

# Can ChatGPT-4o really pass medical science exams? A pragmatic analysis using novel questions.

1 Philip M. Newton\*, Christopher J. Summers, Uzman Zaheer, Maira Xiromeriti, Jemima R.  
2 Stokes, Jaskaran Singh Bhangu, Elis G. Roome, Alanna Roberts-Phillips, Darius Mazaheri-  
3 Asadi, Cameron D. Jones, Stuart Hughes, Dominic Gilbert, Ewan Jones, Keioni Essex, Emily  
4 C. Ellis, Ross Davey, Adrienne A. Cox and Jessica A. Bassett.

5 Swansea University Medical School, Swansea, Wales, United Kingdom, SA2 8PP.

6 **\*Correspondence:**

7 Corresponding Author; [p.newton@swansea.ac.uk](mailto:p.newton@swansea.ac.uk)

8 ORCID IDs:

9 PMN <https://orcid.org/0000-0002-5272-7979>

10 CJS <https://orcid.org/0009-0000-5336-2492>

11 UZ <https://orcid.org/0009-0008-2148-1532>

12 MX <https://orcid.org/0000-0002-2975-184X>

13 JRS <https://orcid.org/0000-0003-2623-0245>

14 EGR <https://orcid.org/0009-0009-5845-4164>

15 DMA <https://orcid.org/0009-0002-7999-3123>

16 ECE <https://orcid.org/0009-0005-6493-9337>

17 RD <https://orcid.org/0000-0001-9852-1653>

18 DG <https://orcid.org/0009-0002-0024-3662>

19 EJ <https://orcid.org/0009-0002-9221-1990>

20 JAB <https://orcid.org/0009-0002-2146-2987>

21 AAC <https://orcid.org/0000-0002-3902-3491>

22 **Keywords:** assessment validity, academic integrity, cheating, evidence-based education, MCQs,  
23 pragmatism

24

## 25 **Abstract**

26 ChatGPT apparently shows excellent performance on high level professional exams such as those  
27 involved in medical assessment and licensing. This has raised concerns that ChatGPT could be used  
28 for academic misconduct, especially in unproctored online exams. However, ChatGPT has also  
29 shown weaker performance on pictures with questions, and there have been concerns that ChatGPT's  
30 performance may be artificially inflated by the public nature of the sample questions tested, meaning  
31 they likely formed part of the training materials for ChatGPT. This led to suggestions that cheating  
32 could be mitigated by using novel questions for every sitting of an exam, and making extensive use  
33 of picture-based questions. These approaches remain untested.

34 Here we tested the performance of ChatGPT-4o on existing medical licensing exams in the UK and  
35 USA, and on novel questions based on those exams.

36 ChatGPT-4o scored 94% on the United Kingdom Medical Licensing Exam Applied Knowledge Test,  
37 and 89.9% on the United States Medical Licensing Exam Step 1. Performance was not diminished  
38 when the questions were rewritten into novel versions, or on completely novel questions which were  
39 not based on any existing questions. ChatGPT did show a slightly reduced performance on questions  
40 containing images, particularly when the answer options were added to an image as text labels.

41 These data demonstrate that the performance of ChatGPT continues to improve and that online  
42 unproctored exams are an invalid form of assessment of the foundational knowledge needed for  
43 higher order learning.

44

## 45 Introduction

46 New generative artificial intelligence (AI) tools such as ChatGPT have attracted enormous attention,  
47 in part for their apparent ability to pass high level professional exams, with the subscription version  
48 of ChatGPT, running GPT-4, scoring an average of 75% on MCQ-based exams across a variety of  
49 disciplines (1). This excellent performance is replicated on specific medical qualifying exams such as  
50 the United States Medical Licensing Exam (USMLE) Step 1 where it scored 86% (2) and the United  
51 Kingdom Medical Licensing Exam Applied Knowledge Test (UK MLA AKT) where it scored 76.3%  
52 (3). These exams test high level problem-solving and are designed to assess the application of core  
53 knowledge to clinical scenarios (4) and represent a broader principle wherein multiple choice  
54 questions can, if written appropriately, assess higher-order learning in a range of disciplines (5).

55 However there have been a number of responses and criticisms of the claim that ChatGPT is  
56 genuinely solving the problems presented in these questions, in part because this seems to lead  
57 logically onto the idea that ChatGPT is able to ‘reason’ which apparently it cannot (6). Instead, critics  
58 propose, tools like ChatGPT are more likely ‘regurgitating’ content which has been in their training  
59 materials (7), a proposal which is supported by the fact that many studies use sample papers which  
60 are in the public domain and have been for some time, for instance the USMLE sample paper cited  
61 above was published in 2021. This regurgitation is not proposed to be verbatim, but instead is,  
62 essentially, a paraphrasing of prior training materials in a way that resembles a student who is  
63 plagiarising a piece of text by changing key words but without understanding the meaning, and so  
64 occasionally getting things (very) wrong (8). Thus, the argument goes, part of the reason why LLMs  
65 can ‘pass’ exams is because of this ‘regurgitation’ of sample papers which have been in the public  
66 domain for some time, and so to counter these apparent threat of ChatGPT to exam security and  
67 integrity educators could use novel questions for each sitting of the exam (9). In addition, there have  
68 been efforts to map the features of exam questions which ChatGPT appears to struggle with,  
69 including an increase in the number of answer items, increasing language complexity or having  
70 multiple correct answers. However none of these appears to have any effect on the numbers of  
71 questions which ChatGPT can answer correctly (10).

