

1 **Performance of ChatGPT on Chinese National Medical Licensing Examinations:**
2 **A Five-Year Examination Evaluation Study for Physicians, Pharmacists and**
3 **Nurses**

4 Hui Zong^{1, †}, Jiakun Li^{1, †}, Erman Wu¹, Rongrong Wu¹, Junyu Lu¹, Bairong Shen^{1, *}

5 1. Department of Urology and Institutes for Systems Genetics, Frontiers Science Center
6 for Disease-related Molecular Network, West China Hospital, Sichuan University,
7 Chengdu, 610212, China.

8

9 † The authors contributed equally.

10 * Corresponding Author: Bairong Shen, PhD, Professor, Institutes for Systems Genetics,
11 Frontiers Science Center for Disease-related Molecular Network, West China Hospital,
12 Sichuan University, No. 37, Guoxue Alley, Chengdu, Sichuan, China. Tel: 86-
13 15995854635, Email: bairong.shen@scu.edu.cn

14

15 **Abstract**

16 **Background:** Large language models like ChatGPT have revolutionized the field of
17 natural language processing with their capability to comprehend and generate textual
18 content, showing great potential to play a role in medical education.

19 **Objective:** This study aimed to quantitatively evaluate and comprehensively analysis
20 the performance of ChatGPT on three types of national medical examinations in China,
21 including National Medical Licensing Examination (NMLE), National Pharmacist
22 Licensing Examination (NPLE), and National Nurse Licensing Examination (NNLE).

23 **Methods:** We collected questions from Chinese NLMLE, NPLE and NNLE from year
24 2017 to 2021. In NMLE and NPLE, each exam consists of 4 units, while in NNLE, each
25 exam consists of 2 units. The questions with figures, tables or chemical structure were
26 manually identified and excluded by clinician. We applied direct instruction strategy
27 via multiple prompts to force ChatGPT to generate the clear answer with the capability
28 to distinguish between single-choice and multiple-choice questions.

29 **Results:** ChatGPT failed to pass the threshold score (0.6) in any of the three types of

30 **NOTE:** This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.
examinations over the five years. Specifically, in the NMLE, the highest recorded score

31 was 0.5467, which was attained in both 2018 and 2021. In the NPLE, the highest score
32 was 0.5599 in 2017. In the NNLE, the most impressive result was shown in 2017, with
33 a score of 0.5897, which is also the highest score in our entire evaluation. ChatGPT's
34 performance showed no significant difference in different units, but significant
35 difference in different question types. ChatGPT performed well in a range of subject
36 areas, including clinical epidemiology, human parasitology, and dermatology, as well
37 as in various medical topics such as molecules, health management and prevention,
38 diagnosis and screening.

39 **Conclusions:** These results indicate ChatGPT failed the NMLE, NPLE and NNLE in
40 China, spanning from year 2017 to 2021. but show great potential of large language
41 models in medical education. In the future high-quality medical data will be required
42 to improve the performance.

43

44 **Keywords:** Medical Education, Medical Examination, Artificial Intelligence, Natural
45 Language Processing, ChatGPT

46

47 **Introduction**

48 In the last decade, artificial intelligence (AI) technology has undergone a rapid
49 evolution, achieving noteworthy breakthroughs in numerous fields [1, 2]. Recently, one
50 such breakthrough that has garnered considerable attention is ChatGPT [3], an AI
51 chatbot powered by generative pre-trained transformer (GPT) architecture, specifically
52 GPT-3.5 with 175 billion parameters. This innovative technology is developed through
53 human feedback reinforcement learning and trained on extensive textual data.
54 Remarkably, ChatGPT exhibits remarkable capabilities in various tasks, including but
55 not limited to intelligent dialogue [4], knowledge question answering[5], and text
56 generation[6], thus showcasing unprecedented potential for further development.

57 In medical domain, there has been growing interest in exploration of large language
58 models for tasks such as biomedical question answering (BioGPT [7]), and automatic
59 dialogue generation (DialoGPT [8, 9]). Regrettably, these studies have so far
60 demonstrated limited practical utility in clinical practice. However, ChatGPT, with its

61 powerful language understanding and generation capabilities, showing significant
62 potential in the fields of clinical response generation [5, 6], clinical decision support [4,
63 10, 11], medical education [12, 13], literature information retrieve [14], scientific
64 writing [15-18], and beyond. Recent studies have demonstrated that ChatGPT can pass
65 the United States Medical Licensing Exam (USMLE) [19, 20] , Radiology Board-style
66 Examination [21], UK Neurology Specialty Certificate Examination [22], and Plastic
67 Surgery In-Service Exam [23], with results that are comparable to those of human
68 experts. Nevertheless, Other studies have also indicated that ChatGPT failed to pass the
69 Family Medicine Board Exam [24], and Pharmacist Qualification Examination [25].
70 Possible explanations for this performance difference include language and cultural
71 differences, variations in examination content [26]. These studies highlighted the
72 ChatGPT's ability to comprehend the complex language used in medical contexts and
73 its potential for use in medical education. However, current researches are limited in
74 two aspects. Firstly, it largely focuses on the English language, and secondly, it
75 predominantly emphasizes the physician's examination. Additional investigation is
76 necessary to explore the potential of ChatGPT in other non-English languages and
77 various medical examinations, which can deliver substantial benefits for its expanded
78 application in the medical domain.

