

1 **ChatGPT takes the FCPS exam in Internal Medicine**

2

3 Hina Qazi<sup>1</sup>, Syed Ahsan Ali<sup>2</sup>, Muhammad Irfan<sup>2</sup>, M. A. Rehman Siddiqui<sup>1</sup>

4

5 <sup>1</sup> Department of Ophthalmology and Visual Sciences,

6 <sup>2</sup> Department of Internal Medicine,

7 Aga Khan University Hospital,

8 Stadium Road, Karachi.

9

10 **Corresponding Author:**

11 M. A. Rehman Siddiqui

12 Email: [rehman.siddiqui@gmail.com](mailto:rehman.siddiqui@gmail.com)

13

14

15 HQ drafted the manuscript, collected the data, and statistically analyzed it.

16 AA and MI critically evaluated the manuscript.

17 RS conceptualized the study design and critically revised the manuscript.

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

## ABSTRACT

Large language models (LLMs) have exhibited remarkable proficiency in clinical knowledge, encompassing diagnostic medicine, and have been tested on questions related to medical licensing examinations. ChatGPT has recently gained popularity because of its ability to generate human-like responses when presented with exam questions. It has been tested on multiple undergraduate and subspecialty exams and the results have been mixed. We aim to test ChatGPT on questions mirroring the standards of the FCPS exam, the highest medical qualification in Pakistan.

We used 111 randomly chosen MCQs of internal medicine of FCPS level in the form of a text prompt, thrice on 3 consecutive days. The average of the three answers was taken as the final response. The responses were recorded and compared to the answers given by subject experts. Agreement between the two was assessed using the Chi-square test and Cohen's Kappa with 0.75 Kappa as an acceptable agreement. Univariate regression analysis was done for the effect of subspecialty, word count, and case scenarios in the success of ChatGPT. Post-risk stratification chi-square and kappa statistics were applied.

ChatGPT 4.0 scored 73% (69%-74%). Although close to the passing criteria, it could not clear the FCPS exam. Question characteristics and subspecialties did not affect the ChatGPT responses statistically. ChatGPT shows a high concordance between its responses indicating sound knowledge and a high reliability.

This study's findings underline the necessity for caution in over-reliance on AI for critical clinical decisions without human oversight. Creating specialized models tailored for medical education could provide a viable solution to this problem.

Keywords:

ChatGPT, artificial intelligence, internal medicine, large language models, medical exams

## 36 **Author Summary**

37 Artificial intelligence is the future of the world. Since the launch of ChatGPT in 2014, it become one of the most widely used application for people in all fields  
38 of life. A wave of excitement was felt among the medical community when the chatbot was announced to have cleared the USMLE exams. Here, we have tested  
39 ChatGPT on MCQs mirroring the standard of FCPS exam questions. The FCPS is the highest medical qualification in Pakistan. We found that with a vast data  
40 base, ChatGPT could not clear the exam in all of the three attempts taken by it. ChatGPT, however, scored a near passing score indicating a relatively sound  
41 knowledge.

42 We found ChatGPT to be a consistent LLM for complex medical scenarios faced by doctors in their daily lives irrespective of the subspecialty, length or word  
43 count of the questions. Although ChatGPT did not pass the FCPS exam, its answers displayed a high level of consistency, indicating a solid understanding of  
44 internal medicine. This demonstrates the potential of AI to support and improve medical education and healthcare services in near future.

45

46 **Introduction:**

47 Artificial Intelligence (AI) is increasingly establishing its role in a variety of specialties, including the medical field. It has revolutionized the human approach to  
48 various tasks and problems. In recent years, AI has been tried by medical professionals in diagnosis, precision medicine, and research(1). One of the remarkable  
49 inventions of AI is ChatGPT (Chat Generative Pre-trained Transformer) which is a natural language processing (NLP) Chatbot driven by AI technology. It is a  
50 state-of-the-art language model designed to generate human-like text based on the input it receives. This large language model (LLM) can answer a variety of  
51 questions and reply to ‘prompts’ on request. ChatGPT was released in November 2022 and crossed the one million user mark in just five days after it was made  
52 public. (2). It has gained a high surge of interest in the medical field in the last few months, with endless speculations about its predictive capabilities, and at least  
53 10 authorships in peer-reviewed scientific journals(3).