72 Many early papers which tested the performance of ChatGPT on sample exams deliberately excluded  
73 questions containing images, on the basis that older versions of ChatGPT, even GPT-4, could not  
74 process these images. Thus, the reported performance of ChatGPT may be an over-estimation, since  
75 the percentage scored by ChatGPT uses a lower denominator once image-based questions are  
76 excluded (e.g. (11)). This also leads to proposals that educators could author ‘ChatGPT-proof’  
77 questions by including images, along with mathematical calculations and reasoning tests, which it is  
78 proposed that ChatGPT does not perform well at (6).

79 These issues are important in part because of wider questions about the security, but also the  
80 inclusivity and cost, of examinations. In particular the sorts of university-administered knowledge  
81 tests that form part of a STEM curriculum prior to assessment using formal licensing examinations.  
82 Online examinations are cheaper and more flexible than their in-person equivalents, but they  
83 potentially risk cheating; during the COVID-19 pandemic, cheating in online exams appeared to  
84 double, and more students reported cheating than not (12). One apparent solution to this problem is to  
85 increase the use of online proctoring/invigilation systems to monitor student behaviour. However,  
86 these then drive back up the cost of the exams, and the student experience of remote proctoring is  
87 poor, with concerns about privacy, fairness, inclusivity and cost (13,14). An alternative is to avoid  
88 the use of proctoring altogether. A high profile 2023 publication analysed exam performance data  
89 from the COVID lockdown and concluded that unproctored online exams are a ‘valid and

90 meaningful' way of measuring student learning (15), although this analysis has been challenged (16)  
91 and does not include a consideration of ChatGPT. Thus it is important to understand whether  
92 ChatGPT truly can pass exams, including novel questions with images, as part of a consideration  
93 about how best to deploy exams, online or in-person, proctored or not.

94 Pragmatism is a research paradigm which prioritises the asking of questions whose answers will be  
95 useful, rather than perhaps asking more academic or basic questions (17). If ChatGPT truly can pass  
96 high level STEM exams, even with novel questions containing images, then from a pragmatic  
97 standpoint this is important because it essentially settles any debate about whether these  
98 examinations cannot be conducted in an online, unproctored format. From the pragmatic perspective,  
99 it does not matter *how* ChatGPT is doing this, either by truly solving problems or through some  
100 sophisticated paraphrasing. There is a related pragmatic issue, which is that for most STEM subjects  
101 there is a core curriculum; a basic set of knowledge and skills which graduates must be able to  
102 demonstrate in order to graduate, and also to be able to apply knowledge to practice. This cumulative  
103 view of learning has a long history but remains prevalent today through the use of instruments such  
104 as Bloom's Taxonomy (18). In essence, we cannot expect students to undertake learning and practice  
105 at the higher levels of Blooms Taxonomy unless they have the core foundational knowledge to be  
106 applied to those higher levels. Thus educators need to assess that foundational knowledge first,  
107 before it is applied, particularly where there are safety concerns, e.g. for patients. However, it seems  
108 reasonable to propose that there are only so many ways that one can phrase the exam questions which  
109 might assess these core principles. This creates a risk that, if educators strive to write completely  
110 novel questions on every core topic for every exam sitting, just to thwart ChatGPT, then this will  
111 rapidly become impossible. These issues also have relevance for the proposed positive benefits of  
112 ChatGPT. It offers great promise as a tutoring tool for students who are preparing for exams (19) but  
113 educators and learners both need to be confident that the answers given are logical and reasonable  
114 (20).

115 Some of the controversy and discourse about the apparent ability of ChatGPT to pass and perform  
116 well (or not) on exams likely comes from the frequent updating of ChatGPT over a short timescale. A  
117 review of ChatGPT performance on exams from multiple disciplines found that the subscription  
118 version of ChatGPT running GPT-4 outperformed the free version running GPT-3 or 3.5, with the  
119 average difference being 25 percentage points (1). On May 13 2024 OpenAI, the creators of  
120 ChatGPT, released another update, entitled ChatGPT-4o, showing enhanced performance compared  
121 to GPT-4, particularly on the integration of text, visual and audio information (21). The performance  
122 of ChatGPT-4o on medical licensing exams has not yet been examined.

123 Here then we address the following research questions. It is important to be clear that the specific  
124 medical licensing-type exams used here are intended to be a model for STEM exams generally, given  
125 that they are written to a high standard and are aimed at problem-solving and the application of  
126 knowledge (4,5).

- 127 1. How well does ChatGPT-4o perform on sample medical licensing exams in the USA and UK
- 128 2. Is the performance of ChatGPT affected when these sample questions are rewritten into novel  
129 formats, but assessing the same core curricular concepts
- 130 3. How well does ChatGPT perform on completely novel medical-licensing type questions.

131

132



## 134 **Methods**

135 The following question sources were tested.

- 136 1. (Pilot) Wikiversity Fundamentals of Neuroscience Exam (22)
- 137 2. Sample paper 1, UK Medical Licensing Assessment Applied Knowledge Test (23)
- 138 3. USMLE Step 1 Sample paper (24)
- 139 4. Rewritten questions from 2+3
- 140 5. Completely Novel USMLE-style questions.