79 China, with a population of over 1.4 billion, faces a significant medical burden.
80 The provision of healthcare services involves a collaborative effort among physicians,
81 pharmacists, and nurses who work diligently to offer the best possible care to patients.
82 Physicians are responsible for diagnosing and treating illnesses, pharmacists ensure the
83 appropriate medication is dispensed and administered correctly, while nurses attend to
84 patients' daily medical management and care service. Due to limited medical resources,
85 medical professionals in China face immense pressure, but remain committed to
86 providing high-quality services. The advent of ChatGPT offers a promising solution to
87 ease this burden by delivering intelligent, efficient, and precise medical services to
88 physicians, pharmacists, and nurses.

89 Medical examinations, including the Chinese National Medical Licensing
90 Examination (NMLE), the Chinese National Pharmacist Licensing Examination

91 (NPLE), and the Chinese National Nurse Licensing Examination (NNLE) are
92 implemented by the government to improve professional standards, ensure medical
93 safety and enhance healthcare services quality [27]. Through these medical
94 examinations, the medical knowledge, clinical skills, and ethical standards mastered by
95 medical staffs can significantly improve the quality of their services. This, in turn, can
96 reduce the incidence of medical errors and accidents, and protect the fundamental right
97 to health and safety of patients.

98 These medical licensing examinations aim to comprehensively evaluate candidate's
99 knowledge of medical science, clinical examination, disease diagnosis, surgical
100 treatment, patient prognosis, policies, and regulations, among other areas. Successfully
101 passing these examinations is a prerequisite for obtaining professional certification for
102 physicians, pharmacists, and nurses.

103 In this study, we aimed to quantitatively evaluate the performance of ChatGPT on
104 three types of national medical examinations in China, namely NMLE, NPLE and
105 NNLE. To enhance the reliability of our findings, we meticulously collected a
106 substantial corpus of real-world medical question-answer data from examinations
107 conducted from the year 2017 to 2021. We also conducted a comparative analysis of
108 the performance of different units. For cases where incorrect responses were generated,
109 we solicited feedback from domain experts and performed thorough assessment and
110 error analysis. Our study yields valuable insights for researchers and developers to
111 improve large language models' performance in the medical domain.

112

113 **Methods**

114 **Medical examination datasets**

115 We collected questions from Chinese NMLE, NPLE and NNLE from year 2017 to
116 2021. In NMLE, each exam consists of 4 units, each unit has 150 questions, for a total
117 of 600 questions. In NPLE, each exam consists of 4 units, each unit has 120 questions,
118 for a total of 480 questions. In NNLE, each exam consists of 2 units, each unit has 120
119 questions, for a total of 240 questions. Base on the requirements of the examination, a
120 correct response rate of 60% or higher is considered to meet the passing criteria. The

121 questions with figures, tables or chemical structure were manually identified and
122 excluded by a clinician with five years of clinical experience.

123 **Model setting**

124 We employed ChatGPT, an artificial intelligence chatbot built upon the generative
125 pre-trained transformer technology. The official API was utilized to invoke the chatbot,
126 with gpt-3.5-turbo as model parameter and default values for other parameters. As
127 shown in Figure 1, the input question consisted of the background description and
128 choices. To elicit diverse responses, we applied direct instruction strategy via prompt,
129 such as “Please return the most correct answer”, “Only one best option can be selected”,
130 “The correct choice is” and “This is a multiple choices question, please return the
131 correct answer”. These prompts force the model to generate the clear answer, as well as
132 the capability to distinguish between single-choice and multiple-choice questions.

133 **Evaluation**

134 For each question, the response of ChatGPT was reviewed by an experienced
135 clinicians to determine the predict answer, which was then compared with the true
136 answer. The score was calculated based on whether the answers match or not. A score
137 of 1 was awarded if there is agreement between the predict answer and true answer,
138 whereas a score of 0 was given if there is disagreement. The evaluation process has
139 been conducted on all data sets of NMLE, NPLE and NNLE over past five years.

140 **Data analysis**

141 Data process was performed in Python (version 3.9.13, Python Software
142 Foundation) using Jupyter Notebook. Statistical analysis was performed using
143 GraphPad Prism 9 Software. The significance of differences among groups was set at
144 $p < 0.05$.