54 AI has been indicated as a safe and effective tool for triage in ER demonstrating high sensitivity in deciding the need for patient admission and providing  
55 comprehensive diagnoses, in addition to offering treatment strategies(4). ChatGPT exhibits promising potential in contributing to clinical decision-making skills.  
56 A recent study by Rao et al. reported clinical accuracy of up to 71.7% for ChatGPT in making clinical judgments, diagnosis, and management decisions.  
57 However, it lagged in listing appropriate differential diagnoses which are the essence of a physician’s role (5). It has been studied as a feasible tool for radiologic  
58 decision-making, with the potential to improve clinical workflow and responsible use of radiology services. (6)

59 Recently, more healthcare consumers have turned to the internet to seek health-related advice. It is of a major concern whether the information accessed is  
60 accurate and comparable to that of a physician (7). When comparing the safety of responses of ChatGPT to eye-related patient queries to that of qualified  
61 ophthalmologists, the AI did not differ much from humans in the likelihood of causing harm(8).

62 ChatGPT has rapidly demonstrated its ability to answer exam-style questions and provide explanations, raising questions about its potential role in education and  
63 assessment. (9) There have been several recent studies in which the AI chatbot has taken various exams in different specialties. It is seen to be capable of  
64 achieving a passing grade when tested on exams such as the United States Medical Licensing Examination and the European Exam in Core Cardiology (10, 11).

65 ChatGPT was also able to pass BLS and ACLS exams with the questions being presented as open-ended questions with outstanding results. However, it failed  
66 when the same questions were presented in a multiple-choice question format(12). These facts taken together make a compelling case for the potential  
67 applications of ChatGPT as an interactive medical education tool to support learning. However, there are some studies indicating that it did not perform well in  
68 exams like Taiwan's Family Medicine Board Exam and AHA MCQs(12, 13). When pitted against practice questions of specialized board exams like the  
69 neonatal-perinatal medicine board examination, which is taken by practicing pediatricians specializing in neonatology, ChatGPT managed to score only 46%.  
70 These suggest further testing of the AI model before incorporating it into medical education and clinical workflows. (14)

71 The data on the performance of ChatGPT in examinations of low-income countries, with different ethnic and cultural backgrounds, is scarce. This study tested the  
72 ability of ChatGPT on exam-style questions that mirror the standards of FCPS theory examinations conducted by CPSP. FCPS examinations are given to  
73 candidates enrolled in Fellowship of College of Physicians and Surgeons (FCPS) programs after completion of their training. The theory examination consists of  
74 two parts i.e. Paper 1 and Paper 2 with 100 multiple choice questions (MCQs) each covering the thorough knowledge of internal medicine, testing  
75 pathophysiology, clinical reasoning, and guideline-recommended medical management. The passing criteria is 75%. This study is designed to assess the clinical  
76 usefulness of ChatGPT(4.0) and evaluate its performance on questions that lie within the scope of the final exit exam of Fellowship Of College of Physicians and  
77 Surgeons (FCPS part 2) examinations conducted by the College of Physicians and Surgeons (CPSP).

78

79

80

81

82

83 **Methods and materials:**

84 **ChatGPT-4**

85 This comprehensive study was designed to test ChatGPT (OpenAI; San Francisco, CA)'s performance on questions of FCPS part 2 level in the specialty of  
86 internal medicine. We utilized the latest version of ChatGPT, which is based on the GPT-4 architecture, and is an open-access LLM developed by OpenAI  
87 (<https://openai.com>), launched in 2023. As an LLM, it is designed to generate human-like responses. This AI chatbot has been trained on extensive data,  
88 including medical texts and journals, up to September 2021.

89 **FCPS 2 exam questions:**

90 The FCPS part 2 theory exam consists of multiple-choice questions with 5 choices for each question. As a source of MCQs, publically available exam  
91 questions of FCPS part 2 level were used. A total of 111 single-answer MCQs which reflected the standards of the FCPS 2 exam were obtained and screened.  
92 Approximately ten questions from each subspecialty namely endocrinology, cardiology, nephrology, neurology, pulmonology, infectious diseases, hematology,  
93 and oncology were selected. Success was defined as ChatGPT achieving the passing criteria of 75% or higher.