141 ***Rewriting of existing questions in the public domain.*** Each question from sources 1-3 was rewritten  
142 by a member of the research team. Each question was rewritten three times with each rewrite  
143 undertaken by a different research team member. Rewriting instructions were to create an original  
144 question but which assessed the same learning, specifically to ‘change as much as possible about the  
145 question, without changing the underlying learning. Change all the text where possible’. Suggestions  
146 of specific items to change include demographic details in the scenarios, answer options, answer  
147 order. Each team member was also provided with a summary of common issues found when writing  
148 USMLE-style questions (4) and asked to avoid any of the identified writing flaws. All rewritten items  
149 were checked for accuracy and originality by registered doctors (CS, RD) or a subject matter expert  
150 (PMN) and adjusted where necessary, for example if the revised question could be made even more  
151 different to the original question.

152 An initial pilot was undertaken using five questions on neuroscience from ‘Wikiversity’ website.  
153 These were considered ‘lower order’ questions, assessing basic factual knowledge of neurological  
154 disease. The questions have been in the public domain since 2013. Each question was rewritten into  
155 three different forms by a member of the research team, who then discussed the process and  
156 feasibility of scaling the methodology to a larger exam. All four versions of each question were then  
157 pilot tested using GPT-4 on 23/04/24 and 24/04/24.

158 ***Analysis of existing medical licensing exams and rewrites.*** Each question was tested using a single  
159 shot method in the way that would be expected to be the most likely approach taken by a student who  
160 was seeking to cheat on an MCQ exam, i.e. the text was highlighted in the pdf (original questions) or  
161 word document (rewrites), copied and then pasted directly into ChatGPT-4o with no attempt to  
162 format the text. Where the question included a picture, this was copied using screen clipping, saved  
163 and uploaded as a .png file with only the country and the question number as the filed name (e.g.  
164 ‘UK32’). No additional prompts were given apart from the content of the question. Each question  
165 was asked in a new chat and no memory functions were activated. For the USMLE questions, a  
166 ‘temporary chat’ was activated for each question. No responses were given to ChatGPT. ChatGPT’s  
167 first response was recorded each time as correct/incorrect. ChatGPT-4o tests were undertaken May  
168 14-24 2024.

169 ***Creation and analysis of novel questions.*** Two sets of completely novel questions were generated,  
170 totalling 90 questions in all. A first set of forty novel questions were created in the style of questions  
171 for the UK MLA AKT and USMLE, by an author who is experienced in the creation of these  
172 assessment items (CS), according to guidance from the United States National Board of Medical  
173 Examiners (4). Ten of these questions included novel images that were either created for this study or  
174 were images from the private collection of one of the authors (CS). None of these images are  
175 available in the public domain. All images were obtained with appropriate consent and anonymised  
176 prior to use in keeping with paragraph 10 of the GMCs professional standards on making and using

177 visual recordings of patients (25). These questions were mapped to curricula items from the MLA  
178 content map (26) and were of a comparative style and difficulty to the MLA. A second set of  
179 questions was written by an author (PMN) using guidance for the creation of multiple-choice  
180 questions which assess higher order learning in STEM. These guidelines include identifying assumed  
181 knowledge, creating problem-solving scenarios and the use of actions as answer options (5). Some of  
182 these questions included images sourced from Wikimedia Commons. During this process the authors  
183 observed a trend that ChatGPT appeared to struggle with anatomical images that had novel text  
184 labels, e.g. a brain section with the labels A-H added, with arrows to specific brain regions that  
185 corresponded to question answers. To probe this further, an additional set of questions was generated  
186 so that there were a total of 14 pairs of questions which assessed the same learning but either using a  
187 labelled image, or text equivalent. Finally, ChatGPT was then asked simply to identify the labels on  
188 the images from these questions where possible. Each question was asked in a new 'temporary chat'.  
189 ChatGPT-4o tests were undertaken May 24-Jun 18 2024.

190

## 191 Results

192 **Summary.** We tested a total of 705 assessment items, of which ChatGPT answered 635 (90%)  
193 correctly. 111 of these questions contained images, of which ChatGPT answered 76 (68.5%)  
194 correctly. A breakdown of these items is below.

195 **Wikiversity Pilot.** GPT-4 correctly answered all versions of all questions, both the originals and the  
196 rewritten versions.

197 **United Kingdom Medical Licensing Assessment, Applied Knowledge Test.** ChatGPT-4o answered  
198 94 of 100 questions on the original paper. Five of the questions included pictures. ChatGPT answered  
199 four of these correctly. ChatGPT then scored 93%, 91% and 95% on the three collections of rewrites.  
200 One question, on herpes zoster ophthalmicus, was answered incorrectly on all four occasions. In all  
201 other cases there was no consistent pattern. Some questions that ChatGPT had answered incorrectly  
202 were answered correctly once rewritten, but the converse was also true. 85% of questions were  
203 answered correctly in all four versions (original and all three rewrites). The full dataset and questions  
204 are in Supplementary data S1.