145

146 **Results**

147 **Overall performance**

148 As shown in Figure 2, ChatGPT failed to pass the threshold score (0.6) in any of
149 the three types of examinations over the five years. Specifically, in the Chinese NMLE,
150 the highest recorded score was 0.5467, which was attained in both 2018 and 2021. In

151 the Chinese NPLE, the highest score was 0.5599 in 2017. In the Chinese NNLE, the
152 most impressive result was shown in 2017, with a score of 0.5897, which is also the
153 highest score in our entire evaluation. Conversely, the 2019 NPLE exam resulted in the
154 lowest score, with a recorded value of 0.4356.

155 **Detailed performance**

156 The score of each unit in the Chinese NMLE is shown in Table 1. The performance
157 of ChatGPT has varied across different units and years. In 2017 and 2020, ChatGPT
158 performed best in Unit 2. In 2018 and 2019, ChatGPT performed best in Unit 1. In the
159 2021, ChatGPT performed best in both Unit 2 and Unit 3. In 2018 and 2021, ChatGPT
160 correctly answered 328 out of 600 questions. This is because the complexity and
161 difficulty of question in each unit were vary from year to year. On average, ChatGPT
162 achieved the highest score in Unit 2 (84.6), followed by Unit 1 (79.8), Unit 3 (78.2),
163 Unit 4 (75.4).

164 **Table 1.** The score of each unit in Chinese National Medical Licensing Examination.

Year	2017	2018	2019	2020	2021	Average
Unit1 score	81	87	84	71	76	79.8
Unit2 score	96	83	75	83	86	84.6
Unit3 score	79	79	72	75	86	78.2
Unit4 score	65	79	75	78	80	75.4
Total score	321	328	306	307	328	318
Questions	600	600	597	600	600	-
Accuracy	53.50%	54.67%	51.26%	51.17%	54.67%	53.05%

165

166 The score of each unit in the Chinese NPLE is shown in Table 2. In NPLE, each
167 unit has 120 questions, and each exam has 480 questions. We identified and removed
168 the questions included figures, tables or chemical structure. Such questions appeared
169 the most in 2018 (30), followed by 2020 (22), 2017 (21), 2021 (17) and 2019 (14). On
170 average, the ChatGPT performed best in Unit 4 (62), followed by Unit 2 (58), Unit 3
171 (53) and Unit 1(52). In the year 2017, ChatGPT achieved highest score, and correctly

172 answered 257 out of 459 questions.

173 **Table 2.** The score of each unit in Chinese National Pharmacist Licensing Examination.

Year	2017	2018	2019	2020	2021	Average
Unit1 score	57	55	51	48	49	52
Unit2 score	63	54	52	60	61	58
Unit3 score	60	56	49	53	47	53
Unit4 score	77	58	51	59	65	62
Total score	257	223	203	220	222	225
Questions	459	450	466	458	463	-
Accuracy	55.99%	49.56%	43.56%	48.03%	47.95%	49.02%

174

175 The Table 3 shown the detailed score of each unit of Chinese NNLE. There are
176 totally 26 questions included figures, tables or chemical structure were removed. In
177 2017 and 2018, ChatGPT performed better in Unit2 than Unit1. Conversely, in 2019,
178 2020 and 2021, ChatGPT performed better in Unit1 than Unit2. On average, ChatGPT's
179 performance of the two units had no noticeable difference.

180 **Table 3.** The score of each unit in Chinese National Nurse Licensing Examination.

Year	2017	2018	2019	2020	2021	Average
Unit1 score	65	57	72	54	67	63
Unit2 score	73	74	57	53	63	64
Total score	138	131	129	107	130	127
Questions	234	232	238	232	238	-
Accuracy	58.97%	56.47%	54.2%	46.12%	54.62%	54.08%

181

182 In comparison, ChatGPT exhibited better proficiency in NNLE (54.08%), with
183 NMLE (53.05%) and NPLE (49.02%) following behind. The result corresponds to the
184 complexity and difficulty of the exam questions.

185 **Performance on different units and question types**

186 Figure 3 demonstrated the comparative analysis of ChatGPT's performance

187 differences across units and question types. The results shown there was no significant
188 difference in across different units in NMLE (Figure 3A), NPLE (Figure 3B), and
189 NNLE (Figure 3C). However, in the case of NPLE (Figure 3D), ChatGPT demonstrated
190 higher performance in single-choice questions compared to multiple-choice questions,
191 with a highly statistical difference ($p < 0.0001$).

192 **Performance on different subjects and topics**

193 To better understand why ChatGPT failed in the Chinese medical examination, we
194 took the 2021 NMLE exam as an example, and labeled the medical subjects and topics
195 for each question (Figure 4). The result revealed that ChatGPT excelled in clinical
196 epidemiology, human parasitology, and dermatology, with all questions answered
197 correctly. However, the model faltered in subjects such as pathology, pathophysiology,
198 public health regulations, physiology, and anatomy, with a correct rate of less than 50%.
199 Additionally, we observed that ChatGPT performed admirably in topics related to
200 molecule, health management and prevention, diagnosis and screening, but its
201 performance was lackluster in topics such as clinical manifestations, indicator values,
202 structural location, cell, and tissue. Interestingly, we found no significant difference in
203 performance on case-based questions and non-case-based questions. The result
204 provides deep insight into the strengths and weaknesses of ChatGPT in medical
205 examinations, and pave the way for future research to improve the model's capabilities
206 in this domain.