94 **Input to ChatGPT:**

95 The selected questions were manually pasted into ChatGPT, preceded by the statement, "Please choose the best answer and explain your reasoning:" in a new  
96 separate line. The multiple-choice answers were provided, with each option pasted in a separate new line. A new chat session was started each time to avoid  
97 retention bias. About the output from ChatGPT, responses that were incorrect or inconclusive were both considered incorrect for this study. For questions for  
98 which ChatGPT provided 2 answers, another prompt was given to choose one best answer, and that answer was taken as the final answer. To evaluate  
99 ChatGPT's performance, each question from the question bank was submitted in the form of a text prompt. To evaluate the consistency and stability of these  
100 responses, all questions were submitted to ChatGPT three times on different dates from 15 November to 20 November 2023. The average of those three

101 responses was taken as the final answer. All questions containing visual elements such as clinical images, charts, and tables were excluded. Upon receiving the  
102 responses generated by ChatGPT, they were scored against the correct answers provided by the subject experts. Responses from ChatGPT were compared with  
103 the answers provided by the subject experts along with their clinical justification. Two subject experts answered these questions in a masked manner. For the  
104 questions in which there were differences in opinion, the two subject experts discussed the responses to reach a consensus. Answers on which the subject experts  
105 did not have a consensus, were planned to be excluded.

### 106 **Statistical Analysis**

107 The statistical analysis was done using SPSS version 23. ChatGPT responses of each attempt and the final key were compared to assess agreement between the  
108 two variables using the Chi-square test and Cohen's Kappa. The Cohen's kappa value of 0.75 was considered an acceptable agreement. Univariate regression  
109 analysis was done for the effect of subspeciality, word count, and case scenarios in the success of ChatGPT. Stratified analysis was done for subspecialties and  
110 for answers on which physicians initially agreed. Post-risk stratification chi-square and kappa statistics were applied. A p-value less than 0.05 was considered  
111 statistically significant at the confidence interval of 95%.

### 112 **Ethical considerations**

113 As this study did not involve any human or animal interaction, no approval from the institutional review board or informed consent was required.

