

# 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 questions with pictures, 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 of  
33 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, requiring the application of core knowledge to  
53 clinical scenarios (4) and represent a broader principle wherein multiple choice questions can, if  
54 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 the 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 identify features of exam questions which ChatGPT might struggle with, for example  
69 an increase in the number of answer items, increasing language complexity or having multiple correct  
70 answers. However none of these appears to have any effect on the numbers of questions which  
71 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 increase the risk of cheating. During the COVID-19 pandemic, the percentage of students  
84 who admitted to cheating in online exams appeared to double, and more students reported cheating  
85 than not (12). One apparent solution to this problem is to increase the use of online  
86 proctoring/invigilation systems to monitor student behaviour. However, these then drive back up the  
87 cost of the online exams, and the student experience of remote proctoring is poor, with concerns  
88 about privacy, fairness, inclusivity and cost (13,14). An alternative is to avoid the use of proctoring  
89 altogether. A high profile 2023 publication analysed exam performance data from the COVID

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

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

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

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

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

133

134

135

## 136 **Methods**

137 The following question sources were tested.

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

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

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

160 ***Analysis of existing medical licensing exams and rewrites.*** Each question was tested using a single  
161 shot method in a way that would be expected to be the most likely approach taken by a student who  
162 was seeking to cheat on an MCQ exam, i.e. the text was highlighted in the pdf (original questions) or  
163 word document (rewrites), copied and then pasted directly into ChatGPT-4o with no attempt to  
164 format the text. Where the question included a picture, this was copied using screen clipping, saved  
165 and uploaded as a .png file with only the country and the question number as the file name (e.g.  
166 ‘UK32’). No additional prompts were given apart from the content of the question. Each question  
167 was asked in a new chat and no memory functions were activated. For the USMLE questions, a  
168 ‘temporary chat’ was activated for each question. No responses were given to ChatGPT. ChatGPT’s  
169 first response was recorded each time as correct/incorrect. ChatGPT-4o tests were undertaken May  
170 14-24 2024.

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

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

192



## 193 Results

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

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

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

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

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

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





237

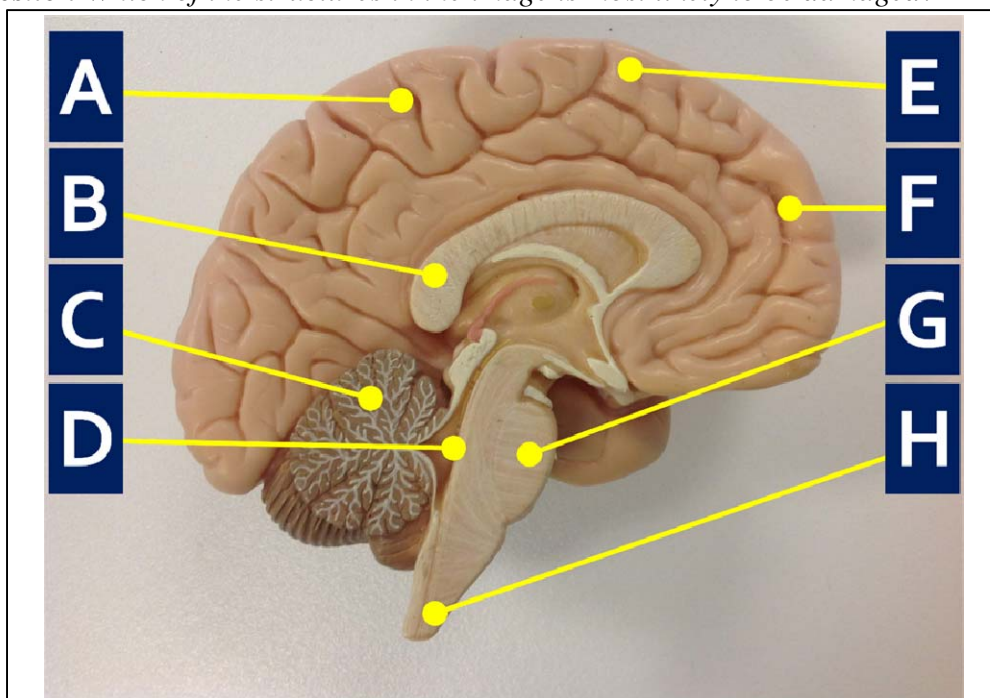
238 **Common scenario** An elderly gentleman is rushed to hospital after being found on the floor at home.  
239 He appears to be able to breathe and his heartrate is elevated but stable. However he appears to be  
240 completely paralysed and does not respond when asked questions. His pupils are pinpoint.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

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

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

253  
254 **Image Question** Which of the structures in the image is most likely to be damaged?