207

208 **Discussion**

209 In this study, we evaluated the performance of ChatGPT, an artificial intelligence
210 chatbot, in answering medical exam questions from Chinese NMLE, NPLE, and NNLE
211 from year 2017 to 2021.

212 **ChatGPT failed NNLE, NMLE and NPLE in China**

213 The results of our study revealed that ChatGPT was unsuccessful in meeting the
214 requirements of the three primary medical licensure assessments namely, NMLE,
215 NPLE and NNLE in China, spanning from 2017 to 2021. There are several possible
216 reasons for this.

217 Firstly, According to OpenAI, ChatGPT has been trained with the vast majority of
218 the data is in English, with only a small amount of data in other languages, like Chinese.
219 More richer training dataset allows the model to learn more knowledge. Recent studies
220 shown that ChatGPT passed United States Medical Licensing Examination [19, 20].
221 However, it failed to pass the Taiwanese Pharmacist Licensing Examination [25], and
222 performed worse than medical students on Korean-based parasitology examination [28].
223 These findings suggest that ChatGPT may require additional training data in non-
224 English languages to enhance its performance in non-English medical exams.

225 Secondly, there are differences in medical policies, regulations, and management
226 agencies across countries with different languages or cultures. In the Chinese NMLE,
227 some questions relate to healthcare policies, while in the Chinese NPLE, the entire unit
228 4 is officially designated as pharmaceutical management and regulations. These
229 questions cover topics such as drug production, market circulation, pharmaceutical
230 management, and legal regulations. The aim is to assess awareness and compliance
231 ability in clinical practice. Generally, these questions are relatively short in length and
232 clear in meaning. While ChatGPT has acquired a wealth of knowledge on healthcare
233 policies from English-speaking countries due to its extensive English dataset, it may
234 encounter difficulties in correctly understanding the healthcare policies of non-English-
235 speaking countries, leading to erroneous responses to related questions. Additionally,
236 healthcare policies undergo regular updates over time, making such questions more
237 challenging.

238 Thirdly, in some questions, the task requires reading the question and selecting the
239 most suitable answer from five given choices. However, there may be more than one
240 choice that can answer the question, and in such cases, ChatGPT may generate multiple
241 answers, including the correct answer. As this is a single-choice question, ChatGPT is
242 forced to select a single choice as answer, which can limit its content comprehension
243 ability and lead to incorrect answer.

244 **The potential of large language model in medical education**

245 As a significant milestone in the development of artificial intelligence, ChatGPT
246 driven by large language model has powerful capability in language understanding and

247 content generation. With its remarkable potential, ChatGPT could be a valuable
248 resource in acquiring medical knowledge and learning clinical skills for students, and
249 serve as an informative assistant in preparing teaching materials and evaluating course
250 projects for teachers.

251 In our study, ChatGPT has achieved an accuracy rate of over 50% in most of the
252 exams, indicating a significant potential for ChatGPT in medical education. Previous
253 study shown in the Chinese Rural General Medical Licensing Examination, only 55%
254 of students were able to pass the written examination [29]. In China, the significant
255 healthcare burden necessitates a vast number of licensed clinical staff and healthcare
256 providers. However, the rigorous examinations lead to low pass rates, exacerbating the
257 shortage of licensed practitioners, especially in rural areas. The large language model
258 presents a promising avenue for enhancing medical education and advancing healthcare
259 reform, with the potential to reduce medical burden.

260 Finally, the advancement of artificial intelligence (AI), specifically large language
261 models, in medical education needs public benchmarking datasets and fair evaluation
262 metrics for performance assessment. There is also a need to interact with human experts
263 from multiple dimensions and obtain continuous feedback. In addition, the use of such
264 model must also consider data privacy, cognitive bias, and comply with regulations.

265 **Limitations**

266 Our study has some limitations. First, the questions of China NMLE, NPLE, and
267 NNLE are all multiple-choice format. While this format meets our study purposes, it
268 did not fully showcase the content generation capabilities of ChatGPT. In the future, it
269 would be beneficial to include more open-ended questions. Second, we evaluated the
270 performance of ChatGPT in medical examinations with zero-shot learning. However,
271 better performance may be achieved by incorporating knowledge-enhanced training
272 methods. Third, the different variations of prompt may impact ChatGPT's response,
273 leading to diverse answers. Therefore, it is imperative to develop innovative techniques
274 that can generate more consistent and trustworthy responses in the future. Finally,
275 further investigation is needed to determine the underlying factors contributing to this
276 substandard performance, and to explore broader application of ChatGPT in medical

277 education and clinical decision-making support.