114

115

116

117

118 Flowchart:

119

120

121

122

123

124

125

126

127

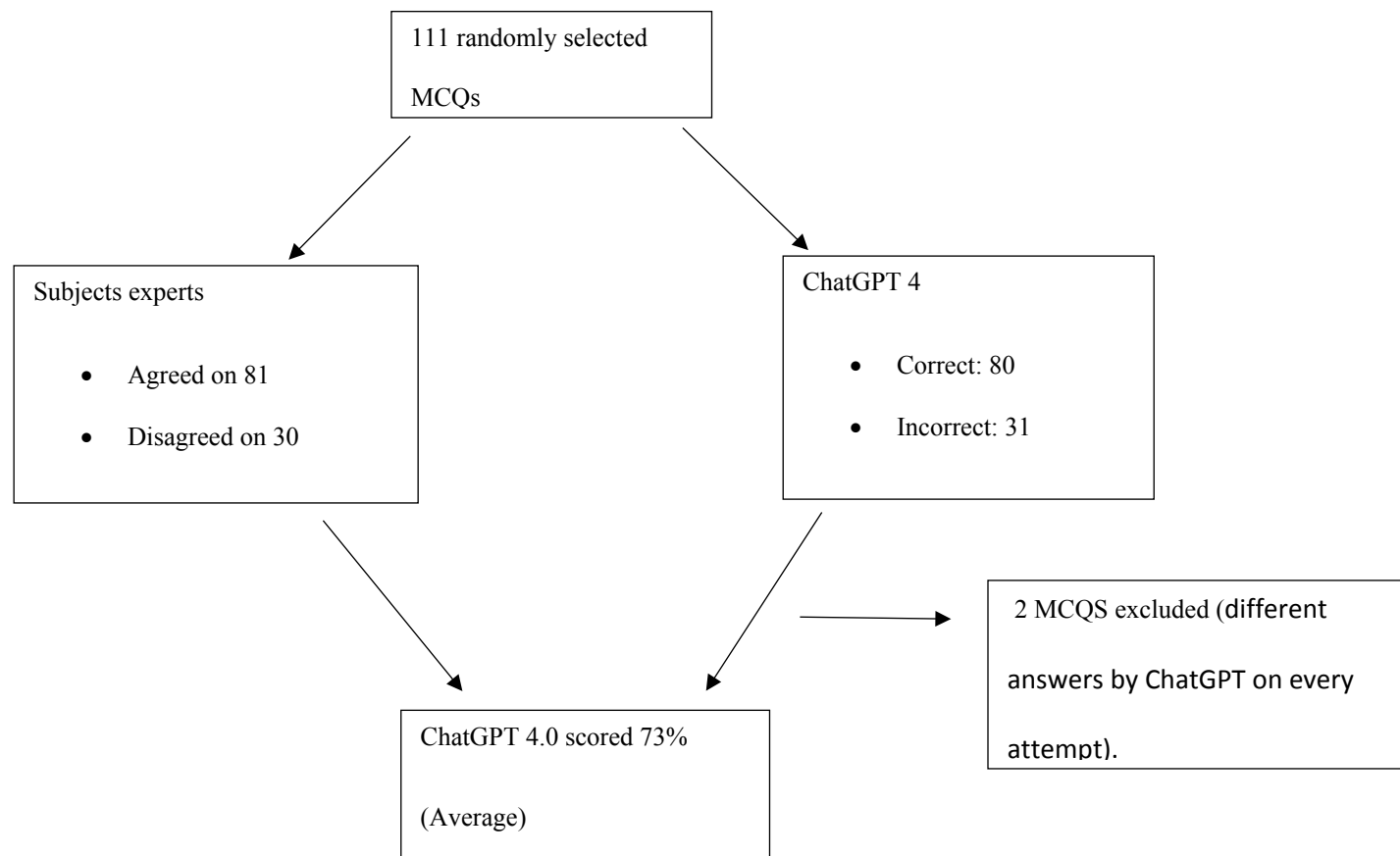
128

129

130

131

132





133 **Results:**

134 This study tested the performance of ChatGPT 4.0 on questions that fall within the scope of the FCPS 2 theory examination in internal medicine. Hundred and  
135 eleven publically available questions, comprising various subspecialties were tested on ChatGPT (Table 1). Case scenarios comprised a major portion of the  
136 exam i.e. 83 questions, while the rest were clinical facts. None of the questions were related to the local policies. The word count of the questions ranged from 8  
137 words to 241 words, with a median of 72 words. With the median as a reference, questions with a word count less than 72 were labeled as short questions and  
138 those with more than 72 words were labeled as long questions. Out of 111 questions, subject experts gave similar answers to 81 questions while they differed on  
139 30 questions, indicating the level of difficulty of these questions. They later reached a consensus for every answer after discussion and no disagreement was  
140 reported. Their responses were taken as the final key. Two questions were excluded from the statistical analysis as ChatGPT came up with a different answer  
141 every time and their average could not be calculated. ChatGPT took three attempts on all questions. In the first attempt. ChatGPT scored 74% (Kappa 0.674).In  
142 all three attempts, it could not cross the 75% passing criteria, Of the 81 questions on which the subject experts initially agreed, ChatGPT scored 74%. On the  
143 questions, they agreed to after discussion, ChatGPT scored 73%. ChatGPT could not clear its second attempt too, and scored 72.0 % ( Kappa 0.672). ChatGPT  
144 failed the third attempt with a score of 69.7% (Kappa 0.630). The average of the three attempts was taken as the final score and was 73% (Kappa 0.675)(p 0.001).

145

146

147

148

149

150

<b>Subspecialties</b>	<b>Frequency</b>	<b>Percent</b>
Nephrology	9	8.3
Neurology	12	11.1
Cardiology	10	9.3
Dermatology	1	0.9
Endocrinology	13	12
Forensic	1	0.9
Gastroenterology	14	13
Hematology/oncology	15	13.9
Infectious Diseases	9	8.3
Pulmonology	11	10.2
Rheumatology	13	12
<b>Questions on physicians Initially Agreed</b>		
No	29	26.9
Yes	79	73.1
<b>Question length</b>		
Short	53	49.1
Long	55	50.9
<b>Case scenarios</b>		
No	25	23.1
Yes	83	76.9
<b>Basic medical facts</b>		
No	85	78.7
Yes	23	21.3

151

Total	108	100.0
-------	-----	-------

152

153

154 *Table 1: Qualitative analysis of question characteristics.*

155

156 Statistically, the consistency between attempts was seen by assessing the inter-attempt Kappa agreement and the chi-square test. The first attempt when compared  
157 to the second attempt showed an almost perfect agreement of 0.99 or 99%. The second attempt when compared to the third attempt gave an agreement of 0.92 or  
158 92%. The third attempt when compared to the first attempt showed concordance of 88% or an agreement of 0.88.

159 On univariate regression analysis, no statistical difference was found concerning subspecialty, word count, and case scenarios in the success of ChatGPT (p  
160 0.175). Therefore, binary logistic regression analysis could not be done. Stratified analysis was done for subspecialties and for answers on which physicians  
161 initially agreed. Post-risk stratification chi-square and kappa statistics were applied. Pulmonology had the highest agreement of 100% while nephrology had the  
162 lowest of 8.3% but this difference was statistically non-significant. The qualitative analysis of the responses by ChatGPT is given in Table 2..