- 255 A
- 256 B
- 257 C
- 258 D
- 259 E
- 260 F
- 261 G
- 262 H



269

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

## 273 Discussion

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

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

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

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

317 One intuitive conclusion from these findings with images is that such questions could be used to  
318 thwart ChatGPT and so deter cheating in online exams. However, we caution against this  
319 interpretation. Writing an entire exam based on these types of questions seems implausible and  
320 unlikely to be valid. This limitation likely applies to other methods identified as a way of ‘defeating’  
321 ChatGPT. For example, an older study, using an unidentified version of ChatGPT, showed that  
322 ChatGPT overselects answer options ‘all of the above’ or ‘none of the above’, meaning that when  
323 these answer options are present but are incorrect, ChatGPT shows a much lower performance  
324 compared to when these answer options are absent or when they are present but are the correct  
325 answer. However, designing questions which incorporate this flaw also seems likely to be a short-  
326 term measure that may well result in poorer quality questions and weaker curriculum coverage. These  
327 types of answer options are also advised against when writing high quality assessment items (5).

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

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

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

359 Similar benefits could also be obtained for students. The research team here noted the accuracy and  
360 value of the explanations provided by ChatGPT when answering the questions, and these naturally

361 suggest the potential of ChatGPT, and the aforementioned custom GPTs, as study tools for students.  
362 Such an approach has been successfully used in ophthalmology (37) and anatomy learning (38).

### 363 **Conclusion**

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

### 371 **Conflict of Interest Statement**

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

373

## 374 **References**

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



- 411 14. Meulmeester FL, Dubois EA, Krommenhoek-van Es C (Tineke), de Jong PGM, Langers AMJ.  
412 Medical Students' Perspectives on Online Proctoring During Remote Digital Progress Test. *Med*  
413 *Sci Educ*. 2021 Sep 30;31(6):1773–7.
- 414 15. Chan JCK, Ahn D. Unproctored online exams provide meaningful assessment of student  
415 learning. *Proc Natl Acad Sci*. 2023 Aug;120(31):e2302020120.
- 416 16. Newton PM. The validity of unproctored online exams is undermined by cheating. *Proc Natl*  
417 *Acad Sci*. 2023 Oct 10;120(41):e2312978120.
- 418 17. Newton PM, Da Silva A, Berry S. The Case for Pragmatic Evidence-Based Higher Education: A  
419 Useful Way Forward? *Front Educ* [Internet]. 2020 [cited 2021 May 8];5. Available from:  
420 <https://www.frontiersin.org/articles/10.3389/feduc.2020.583157/full>
- 421 18. Newton PM, Da Silva A, Peters LG. A Pragmatic Master List of Action Verbs for Bloom's  
422 Taxonomy. *Front Educ* [Internet]. 2020 [cited 2020 Jul 14];5. Available from:  
423 <https://www.frontiersin.org/articles/10.3389/feduc.2020.00107/full>
- 424 19. Koga S. The Potential of ChatGPT in Medical Education: Focusing on USMLE Preparation. *Ann*  
425 *Biomed Eng*. 2023 Oct 1;51(10):2123–4.
- 426 20. Daungsupawong H, Wiwanitkit V. ChatGPT-4 Performance on USMLE Step 1 Style Questions  
427 and Its Implications for Medical Education: Correspondence. *Med Sci Educ* [Internet]. 2024 Apr  
428 5 [cited 2024 Jun 3]; Available from: <https://doi.org/10.1007/s40670-024-02033-9>
- 429 21. OpenAI. Hello GPT-4o [Internet]. [cited 2024 Jun 3]. Available from:  
430 <https://openai.com/index/hello-gpt-4o/>
- 431 22. Wikiversity. Fundamentals of Neuroscience/Exams - Wikiversity [Internet]. 2013 [cited 2024  
432 Feb 10]. Available from: [https://en.wikiversity.org/wiki/Fundamentals\\_of\\_Neuroscience/Exams](https://en.wikiversity.org/wiki/Fundamentals_of_Neuroscience/Exams)
- 433 23. Medical Schools Council. Practice exam for the MS AKT | Medical Schools Council [Internet].  
434 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)  
435 [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)
- 436 24. United States Medical Licensing Examination. Step 1 Sample Test Questions | USMLE  
437 [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)  
438 [exam/step-1-materials/step-1-sample-test-questions](https://www.usmle.org/prepare-your-exam/step-1-materials/step-1-sample-test-questions)
- 439 25. GMC. Making and using visual and audio recordings of patients (summary) [Internet]. General  
440 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)  
441 [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)  
442 [audio-recordings-of-patients](https://www.gmc-uk.org/professional-standards/professional-standards-for-doctors/making-and-using-visual-and-audio-recordings-of-patients)
- 443 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)  
444 [uk.org/education/medical-licensing-assessment/mla-content-map](https://www.gmc-uk.org/education/medical-licensing-assessment/mla-content-map)
- 445 27. Bretag T, Harper R, Burton M, Ellis C, Newton P, Rozenberg P, et al. Contract cheating: a  
446 survey of Australian university students. *Stud High Educ*. 2019 Nov 2;44(11):1837–56.



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