278

279 **Conclusion**

280 In conclusion, we evaluated the performance of ChatGPT on three types of national
281 medical examinations in China, including NMLE, NPLE, and NNLE from year 2017
282 to 2021. The results indicated ChatGPT failed to meet the official pass criteria of 60%
283 correct response rate in any of the three types of examinations over the five years. The
284 performance of ChatGPT varied across different units and years, with the highest score
285 achieved in NNLE of year 2017. ChatGPT exhibited relatively better proficiency in
286 NNLE, with NMLE and NPLE following closely behind. ChatGPT performed well in
287 a range of subject areas, including clinical epidemiology, human parasitology, and
288 dermatology, as well as in various medical topics such as molecules, health
289 management and prevention, diagnosis and screening.

290

291 **Funding**

292 This work was supported by the National Natural Science Foundation of China
293 (32270690 and 32070671).

294

295 **Data Availability**

296 The data analyzed and reported in this study are available from the authors upon
297 reasonable request.

298

299 **Authors' Contributions**

300 HZ, JL, EW, RW, JL and BS involved in the study conceptualization. JL collected and
301 preprocessed the data. HZ conducted data analysis, results interpretation and
302 manuscript preparation. HZ, JL EW, RW, JL and BS contributed to the review and
303 editing of the manuscript. BS supervised the study. All authors read and approved the
304 final manuscript.

305

306 **Acknowledgments**

307 We thank OpenAI, which allows free access to ChatGPT.

308

309 **Conflict of Interest**

310 The authors declare no conflict of competing interests.

311

312 **References**

313

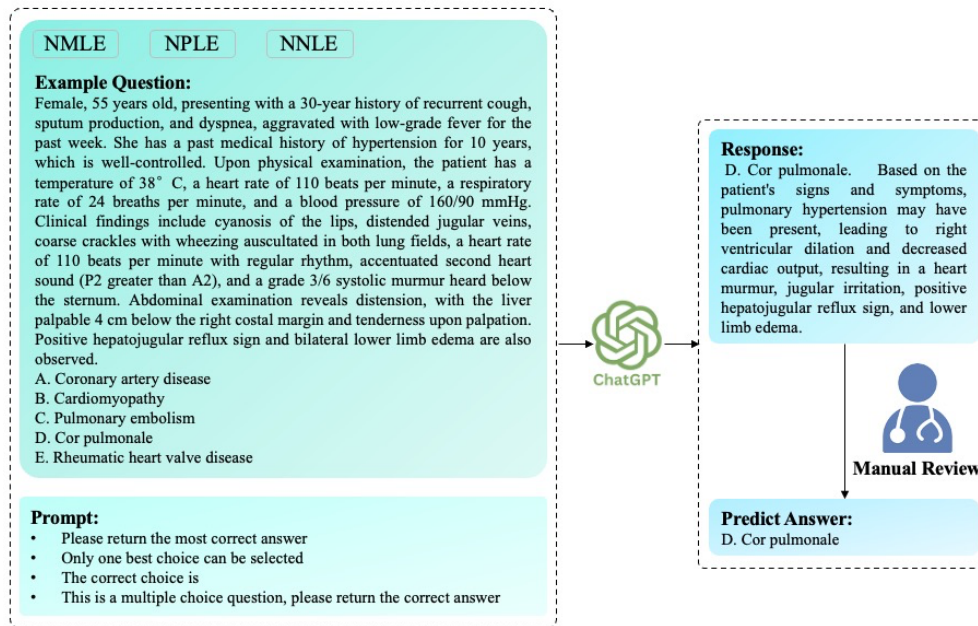
- 314 1. Bhinder, B., et al., *Artificial Intelligence in Cancer Research and Precision Medicine*. Cancer
315 Discov, 2021. **11**(4): p. 900-915.
- 316 2. Moor, M., et al., *Foundation models for generalist medical artificial intelligence*. Nature,
317 2023. **616**(7956): p. 259-265.
- 318 3. van Dis, E.A.M., et al., *ChatGPT: five priorities for research*. Nature, 2023. **614**(7947): p.
319 224-226.
- 320 4. Sarink, M.J., et al., *A study on the performance of ChatGPT in infectious diseases clinical*
321 *consultation*. Clin Microbiol Infect, 2023.
- 322 5. Lee, T.C., et al., *ChatGPT Answers Common Patient Questions About Colonoscopy*.
323 Gastroenterology, 2023.
- 324 6. Young, J.N., et al., *The utility of ChatGPT in generating patient-facing and clinical*
325 *responses for melanoma*. J Am Acad Dermatol, 2023.
- 326 7. Luo, R., et al., *BioGPT: generative pre-trained transformer for biomedical text generation*
327 *and mining*. Brief Bioinform, 2022. **23**(6).
- 328 8. Zhang, Y., et al. *DIALOGPT : Large-Scale Generative Pre-training for Conversational*
329 *Response Generation*. 2020. Online: Association for Computational Linguistics.
- 330 9. Das, A., et al. *Conversational Bots for Psychotherapy: A Study of Generative Transformer*
331 *Models Using Domain-specific Dialogues*. 2022. Dublin, Ireland: Association for
332 Computational Linguistics.
- 333 10. Komorowski, M., M. Del Pilar Arias Lopez, and A.C. Chang, *How could ChatGPT impact*
334 *my practice as an intensivist? An overview of potential applications, risks and limitations*.
335 Intensive Care Med, 2023.
- 336 11. Munoz-Zuluaga, C., et al., *Assessing the Accuracy and Clinical Utility of ChatGPT in*
337 *Laboratory Medicine*. Clin Chem, 2023.
- 338 12. Yang, H., *How I use ChatGPT responsibly in my teaching*. Nature, 2023.
- 339 13. Abd-Alrazaq, A., et al., *Large Language Models in Medical Education: Opportunities,*
340 *Challenges, and Future Directions*. JMIR Med Educ, 2023. **9**: p. e48291.
- 341 14. Jin, Q., R. Leaman, and Z. Lu, *Retrieve, Summarize, and Verify: How Will ChatGPT Affect*
342 *Information Seeking from the Medical Literature?* J Am Soc Nephrol, 2023.
- 343 15. Kooroor, J.G., A.K. Gupta, and S. Bacchi, *ChatGPT: effective writing is succinct*. BMJ, 2023.
344 **381**: p. 1125.
- 345 16. Shafiee, A., *Matters arising: authors of research papers must cautiously use ChatGPT for*
346 *scientific writing*. Int J Surg, 2023.
- 347 17. Gao, C.A., et al., *Comparing scientific abstracts generated by ChatGPT to real abstracts*