163

164

			Final Key by subject experts					Total
			Option A	Option B	Option C	Option D	Option E	
<b>Answers by ChatGPT (Average of 3 attempts)</b>	<b>Option A</b>	Percentages within an average attempt	84.2%	10.5%	5.3%	0.0%	0.0%	100.0%
	<b>Option B</b>	Percentages within an average attempt	10.3%	65.5%	3.4%	6.9%	13.8%	100.0%
	<b>Option C</b>	Percentages within an average attempt	9.1%	0.0%	86.4%	0.0%	4.5%	100.0%
	<b>Option D</b>	Percentages within an average attempt	22.2%	0.0%	0.0%	66.7%	11.1%	100.0%
	<b>Option E</b>	Percentages within an average attempt	10.0%	10.0%	10.0%	0.0%	70.0%	100.0%
<b>Total percentages</b>			25.0%	21.3%	21.3%	13.0%	19.4%	100.0%

165 *Table 2: Qualitative analysis of responses generated by ChatGPT.*

166 **Discussion:**

167 ChatGPT garnered significant media attention and generated a buzz of excitement when it was reported to have cleared the USMLE exam, which is a licensing  
168 exam taken by medical undergraduates(10). ChatGPT was also reported to have passed multiple international licensing exams and sub-specialty exams (table  
169 3and table 4)(Fig.1). Along with enthusiasm, it has sparked skepticism about its accuracy, reliability, and its potential in the field of medicine.

170

Authors	Study date	Country	Journal	Exam	Specialty	Question type	Results
Hopkins et al.	March 2023	US	Journal of Neurosurgery	Congress of Neurological Surgeons (CNS) Self-Assessment Neurosurgery (SANS)	Neurology	MCQs	Passed
Birkett et al.	May 2023	UK	Br J Anaesth	FRCA	Anesthesia	MCQs	Failed
Shay et al.	May 2023	UK	Br J Anaesth	Anaesthesiology board examination questions	Anesthesia	MCQs	Passed
Rohaid et al	December 2023	USA	Neurosurgery	Self-Assessment Neurosurgery Examinations (SANS) American Board of Neurological Surgery Self-Assessment Examination I	Neurosurgery	MCQs	Passed
Mihalache et al.	April 2023	Canada	JAMA	OphthoQuestions free trial for ophthalmic board certification preparation	Ophthalmology	MCQs	Failed
Lem et al.	Aug 2023	USA	Clin Orthop Relat Res	American Board of Orthopaedic Surgery Examination ( In training)	Orthopaedics	MCQs	Failed
Suchman et al.	May 2023	USA	Am J Gastroenterol	American College of Gastroenterology Self-Assessment Test	Gastroenterology	MCQs	Failed
Bhayana et al.	April 2023	Canada	Radiology	Canadian Royal College and American Board of Radiology examinations	Radiology	MCQs	Failed
Skalidis et al.	May 2023	Switzerland	Eur Heart J Digit Health	European Exam in Core Cardiology	Cardiology	MCQs	Passed
Passby et al.	June 2023	UK	Clin Exp Dermatol	Dermatology Specialty Certificate Examination	Dermatology	MCQs	Passed
Beam et al.	July 2023	USA	JAMA Pediatr	Neonatal Board Examination	Neonatology	MCQs	Failed
Weng et al.	Aug 2023	Taiwan	J Chin Med Assoc	Family Medicine Board Exam	Family Medicine	MCQs	Failed
Teegbay et al.	Sept 2023	USA	J Acad Ophthalmol (2017	OKAP	Ophthalmology	MCQs	Passed
Gencer et al.	Aug 2023	Turkey	Am J Med Sci	Turkish-language thoracic surgery exam	Thoracic surgery	MCQs	Passed
Sahin et al.	Dec 2023	Turkey	Comput Biol Med.	Turkish Neurosurgical Society Proficiency Board Exams (TNSPBE)	Neurosurgery	MCQs	Passed
Kufel et al.	Sept.2023	Poland	Pol J Radiol.	Radiology NSE	Radiology	MCQs	Failed
Kinoshita et al.	Oct 2023	Japan	J Anaesth	JSA-certified anesthesiologist exam	Anesthesia	MCQS	Failed
Huynh et al.	July 2023	USA	Urol Pract.	2022 Self-assessment Study Program	Urology	MCQs + open-ended	Failed
Panthier et al.	Sept 2023	France	J Fr Ophthalmol.	European Board of Ophthalmology examination (French version)	Ophthalmology	MCQs, short answer questions, true and false questions, clinical scenarios, and theoretical questions	Passed