205 **United States Medical Licensing Exam Step 1.** ChatGPT-4o scored 89.9% (107/119) of the original  
206 questions correctly. Of the original 119, there were images in 23 of them, of which 16 (69.6%) were  
207 answered correctly. This perhaps suggested that ChatGPT might struggle more with the picture  
208 questions in this exam. Given that ChatGPT-4o had already demonstrated no impairment of  
209 performance when rewriting questions from the UK MLA AKT into a novel format, we decided to  
210 rewrite only a sample of 27 of the USMLE questions, but to probe further the apparent diminished  
211 performance on questions contained pictures by including 13 picture questions, of which 5 had been  
212 answered incorrectly in the original paper. Of the sample of 27, ChatGPT scored 74.1% (20/27) on  
213 the original versions, and then 85.2% (23/27), 70.4% (19/27) and 85.2% (23/27) on the rewrites. Only  
214 one question was answered incorrectly in all four versions. This was a picture question which  
215 involved answering questions based on a graph, while the other four picture questions which  
216 ChatGPT had answered incorrectly were then answered correctly at least once during the rewrites.  
217 55.6% (15/27) of questions were answered correctly on all four occasions. The full dataset and  
218 questions are in Supplementary data S1.

219 **Novel questions:** A total of 90 novel questions were generated, of which ChatGPT answered 75  
220 correctly. 28 of these questions were in pairs (2x14) which assessed the same learning. One version  
221 of the question contained a labelled image where the labels were simple letters (A,B,C etc) and these  
222 were the answer options, whereas the paired question contained answer options in text form. An  
223 example of this format is in Figure 1. ChatGPT answered 13/14 of the text version of these questions,  
224 but only 2/14 of the labelled image questions. A summary of the analysis is in Supplementary data  
225 S1. The novel questions may be shared upon request but are not published here due to the images  
226 contained within.

227 **Identification of labels on images.** Ten of the labelled images were structured in a way that it was  
228 reasonable to upload them to ChatGPT-4o with the prompt 'Can you identify all the labels (A-X) on  
229 the uploaded image?' where 'X' was either E, F, G or H depending on the number of labels. Of a  
230 total of 66 labels across the 10 images, ChatGPT correctly labelled 25 items. For all 10 images  
231 ChatGPT correctly identified the main structure in the image (e.g. brain, kidney) but not the labelled  
232 subregions.

233



234

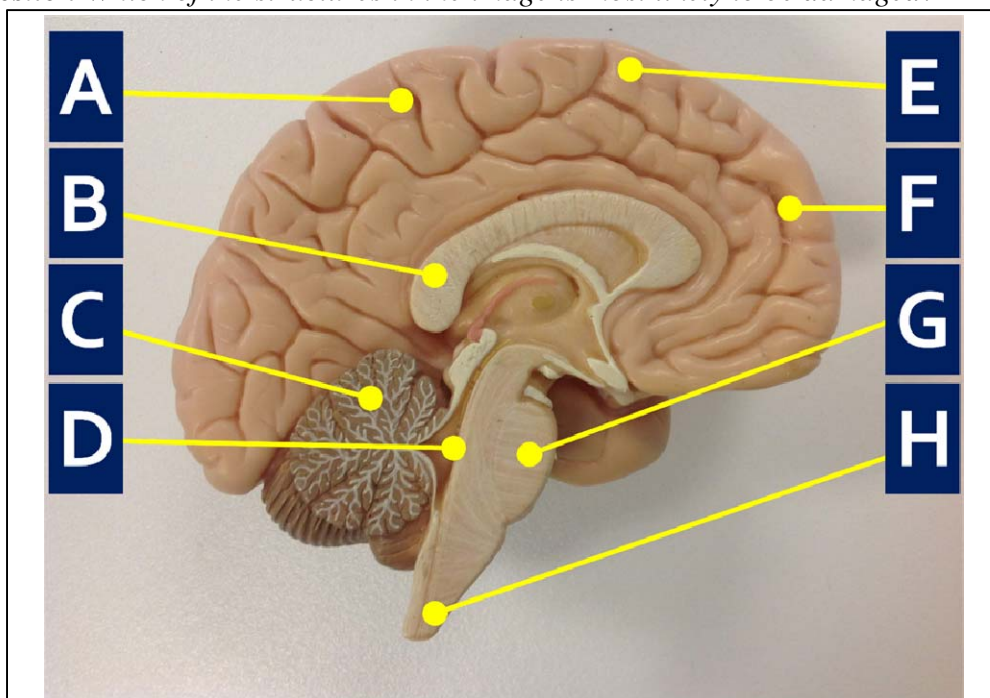
235 **Common scenario** An elderly gentleman is rushed to hospital after being found on the floor at home.  
236 He appears to be able to breathe and his heartrate is elevated but stable. However he appears to be  
237 completely paralysed and does not respond when asked questions. His pupils are pinpoint. He does  
238 not blink when something goes near his eyes, but when a light is shone into his eyes, they move  
239 horizontally to follow the light.

240 **Text question** Damage to which structure in the brain is most likely to result in the above  
241 presentation?

- 242 A. Primary Motor Cortex
- 243 B. Hippocampus
- 244 C. Cerebellum
- 245 D. Nucleus Accumbens
- 246 E. Globus Pallidus
- 247 F. Substantia Nigra
- 248 G. Pons
- 249 H. Medulla

250  
251 **Image Question** Which of the structures in the image is most likely to be damaged?