- 348 *with detectors and blinded human reviewers.* NPJ Digit Med, 2023. **6**(1): p. 75.
- 349 18. Salvagno, M., F.S. Taccone, and A.G. Gerli, *Can artificial intelligence help for scientific*
350 *writing?* Crit Care, 2023. **27**(1): p. 75.
- 351 19. Kung, T.H., et al., *Performance of ChatGPT on USMLE: Potential for AI-assisted medical*
352 *education using large language models.* PLOS Digit Health, 2023. **2**(2): p. e0000198.
- 353 20. Gilson, A., et al., *How Does ChatGPT Perform on the United States Medical Licensing*
354 *Examination? The Implications of Large Language Models for Medical Education and*
355 *Knowledge Assessment.* JMIR Med Educ, 2023. **9**: p. e45312.
- 356 21. Bhayana, R., S. Krishna, and R.R. Bleakney, *Performance of ChatGPT on a Radiology Board-*
357 *style Examination: Insights into Current Strengths and Limitations.* Radiology, 2023. **307**(5):
358 p. e230582.
- 359 22. Giannos, P., *Evaluating the limits of AI in medical specialisation: ChatGPT's performance*
360 *on the UK Neurology Specialty Certificate Examination.* BMJ Neurol Open, 2023. **5**(1): p.
361 e000451.
- 362 23. Humar, P., et al., *ChatGPT is Equivalent to First Year Plastic Surgery Residents: Evaluation*
363 *of ChatGPT on the Plastic Surgery In-Service Exam.* Aesthet Surg J, 2023.
- 364 24. Weng, T.L., et al., *ChatGPT failed Taiwan's Family Medicine Board Exam.* J Chin Med Assoc,
365 2023.
- 366 25. Wang, Y.M., H.W. Shen, and T.J. Chen, *Performance of ChatGPT on the Pharmacist*
367 *Licensing Examination in Taiwan.* J Chin Med Assoc, 2023.
- 368 26. Seghier, M.L., *ChatGPT: not all languages are equal.* Nature, 2023. **615**(7951): p. 216.
- 369 27. Wang, X., *Experiences, challenges, and prospects of National Medical Licensing*
370 *Examination in China.* BMC Med Educ, 2022. **22**(1): p. 349.
- 371 28. Huh, S., *Are ChatGPT's knowledge and interpretation ability comparable to those of*
372 *medical students in Korea for taking a parasitology examination?: a descriptive study.* J
373 Educ Eval Health Prof, 2023. **20**: p. 1.
- 374 29. Han, X., et al., *Performance of China's new medical licensing examination for rural general*
375 *practice.* BMC Med Educ, 2020. **20**(1): p. 314.