172 Table 3: List of sub-specialty exams taken by ChatGPT

174

176

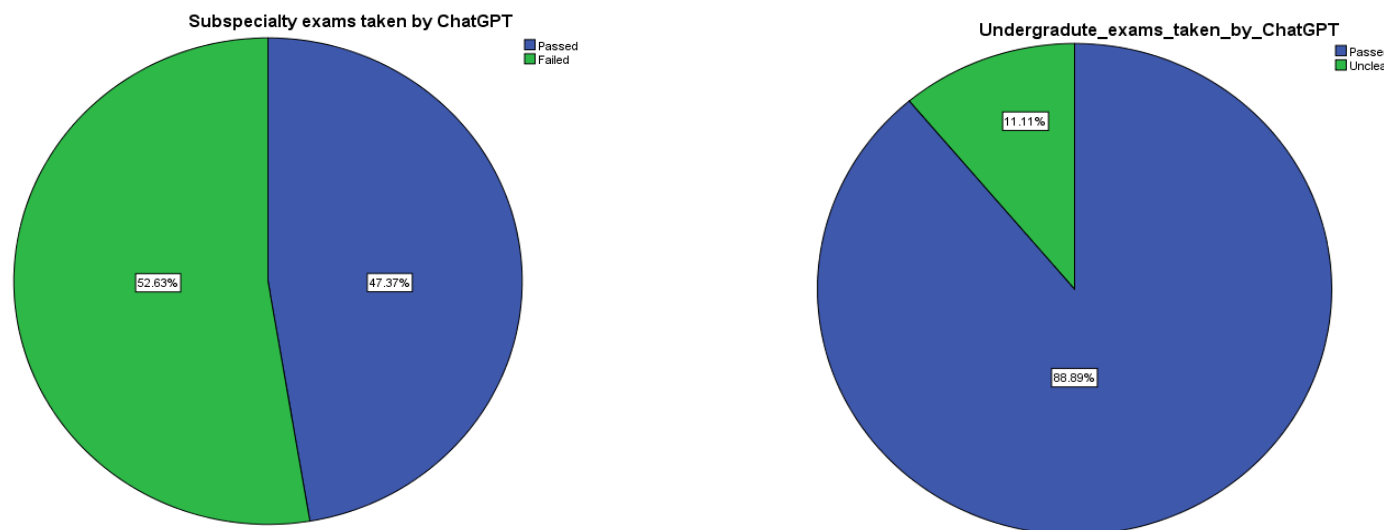
Authors	Study date	Country	Journal	Exam	Question type	Results
Gilson et al.	Feb 2023	USA	JMIR Med Educ	USMLE	MCQs	Pass
Kung et al.	Feb 2023	USA	PLOS Digit Health	USMLE	MCQs	Pass
Aljindan et al.	Sept 2023	Saudia	Cureus	Saudi Medical Licensing Exam	MCQs	Pass
Wang et al.	Aug 2023	China	J Med Syst	Chinese National Medical Licensing Examination	MCQs	Below performing students
Huang et al.	Oct 2023	Taiwan	Healthcare (Basel)	Registered Nurse License Exam	MCQs	Pass
Rosol et al.	Nov 2023	Poland	Sci Rep.	Polish Medical Final Examination	MCQs	Pass
Kasai et al.	April 2023	Japan	Springer	Japanese Medical Licensing Examinations	MCQs	Pass
Lai et al.	Sept 2023	UK	Front Med (Lausanne).	UKMLA	MCQs	Pass
Bonetti et al.	July 2023	Italy	Ann Biomed Eng.	Italian Residency Admission	MCQs	Pass

Table 4.:List of Undergraduate exams taken by ChatGPT



177

178



179

*Fig 1: Pie chart indicating the percentages of passed and failed exams taken by ChatGPT*

180

181 The results of this study indicate that while ChatGPT demonstrates a reasonable level of knowledge in internal medicine, as evidenced by its nearly passing  
182 scores in the FCPS 2 examinations, it does not yet meet the benchmark for clinical reliability. The FCPS exam is a highly specialized exam taken by postgraduate  
183 trainees after the completion of their 4-year training, to be a qualified consultant in their field. Other studies have reported similar results, where ChatGPT failed  
184 to clear specialty exams (Table 2). For example, ChatGPT failed the Taiwan Family Medicine exam (13). There could be two reasons as to why it failed. Either it  
185 could not accurately answer family medicine-related questions or it was unable to respond correctly to prompts in Chinese language. ChatGPT could also not

186 perform well on the American Academy of Gastroenterology exam questions scoring 62.4 percent compared to the required 70 percent to pass the exams(16).  
187 Similarly, ChatGPT showed poorer responses on American Board of Orthopedic Surgery style questions when compared with responses from in-training  
188 residents, scoring only 47% (17).