- 252 A
- 253 B
- 254 C
- 255 D
- 256 E
- 257 F
- 258 G
- 259 H



266

267 **Figure 1.** An example of a novel-higher order MCQ written using established guidelines (5), with  
268 text options as answers (which ChatGPT answers correctly), or a labelled image which ChatGPT  
269 answers incorrectly. Note that the answer options do not correspond exactly.

## 270 Discussion

271 ChatGPT-4o showed a very high level of performance on the papers tested, even when the questions  
272 were rewritten so that they assessed the same learning but with different wording. This level of  
273 performance was also found on completely novel questions written in the style of professional  
274 licensing exams. Our analysis included many questions which are based on images, and almost all  
275 questions were designed to assess higher-order problem-solving (4,5).

276 A repeated finding from the research on academic misconduct demonstrates that one of the strongest  
277 factors contributing to an increased likelihood in the occurrence of academic dishonesty is the ease  
278 with which it can be committed (12,27). Cheating in online exams was already high before the  
279 emergence of ChatGPT (12) and our findings demonstrate that any student using ChatGPT would  
280 likely receive an excellent mark even if they had no prior knowledge whatsoever, thus further  
281 increasing any temptation to cheat. Thus it seems reasonable to propose that our findings mean online  
282 unproctored summative exams are now no longer a valid form of assessment, a conclusion which is  
283 in contrast to findings published following an analysis of exam performance during the COVID  
284 pandemic, but before the emergence of ChatGPT (15).

285 The high performance levels of ChatGPT may also increase the temptation to cheat using ChatGPT  
286 even in proctored exams, particularly if they are taken online; data suggest that proctoring  
287 considerably reduces cheating in online exams but does not eliminate it completely (12). We are not  
288 aware of any current data on the extent to which students are using ChatGPT to cheat in online  
289 exams, proctored or unproctored, although this is the subject of ongoing work. A study conducted in  
290 Vietnam in May 2023 showed that 23.7% of undergraduates cheated using ChatGPT, although the  
291 assessment formats were not specified (28). A study conducted at around the same time in US high  
292 schools found similar numbers in one school, though lower in two others (29). These figures seem  
293 likely to increase as ChatGPT becomes better known and more widely available, along with similar  
294 tools such as Claude.AI.

295 One intuitive response to these challenges is to design questions which ChatGPT finds harder to  
296 answer. This ‘arms race’ approach is partly the genesis of this paper, based in part on suggestions  
297 that earlier studies observed that ChatGPT could not process image-based questions at all, and other  
298 studies suggesting that ChatGPT is a ‘copy and paste’ machine whose impact can be minimized by  
299 using novel questions for each sitting of an exam (9). We did find that ChatGPT struggled more on a  
300 very specific type of MCQ, where the answer items were single letter labels and arrows on images.  
301 There is more than one possible explanation for this. These questions are designed to require  
302 ‘assumed knowledge’ and so to be harder to answer than factual recall questions (5). For example,  
303 the picture item shown in figure 1 requires the test taker to know that the scenario represents the  
304 clinical condition Locked-In Syndrome, and then to know that this is associated with damage to the  
305 pons, and then to be able to identify the anatomical location of the pons on a picture of a model.  
306 ChatGPT consistently struggled with these and so one interpretation is that it is this ‘multi-step’  
307 approach that trips up ChatGPT. However, ChatGPT was consistently correct on the text versions of  
308 these questions and would give detailed descriptions of the answer option and was clearly able to  
309 identify, in text form, where the pons is located (for example). But when simply asked to identify the  
310 labels on these images ChatGPT struggled, indicating that it is the processing of these specific types  
311 of text-labelled images which ChatGPT struggles with.

312 One intuitive conclusion from these findings with images is that such questions could be used to  
313 thwart ChatGPT and so deter cheating in online exams. However, we caution against over-

314 interpreting this finding as identifying a ‘ChatGPT-proof’ question formats. Writing an entire exam  
315 based on these questions seems implausible and unlikely to be valid. This limitation likely applies to  
316 other methods identified as a way of ‘defeating’ ChatGPT. An older study, using an unidentified  
317 version of ChatGPT, showed that ChatGPT overselects answer options ‘all of the above’ or ‘none of  
318 the above’, meaning that when these answer options are present but are incorrect, ChatGPT shows a  
319 much lower performance compared to when these answer options are absent or when they are present  
320 but are the correct answer. However, designing questions which incorporate this flaw also seems  
321 likely to be a short-term measure that may well result in poorer quality questions and weaker  
322 curriculum coverage. These types of answer options are also advised against when writing high  
323 quality assessment items (5).

324 Any reduction in the use of online unproctored exams will clearly not eradicate academic  
325 misconduct. There are a wide range of dishonest behaviours undertaken by medical and other  
326 students (30), and the performance of ChatGPT on assessment formats such as essays is also very  
327 strong (31). Essays are, by design, asynchronous and unmonitored, meaning that it would be almost  
328 impossible to prevent a student from using ChatGPT to complete assignments in these formats.  
329 Detection tools have been developed and these appear to show good accuracy for raw text generated  
330 by tools such as ChatGPT (32) but they can be easily circumvented (33) and even a very small rate of  
331 false-positives is problematic since there is no independent source to match a student assignment to,  
332 unlike with ‘conventional’ plagiarism, meaning that problematic, adversarial situations can quickly  
333 arise when students are accused of cheating on essays using ChatGPT (34).