376

377

378 **Figure Legends**

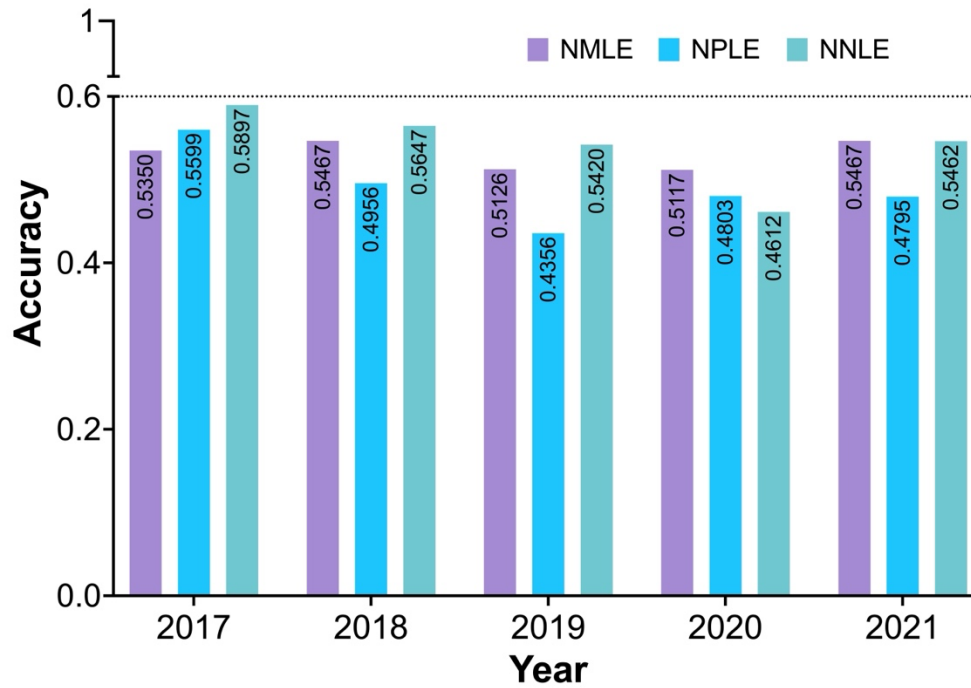


379

380

381 **Figure 1.** The overview of interaction with ChatGPT. The question included
382 background description and choices from three national licensing examinations,
383 including Chinese National Medical Licensing Examination (NMLE), National
384 Pharmacist Licensing Examination (NPLE) and National Nurse Licensing Examination
385 (NNLE). The prompt was designed to force a clear answer, as well as the ability to
386 recognize single-choice or multiple-choice questions. The response of ChatGPT were
387 manually reviewed by an experienced clinician to determine the answer. The correct
388 answer to this question is “D. Cor pulmonale”. It should be noted that while English
389 text was shown in the figure, the experiment itself used Chinese text as both the input
390 and output language.

