT's potential  
197 applications and advantages. It excels in non-clinical settings where general knowledge and  
198 language understanding are crucial. Given time, ChatGPT's and other AI model's knowledge is  
199 anticipated to expand. Thus, AI may acquire the capability to dynamically update its knowledge  
200 base in real-time, and use increasingly complex informational sources accurately, to emerge as  
201 an invaluable tool for medical learners and potentially even patients.

202 This introduces the avenue for future researchers to consider the ethical implications of  
203 AI in medicine. As we continue our efforts to attempt the integration of AI into medical  
204 decision-making processes, there remains much skepticism on its utility, and rightfully so. While  
205 AI offers unprecedented capabilities for analyzing vast amounts of patient data and providing

206 diagnostic insights, it also introduces a complex ethical dilemma. Accountability, transparency,  
207 and obsolescence of a profession are at the forefront of this multifaceted dilemma. In its current  
208 infancy, AI is nonthreatening to a physician as a profession as our interactions with patients are  
209 pivotal to providing hands on care. Moreover, empathy and compassion are pillars in the dogma  
210 of healthcare, which are human qualities and not yet replicable by AI. Regarding accountability,  
211 physicians must take ownership of their decisions which can greatly impact the lives of their  
212 patients. An AI in contrast has no accountability for providing its opinion as it is not presently  
213 governed to do as such. The decision-making processes of an AI can also appear opaque, making  
214 it challenging to understand how it arrived at that conclusion. Additionally, there are worries  
215 regarding bias and fairness, as AI systems can inadvertently perpetuate or even amplify existing  
216 biases present in the data used to train them, potentially leading to worsening disparities in  
217 healthcare outcomes. Likewise, the issue of patient autonomy and informed consent becomes  
218 paramount when AI systems are employed in medical decision-making, as patients may not fully  
219 comprehend or have control over the algorithms guiding their care. As healthcare continues to  
220 embrace AI technologies, navigating these moral quandaries will be crucial to ensure the  
221 responsible, ethical, and equitable use of AI in medical practice.

222

## 223 **Conclusion**

224 In conclusion, our findings emphasize the need for caution and meticulous assessment  
225 when deploying language models in specialized fields like otolaryngology or medicine, where  
226 precision is critical, and the stakes are high. ChatGPT showcases remarkable capabilities in  
227 natural language understanding and has been shown to pass a host of different board  
228 examinations(2-8). In our study, ChatGPT scored an average of 54.66% which is similar to the

229 57% correct seen in Hoch et al(9). Considering this, ChatGPT is not yet intelligent enough to  
230 become the trusted gold standard to accessing medical information within Otolaryngology.

231 Additionally, AI systems cannot replicate human elements of care such as empathy,  
232 compassion, and ethical judgement, which are essential tenants of healthcare. Future research  
233 may focus on refining and tailoring language models for specific domains, incorporating real-  
234 time learning mechanisms, and addressing the interpretability challenges associated with  
235 automated systems in complex decision-making processes within the medical field.

236 Consequently, with time, AI language models may evolve into indispensable tools for medical  
237 professionals and potentially even to patients and future research must aim to keep our  
238 understanding of their limits and abilities up to date.

239

240 References:

- 241
- 242 1. Schade M. How ChatGPT and Our Language Models Are Developed.
  - 243 2. L. V. AI models like ChatGPT and GPT-4 are acing everything from the bar exam to AP  
244 Biology. Here's a list of difficult exams both AI versions have passed. Business Insider. 2023.
  - 245 3. Gilson A, Safranek CW, Huang T, Socrates V, Chi L, Taylor RA, et al. How Does ChatGPT  
246 Perform on the United States Medical Licensing Examination? The Implications of Large  
247 Language Models for Medical Education and Knowledge Assessment. *JMIR Med Educ*.  
248 2023;9:e45312.
  - 249 4. Long C, Lowe K, Zhang J, Santos AD, Alanazi A, O'Brien D, et al. A Novel Evaluation Model  
250 for Assessing ChatGPT on Otolaryngology-Head and Neck Surgery Certification Examinations:  
251 Performance Study. *JMIR Med Educ*. 2024;10:e49970.
  - 252 5. Antaki F, Touma S, Milad D, El-Khoury J, Duval R. Evaluating the Performance of ChatGPT  
253 in Ophthalmology: An Analysis of Its Successes and Shortcomings. *Ophthalmol Sci*.  
254 2023;3(4):100324.
  - 255 6. Sinha RK, Deb Roy A, Kumar N, Mondal H. Applicability of ChatGPT in Assisting to Solve  
256 Higher Order Problems in Pathology. *Cureus*. 2023;15(2):e35237.
  - 257 7. Ali R, Tang OY, Connolly ID, Fridley JS, Shin JH, Zadnik Sullivan PL, et al. Performance of  
258 ChatGPT, GPT-4, and Google Bard on a Neurosurgery Oral Boards Preparation Question Bank.  
259 *Neurosurgery*. 2023;93(5):1090-8.
  - 260 8. Ali R, Tang OY, Connolly ID, Zadnik Sullivan PL, Shin JH, Fridley JS, et al. Performance of  
261 ChatGPT and GPT-4 on Neurosurgery Written Board Examinations. *Neurosurgery*.  
262 2023;93(6):1353-65.
  - 263 9. Hoch CC, Wollenberg B, Luers JC, Knoedler S, Knoedler L, Frank K, et al. ChatGPT's quiz  
264 skills in different otolaryngology subspecialties: an analysis of 2576 single-choice and multiple-  
265 choice board certification preparation questions. *Eur Arch Otorhinolaryngol*. 2023;280(9):4271-  
266 8.
  - 267 10. Lum ZC. Can Artificial Intelligence Pass the American Board of Orthopaedic Surgery  
268 Examination? Orthopaedic Residents Versus ChatGPT. *Clin Orthop Relat Res*. 2023;481(8):1623-  
269 30.
  - 270 11. Karimov Z, Allahverdiyev I, Agayarov OY, Demir D, Almuradova E. ChatGPT vs UpToDate:  
271 comparative study of usefulness and reliability of Chatbot in common clinical presentations of  
272 otorhinolaryngology-head and neck surgery. *Eur Arch Otorhinolaryngol*. 2024;281(4):2145-51.
  - 273 12. Balel Y. Can ChatGPT be used in oral and maxillofacial surgery? *J Stomatol Oral*  
274 *Maxillofac Surg*. 2023;124(5):101471.
  - 275