334 The performance of ChatGPT-4o demonstrated here shows a modest improvement of that seen using  
335 GPT-4, which itself shows a much improved performance compared to GPT-3 and GPT-3.5 (1),  
336 although many prior papers excluded image-based questions from their analyses whereas they are  
337 included here. This trend of improving performance seems likely to continue; at the time of writing  
338 (June 2024), OpenAI are rolling out enhanced visual recognition features in GPT-4o to their  
339 subscribers, meaning that users are able to simply point their camera at the question and it will scan  
340 and ‘read’ the text before generating an answer (21).

341 The high performance of ChatGPT-4o on the exams tested here and elsewhere leads naturally to a  
342 question of whether these tools might also be able to *write* such exams. A review on some of the  
343 older versions of these tools concluded that question generation is possible although with some  
344 limitations and proposed further testing (35). It is now possible to upload considerable volumes of  
345 data to ChatGPT and to build custom GPTs which specific instructions tailored to certain tasks, as  
346 designed by the creator; this approach has already shown promise for the creation of USMLE-style  
347 assessment items and may even be able to generate an entire exam and blueprint it to a curriculum,  
348 saving considerable time and cost for educators and universities (36). This possibility arose during  
349 the conduct of the study here wherein some questions that were initially answered incorrectly by  
350 ChatGPT revealed either strong distractors or potential ambiguities in the question stem or associated  
351 image, suggesting weaknesses in the question itself. No questions tested here were eliminated from  
352 analysis for being actually incorrect or of poor quality, but this analysis suggested that such issues  
353 might be easily identified by using ChatGPT as an adjunct to exam creation and standard setting.

354 Similar benefits could also be obtained for students. The research team here noted the accuracy and  
355 value of the explanations provided by ChatGPT when answering the questions, and these naturally  
356 suggest the potential of ChatGPT, and the aforementioned custom GPTs, as study tools for students.  
357 Such an approach has been successfully used in ophthalmology (37) and anatomy learning (38).



359 **Conclusion**

360 ChatGPT-4o shows very high levels of performance on MCQ-based applied knowledge tests,  
361 including questions with images. These data echo but improve further upon findings from earlier  
362 versions of ChatGPT (39) and suggest that educators will find it extremely difficult to write questions  
363 which are ‘ChatGPT-proof’, even if they are completely novel and image-based. The logical  
364 conclusion is that unproctored online exams are no longer a valid form of assessment, even when  
365 assessing higher order learning. These assessments, and lower-level MCQs based exams testing core  
366 foundational knowledge, should only be conducted under secure conditions.

367 **Conflict of Interest Statement**

368 On behalf of all authors, the corresponding author states that there is no conflict of interest

369

## 370 **References**

- 371 1. Newton P, Xiromeriti M. ChatGPT performance on multiple choice question examinations in  
372 higher education. A pragmatic scoping review. *Assess Eval High Educ.* 2024;0(0):1–18.
- 373 2. Garabet R, Mackey BP, Cross J, Weingarten M. ChatGPT-4 Performance on USMLE Step 1  
374 Style Questions and Its Implications for Medical Education: A Comparative Study Across  
375 Systems and Disciplines. *Med Sci Educ.* 2024 Feb 1;34(1):145–52.
- 376 3. Lai UH, Wu KS, Hsu TY, Kan JKC. Evaluating the performance of ChatGPT-4 on the United  
377 Kingdom Medical Licensing Assessment. *Front Med.* 2023 Sep 19;10:1240915.
- 378 4. Billings M, DeRuchie K, Hussie K, Kulesher A, Merrell J, Morales A, et al. Constructing written  
379 test questions for the Health Sciences [Internet]. National Board of Medical Examiners; 2020  
380 [cited 2022 Apr 7]. Available from: [https://www.nbme.org/sites/default/files/2020-](https://www.nbme.org/sites/default/files/2020-11/NBME_Item%20Writing%20Guide_2020.pdf)  
381 [11/NBME\\_Item%20Writing%20Guide\\_2020.pdf](https://www.nbme.org/sites/default/files/2020-11/NBME_Item%20Writing%20Guide_2020.pdf)
- 382 5. Newton PM. Guidelines for Creating Online MCQ-Based Exams to Evaluate Higher Order  
383 Learning and Reduce Academic Misconduct. In: Eaton SE, editor. *Handbook of Academic*  
384 *Integrity* [Internet]. Singapore: Springer Nature; 2023 [cited 2023 Jul 13]. p. 1–17. Available  
385 from: [https://doi.org/10.1007/978-981-287-079-7\\_93-1](https://doi.org/10.1007/978-981-287-079-7_93-1)
- 386 6. Arkoudas K. GPT-4 Can't Reason [Internet]. arXiv; 2023 [cited 2024 Feb 18]. Available from:  
387 <http://arxiv.org/abs/2308.03762>
- 388 7. Yeo YH, Samaan JS, Ng WH, Ting PS, Trivedi H, Vipani A, et al. Assessing the performance of  
389 ChatGPT in answering questions regarding cirrhosis and hepatocellular carcinoma. *Clin Mol*  
390 *Hepatol.* 2023 Jul;29(3):721–32.
- 391 8. Marcus G. Partial Regurgitation and how LLMs really... [Internet]. Marcus on AI. 2024 [cited  
392 2024 Jun 3]. Available from: [https://garymarcus.substack.com/p/partial-regurgitation-and-how-](https://garymarcus.substack.com/p/partial-regurgitation-and-how-llms/comments)  
393 [llms/comments](https://garymarcus.substack.com/p/partial-regurgitation-and-how-llms/comments)
- 394 9. Lo CK. What Is the Impact of ChatGPT on Education? A Rapid Review of the Literature. *Educ*  
395 *Sci.* 2023 Apr;13(4):410.
- 396 10. Ram S, Qian C. A Study on the Vulnerability of Test Questions against ChatGPT-based  
397 Cheating. In: 2023 International Conference on Machine Learning and Applications (ICMLA)  
398 [Internet]. 2023 [cited 2024 Jun 17]. p. 1710–5. Available from:  
399 <https://ieeexplore.ieee.org/abstract/document/10460039>
- 400 11. Abbas A, Rehman MS, Rehman SS. Comparing the Performance of Popular Large Language  
401 Models on the National Board of Medical Examiners Sample Questions. *Cureus.* 16(3):e55991.
- 402 12. Newton PM, Essex K. How Common is Cheating in Online Exams and did it Increase During the  
403 COVID-19 Pandemic? A Systematic Review. *J Acad Ethics* [Internet]. 2023 Aug 4 [cited 2023  
404 Aug 7]; Available from: <https://doi.org/10.1007/s10805-023-09485-5>
- 405 13. Marano E, Newton PM, Birch Z, Croombs M, Gilbert C, Draper MJ. What is the student  
406 experience of remote proctoring? A pragmatic scoping review. *High Educ Q.* n/a(n/a):e12506.