189 Assessing the performance of ChatGPT on FCPS questions, it consistently failed all three attempts, with the lowest score on the last attempt. Eighty-three MCQs  
190 in this study were based on clinical scenarios and ChatGPT answered 62 of them correctly, indicating a relatively sound knowledge of diagnosing medical  
191 conditions, interpreting diagnostic testing and lab reports providing appropriate management plans. Question characteristics such as word count, case scenarios,  
192 and clinical facts statistically did not seem to affect the performance of ChatGPT. This finding is in line with Rohaid Ali et al who tested ChatGPT-4 on  
193 neurosurgery questions and reported no effect of responses on the word count of the questions(18).

194 There could be multiple explanations as to why ChatGPT failed the FCPS exam. ChatGPT was never trained specifically on medical literature and was  
195 developed as a general interactive LLM. The AI model's training on historical data up to September 2021 may not encompass the most recent medical guidelines  
196 and research, and its inability to access paid content of medical journals could have led to less optimal responses. ChatGPT gathers data from various sources,  
197 some of which may be nonmedical, semi-medical, or outdated, resulting in inaccurate responses. The way it works is to predict the next most suitable word in a  
198 string, generating a likely reply based on existing data, without any regard to factual accuracy. The model lacks inherent comprehension of any subject or matter.  
199 ChatGPT may have struggled with the MCQ-style format of the FCPS exam. As a chatbot, it may be more suited to answering open-ended questions rather than  
200 being given a set of options to choose from. A similar finding was also speculated by Zhu et al.(12). In their study, ChatGPT could clear the open-ended  
201 American Heart Association questions but failed on the MCQs of the same. FCPS exam includes a higher proportion of questions that require more nuanced  
202 clinical judgment or interpretation, particularly in areas where there might be multiple acceptable approaches, amongst which the trainee has to choose the best.  
203 ChatGPT might have found it more challenging.

204 The impact of ChatGPT in the educational sector has received mixed reactions. While some appreciate the AI's ability to provide useful insights and demonstrate  
205 reasoning skills to students, others point out issues such as the production of inaccurate information, the challenge in interpreting responses, the possibility of  
206 bias, and ethical concerns(19).

207 There are certain limitations of our study. This study only assessed ChatGPT's performance on MCQs as the theory exam of FCPS only consists of MCQs. This  
208 format might not fully capture the AI's capabilities in other aspects of clinical reasoning, such as case analysis, patient interaction, and practical skills. Secondly,  
209 the determination of correct answers was based on the consensus of subject experts. While this method is standard, it introduces a subjective element and could  
210 potentially overlook the diversity of acceptable clinical opinions or practices.

211 The study suggests that future iterations of AI models like ChatGPT 4.0 could benefit from more targeted training, specifically in areas where the model  
212 currently underperforms. As we write this, AI can complement, but not replace human expertise. It can also be said that shortly, LLMs when used by internists  
213 will make them better internists. On the contrary, some people believe that in the next 10 years, AI chatbots will be primarily making disease diagnoses and  
214 doctors will only be reached out for a second opinion. However, this study's findings underline the necessity for caution in over-reliance on AI for critical clinical  
215 decisions without human oversight. Regular updates with recent medical literature and guidelines are crucial to ensure the model's responses remain current and  
216 clinically relevant.

217

218

219

220

221

222 **Conclusion:**

223 This cross-sectional study tested randomly selected 111 FCPS-level MCQs on ChatGPT 4.0 on three consecutive attempts. Question characteristics such as word  
224 count, case scenarios, and clinical facts did not seem to affect the responses. Although ChatGPT was not able to pass the FCPS exam, it showed a high  
225 concordance within its answers indicating relatively sound knowledge of internal medicine and reflecting the potential of AI in assisting and enhancing medical  
226 education and healthcare services. We advise caution for those using ChatGPT as a medical education tool. As the advancements in AI technology continue,  
227 particularly in areas of clinical interpretation and specific-domain applications of knowledge, it will be interesting to see how this technology continues to  
228 improve and how it might best be applied in medical education. Creating specialized models tailored for medical education could provide a viable solution to this  
229 problem.