# Comparing Otolaryngology Board Exam Scores Between Medical Trainees and ChatGPT

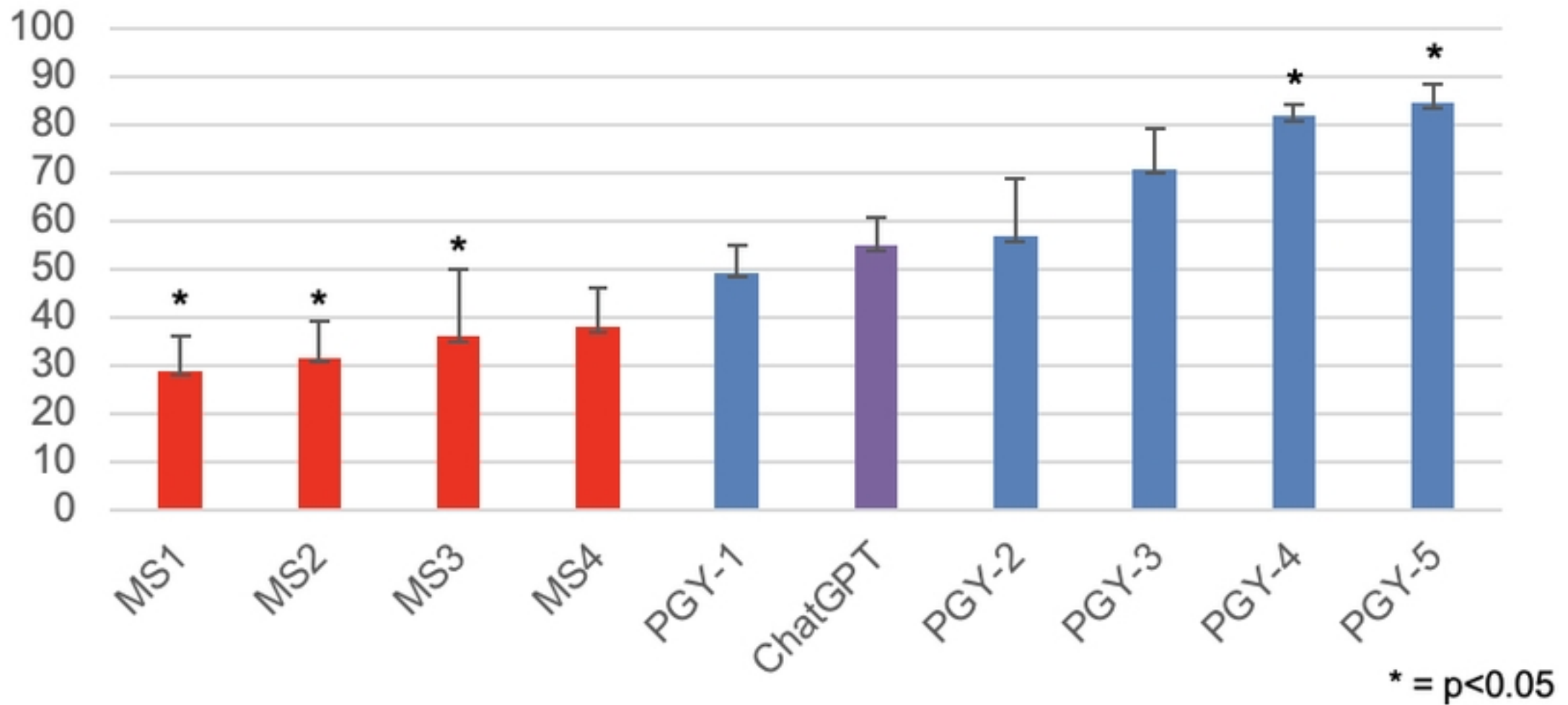


Fig 1