- 407 14. Meulmeester FL, Dubois EA, Krommenhoek-van Es C (Tineke), de Jong PGM, Langers AMJ.  
408 Medical Students' Perspectives on Online Proctoring During Remote Digital Progress Test. *Med*  
409 *Sci Educ*. 2021 Sep 30;31(6):1773–7.
- 410 15. Chan JCK, Ahn D. Unproctored online exams provide meaningful assessment of student  
411 learning. *Proc Natl Acad Sci*. 2023 Aug;120(31):e2302020120.
- 412 16. Newton PM. The validity of unproctored online exams is undermined by cheating. *Proc Natl*  
413 *Acad Sci*. 2023 Oct 10;120(41):e2312978120.
- 414 17. Newton PM, Da Silva A, Berry S. The Case for Pragmatic Evidence-Based Higher Education: A  
415 Useful Way Forward? *Front Educ* [Internet]. 2020 [cited 2021 May 8];5. Available from:  
416 <https://www.frontiersin.org/articles/10.3389/educ.2020.583157/full>
- 417 18. Newton PM, Da Silva A, Peters LG. A Pragmatic Master List of Action Verbs for Bloom's  
418 Taxonomy. *Front Educ* [Internet]. 2020 [cited 2020 Jul 14];5. Available from:  
419 <https://www.frontiersin.org/articles/10.3389/educ.2020.00107/full>
- 420 19. Koga S. The Potential of ChatGPT in Medical Education: Focusing on USMLE Preparation. *Ann*  
421 *Biomed Eng*. 2023 Oct 1;51(10):2123–4.
- 422 20. Daungsupawong H, Wiwanitkit V. ChatGPT-4 Performance on USMLE Step 1 Style Questions  
423 and Its Implications for Medical Education: Correspondence. *Med Sci Educ* [Internet]. 2024 Apr  
424 5 [cited 2024 Jun 3]; Available from: <https://doi.org/10.1007/s40670-024-02033-9>
- 425 21. OpenAI. Hello GPT-4o [Internet]. [cited 2024 Jun 3]. Available from:  
426 <https://openai.com/index/hello-gpt-4o/>
- 427 22. Wikiversity. Fundamentals of Neuroscience/Exams - Wikiversity [Internet]. 2013 [cited 2024  
428 Feb 10]. Available from: [https://en.wikiversity.org/wiki/Fundamentals\\_of\\_Neuroscience/Exams](https://en.wikiversity.org/wiki/Fundamentals_of_Neuroscience/Exams)
- 429 23. Medical Schools Council. Practice exam for the MS AKT | Medical Schools Council [Internet].  
430 2023 [cited 2024 Mar 10]. Available from: [https://www.medschools.ac.uk/medical-licensing-](https://www.medschools.ac.uk/medical-licensing-assessment/preparing-for-the-ms-akt/practice-exam-for-the-ms-akt)  
431 [assessment/preparing-for-the-ms-akt/practice-exam-for-the-ms-akt](https://www.medschools.ac.uk/medical-licensing-assessment/preparing-for-the-ms-akt/practice-exam-for-the-ms-akt)
- 432 24. United States Medical Licensing Examination. Step 1 Sample Test Questions | USMLE  
433 [Internet]. 2021 [cited 2024 Jun 10]. Available from: [https://www.usmle.org/prepare-your-](https://www.usmle.org/prepare-your-exam/step-1-materials/step-1-sample-test-questions)  
434 [exam/step-1-materials/step-1-sample-test-questions](https://www.usmle.org/prepare-your-exam/step-1-materials/step-1-sample-test-questions)
- 435 25. GMC. Making and using visual and audio recordings of patients (summary) [Internet]. General  
436 Medical Council; 2011 [cited 2023 Jun 15]. Available from: [https://www.gmc-](https://www.gmc-uk.org/professional-standards/professional-standards-for-doctors/making-and-using-visual-and-audio-recordings-of-patients)  
437 [uk.org/professional-standards/professional-standards-for-doctors/making-and-using-visual-and-](https://www.gmc-uk.org/professional-standards/professional-standards-for-doctors/making-and-using-visual-and-audio-recordings-of-patients)  
438 [audio-recordings-of-patients](https://www.gmc-uk.org/professional-standards/professional-standards-for-doctors/making-and-using-visual-and-audio-recordings-of-patients)
- 439 26. GMC. MLA content map [Internet]. 2021 [cited 2024 Jun 15]. Available from: [https://www.gmc-](https://www.gmc-uk.org/education/medical-licensing-assessment/mla-content-map)  
440 [uk.org/education/medical-licensing-assessment/mla-content-map](https://www.gmc-uk.org/education/medical-licensing-assessment/mla-content-map)
- 441 27. Bretag T, Harper R, Burton M, Ellis C, Newton P, Rozenberg P, et al. Contract cheating: a  
442 survey of Australian university students. *Stud High Educ*. 2019 Nov 2;44(11):1837–56.