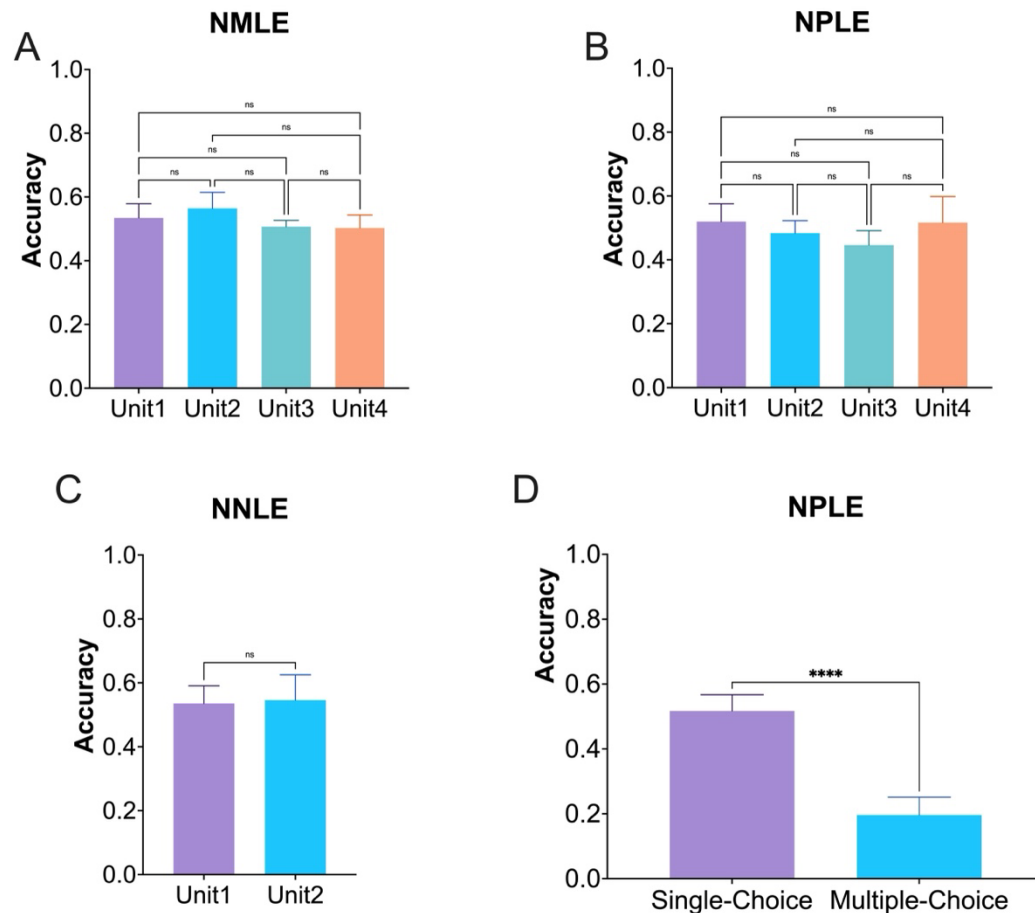
391



392

393 **Figure 2.** The performance of ChatGPT of three national licensing examinations over
394 a period of five years from 2017 to 2022. The examinations included Chinese National
395 Medical Licensing Examination (NMLE), National Pharmacist Licensing Examination
396 (NPLE) and National Nurse Licensing Examination (NNLE).

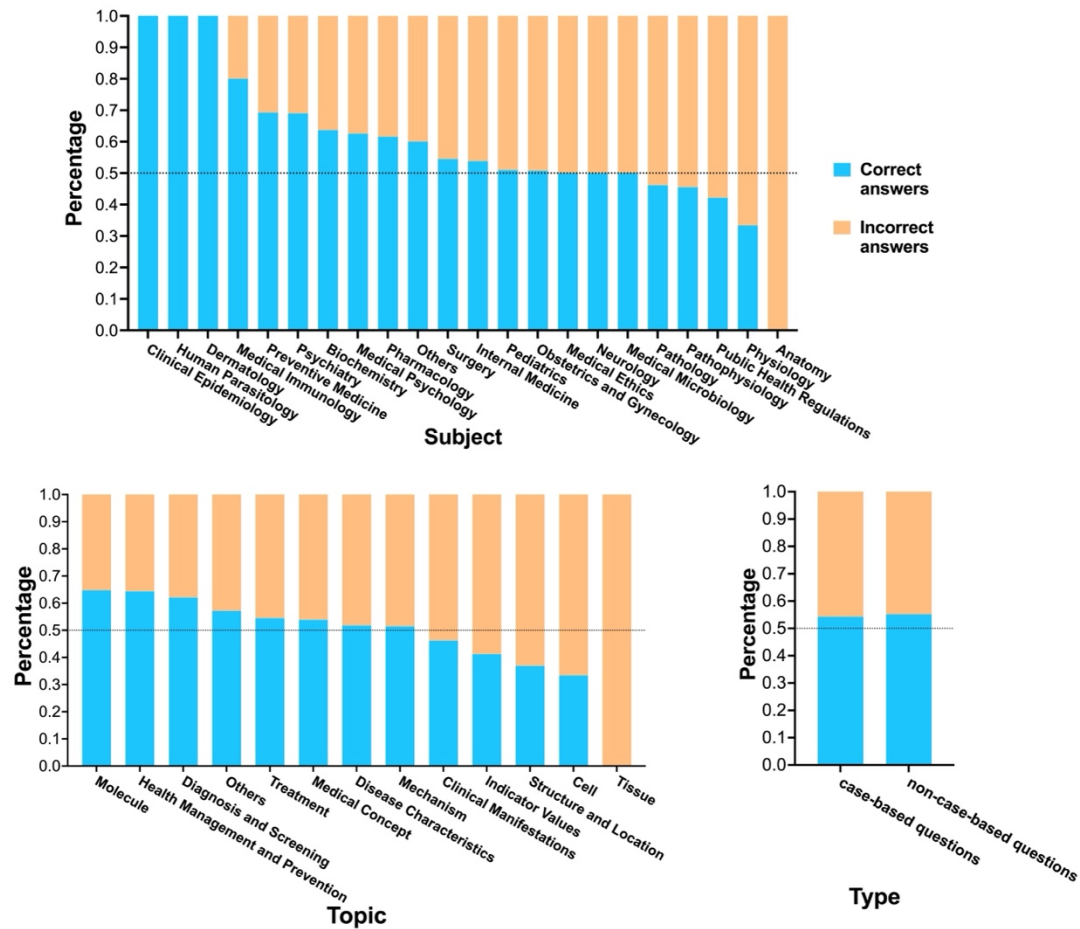
397



398

399 **Figure 3.** The performance of ChatGPT on different units and question types. For
400 different units, there were no significant difference among (A) Chinese National
401 Medical Licensing Examination (NMLE), (B) National Pharmacist Licensing
402 Examination (NPLE), and (C) National Nurse Licensing Examination (NPLE). (D)
403 However, ChatGPT demonstrated higher performance in single-choice questions than
404 multiple-choice questions with a highly significant difference (ns, no significant
405 difference, ****p < 0.0001).

406



407

408 **Figure 4.** The performance of ChatGPT on different subjects, topics and types of

409 questions in the 2021 NMLE exam.

410