230

231

232

234 **Acknowledgment:**

235 The authors have no financial disclosures or support in this work.

236 We would like to thank Ms. Khadijah Abid, senior instructor of research at the Department of Ophthalmology and Visual Sciences, AKUH, for helping with the  
237 statistical analysis of this study.

238

239

240

241

242 **References:**

- 243 1. Topol EJ. High-performance medicine: the convergence of human and artificial intelligence. *Nat Med*. 2019;25(1):44-56.
- 244 2. Open AI. 2023 [cited 2023]. Available from: <https://openai.com/>.
- 245 3. Castelvechi D. Are ChatGPT and AlphaCode going to replace programmers? *Nature*. 2022.
- 246 4. Gebrael G, Sahu KK. Enhancing Triage Efficiency and Accuracy in Emergency Rooms for Patients with Metastatic Prostate Cancer: A  
247 Retrospective Analysis of Artificial Intelligence-Assisted Triage Using ChatGPT 4.0. 2023;15(14).
- 248 5. Rao A, Pang M. Assessing the Utility of ChatGPT Throughout the Entire Clinical Workflow: Development and Usability Study.  
249 2023;25:e48659.
- 250 6. Rao A. Evaluating ChatGPT as an Adjunct for Radiologic Decision-Making. 2023.
- 251 7. Zakar R, Iqbal S, Zakar MZ, Fischer F. COVID-19 and Health Information Seeking Behavior: Digital Health Literacy Survey amongst  
252 University Students in Pakistan. 2021;18(8).
- 253 8. Bernstein IA, Zhang YV, Govil D, Majid I, Chang RT, Sun Y, et al. Comparison of Ophthalmologist and Large Language Model Chatbot  
254 Responses to Online Patient Eye Care Questions. *JAMA Netw Open*. 2023;6(8):e2330320.
- 255 9. Stokel-Walker C. AI bot ChatGPT writes smart essays - should professors worry? *Nature*. 2022.
- 256 10. Kung TH, Cheatham M, Medenilla A, Sillos C, De Leon L, Elepaño C, et al. Performance of ChatGPT on USMLE: Potential for AI-  
257 assisted medical education using large language models. 2023;2(2):e0000198.
- 258 11. Panthier C, Gatinel D. Success of ChatGPT, an AI language model, in taking the French language version of the European Board of  
259 Ophthalmology examination: A novel approach to medical knowledge assessment. *J Fr Ophtalmol*. 2023.
- 260 12. Zhu L, Mou W, Yang T, Chen R. ChatGPT can pass the AHA exams: Open-ended questions outperform multiple-choice format.  
261 Resuscitation. 2023;188:109783.
- 262 13. Weng TL, Wang YM, Chang S, Chen TJ, Hwang SJ. ChatGPT failed Taiwan's Family Medicine Board Exam. *J Chin Med Assoc*.  
263 2023;86(8):762-6.
- 264 14. Beam K, Sharma P, Kumar B, Wang C, Brodsky D, Martin CR, et al. Performance of a Large Language Model on Practice Questions for  
265 the Neonatal Board Examination. *JAMA Pediatr*. 2023;177(9):977-9.
- 266 15. Aljindan FK, Al Qurashi AA, Albalawi IAS, Alanazi AMM, Aljuhani HAM, Falah Almutairi F, et al. ChatGPT Conquers the Saudi  
267 Medical Licensing Exam: Exploring the Accuracy of Artificial Intelligence in Medical Knowledge Assessment and Implications for Modern  
268 Medical Education. *Cureus*. 2023;15(9):e45043.
- 269 16. Suchman K, Garg S, Trindade AJ. Chat Generative Pretrained Transformer Fails the Multiple-Choice American College of  
270 Gastroenterology Self-Assessment Test. *Am J Gastroenterol*. 2023.
- 271 17. Lum ZC. Can Artificial Intelligence Pass the American Board of Orthopaedic Surgery Examination? Orthopaedic Residents Versus  
272 ChatGPT. *Clin Orthop Relat Res*. 2023;481(8):1623-30.
- 273 18. Ali R, Tang OY. Performance of ChatGPT and GPT-4 on Neurosurgery Written Board Examinations. 2023;93(6):1353-65.
- 274 19. Kasneci E, Seßler K, Küchemann S, Bannert M, Dementieva D, Fischer F, et al. ChatGPT for good? On opportunities and challenges of  
275 large language models for education. Learning and individual differences. 2023;103:102274.