- 443 28. Nguyen HM, Goto D. Unmasking academic cheating behavior in the artificial intelligence era:  
444 Evidence from Vietnamese undergraduates. *Educ Inf Technol* [Internet]. 2024 Feb 5 [cited 2024  
445 Feb 18]; Available from: <https://doi.org/10.1007/s10639-024-12495-4>
- 446 29. Lee VR, Pope D, Miles S, Zárate RC. Cheating in the age of generative AI: A high school survey  
447 study of cheating behaviors before and after the release of ChatGPT. *Comput Educ Artif Intell*.  
448 2024 Dec 1;7:100253.
- 449 30. Henning MA, Chen Y, Ram S, Malpas P. Describing the Attributional Nature of Academic  
450 Dishonesty. *Med Sci Educ*. 2019 Jun 1;29(2):577–81.
- 451 31. Herbold S, Hautli-Janisz A, Heuer U, Kikteva Z, Trautsch A. AI, write an essay for me: A large-  
452 scale comparison of human-written versus ChatGPT-generated essays [Internet]. arXiv; 2023  
453 [cited 2023 May 8]. Available from: <http://arxiv.org/abs/2304.14276>
- 454 32. Weber-Wulff D, Anohina-Naumeca A, Bjelobaba S, Foltýnek T, Guerrero-Dib J, Popoola O, et  
455 al. Testing of Detection Tools for AI-Generated Text [Internet]. arXiv; 2023 [cited 2023 Aug 7].  
456 Available from: <http://arxiv.org/abs/2306.15666>
- 457 33. Perkins M, Roe J, Vu BH, Postma D, Hickerson D, McGaughran J, et al. arXiv.org. 2024 [cited  
458 2024 Jun 11]. GenAI Detection Tools, Adversarial Techniques and Implications for Inclusivity in  
459 Higher Education. Available from: <https://arxiv.org/abs/2403.19148v1>
- 460 34. Gorichanaz T. Accused: How students respond to allegations of using ChatGPT on assessments.  
461 *Learn Res Pract* [Internet]. 2023 Jul 3 [cited 2024 May 3]; Available from:  
462 <https://www.tandfonline.com/doi/abs/10.1080/23735082.2023.2254787>
- 463 35. Artsi Y, Sorin V, Konen E, Glicksberg BS, Nadkarni G, Klang E. Large language models for  
464 generating medical examinations: systematic review. *BMC Med Educ*. 2024 Mar 29;24(1):354.
- 465 36. Kiyak YS, Kononowicz AA. Case-based MCQ generator: A custom ChatGPT based on  
466 published prompts in the literature for automatic item generation. *Med Teach* [Internet]. 2024  
467 Feb 6 [cited 2024 Jun 11]; Available from:  
468 <https://www.tandfonline.com/doi/abs/10.1080/0142159X.2024.2314723>
- 469 37. Sevgi M, Antaki F, Keane PA. Medical education with large language models in ophthalmology:  
470 custom instructions and enhanced retrieval capabilities. *Br J Ophthalmol* [Internet]. 2024 May 7  
471 [cited 2024 Jun 11]; Available from: <https://bjo.bmj.com/content/early/2024/05/07/bjo-2023-325046>
- 473 38. Collins BR, Black EW, Rarey KE. Introducing AnatomyGPT: A customized artificial  
474 intelligence application for anatomical sciences education. *Clin Anat* [Internet]. [cited 2024 Jun  
475 11];n/a(n/a). Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/ca.24178>
- 476 39. Sood A, Mansoor N, Memmi C, Lynch M, Lynch J. Generative pretrained transformer-4, an  
477 artificial intelligence text predictive model, has a high capability for passing novel written  
478 radiology exam questions. *Int J Comput Assist Radiol Surg*. 2024 Apr 1;19(4):645–53.
- 479