

1 **How does ChatGPT4 preform on Non-English National Medical Licensing**

2 **Examination? An Evaluation in Chinese Language**

3 Changchang Fang^{*1,2}, M.D., Jitao Ling^{*1}, M.D., Jing Zhou¹, M.D., Yue Wang^{3,4},
4 M.D., Xiaolin Liu^{3,4}, M.D., Yuan Jiang^{3,4}, Ph.D, Yifan Wu¹,M.D., Yixuan Chen¹, M.D.,
5 Zhichen Zhu¹, M.D., Jianyong Ma⁵, M.D., Ziwei Yan⁶, M.PT., Peng Yu^{1,7}, Xiao Liu,
6 M.D.^{3,4,7}

7 1Department of Endocrine, the Second Affiliated Hospital of Nanchang University,
8 Jiangxi, China

9 2Queen Mary College, Nanchang University, Jiangxi, China

10 3Department of Cardiology, Sun Yat-sen Memorial Hospital of Sun Yat-sen
11 University, Guangzhou, China

12 4Guangdong Province Key Laboratory of Arrhythmia and Electrophysiology,
13 Guangzhou, China

14 5Department of Pharmacology and Systems Physiology, University of Cincinnati
15 College of Medicine, Cincinnati, United States

16 6Provincial University Key Laboratory of Sport and Health Science, School of
17 Physical Education and Sport Sciences, Fujian Normal University, Fuzhou, China

18 7Institute for the Study of Endocrinology and Metabolism in Jiangxi, the Second
19 Affiliated Hospital of Nanchang University, Jiangxi, China

20 *Co-first author

21 Word counts: 2952 (excluding abstract and references)

22 **Corresponding author:**

23 ✉Peng Yu, Email: yu8220182@163.com
 24 Department of Endocrine, the Second Affiliated Hospital of Nanchang University,
 25 Jiangxi, China
 26 ✉Xiao Liu, Email: Liux587@mail.sysu.edu.cn
 27 Department of Cardiology, Sun Yat-sen Memorial Hospital of Sun Yat-sen University,
 28 Guangzhou, China
 29 Institute for the Study of Endocrinology and Metabolism in Jiangxi, the Second
 30 Affiliated Hospital of Nanchang University, Jiangxi, China
 31

32 **Abstract**

33 **Background:** ChatGPT, an artificial intelligence (AI) system powered by large-scale
 34 language models, has garnered significant interest in the healthcare. Its performance
 35 dependent on the quality and amount of training data available for specific language.
 36 This study aims to assess the of ChatGPT's ability in medical education and clinical
 37 decision-making within the Chinese context.

38 **Methods:** We utilized a dataset from the Chinese National Medical Licensing
 39 Examination (NMLE) to assess ChatGPT-4's proficiency in medical knowledge
 40 within the Chinese language. Performance indicators, including score, accuracy, and
 41 concordance (confirmation of answers through explanation), were employed to
 42 evaluate ChatGPT's effectiveness in both original and encoded medical questions.
 43 Additionally, we translated the original Chinese questions into English to explore
 44 potential avenues for improvement.

45 **Results:** ChatGPT scored 442/600 for original questions in Chinese, surpassing the
 46 passing threshold of 360/600. However, ChatGPT demonstrated reduced accuracy in
 47 addressing open-ended questions, with an overall accuracy rate of 47.7%. Despite this,
 48 ChatGPT displayed commendable consistency, achieving a 75% concordance rate
 49 across all case analysis questions. Moreover, translating Chinese case analysis
 50 questions into English yielded only marginal improvements in ChatGPT's
 51 performance ($P=0.728$).

52 **Conclusion:** ChatGPT exhibits remarkable precision and reliability when handling
 53 the NMLE in Chinese language. Translation of NMLE questions from Chinese to

54 English does not yield an improvement in ChatGPT's performance.

55 **Introduction**

56 AI (Artificial Intelligence) has gained significant influence in recent years,
 57 simulating human intelligence and cognitive processes to tackle complex problems¹.
 58 Trained on specific datasets, AI systems enhance prediction accuracy and address
 59 complex challenges²⁻⁴, assisting doctors in rapidly searching through medical data,
 60 augmenting creativity, and facilitating error-free decision-making^{5, 6}. ChatGPT is a
 61 Large Language Model that predicts word sequences based on context and generates
 62 novel sequences resembling natural human language. These novel sequences have not
 63 been previously observed by other AI systems⁷.

64 ChatGPT shows promise in medical education, performing well in Certified Public
 65 Accountant (CPA) exams and generating accurate responses to complex inputs⁸.
 66 Applied in the United States Medical Licensing Examination and South Korean
 67 parasitology exams, ChatGPT demonstrates significant advancements, despite
 68 discrepancies with medical students' scores⁹. However, ChatGPT's proficiency relies
 69 on available training data quality and quantity in the languages, and most of them is in
 70 English. With over 1.3 billion speakers, the amount and quality of training data in
 71 Chinese language may not be comparable to English, necessitating further research
 72 into ChatGPT's performance in Chinese medical information. The Chinese National
 73 Medical Licensing Examination (NMLE) is a legally mandated qualification for
 74 doctors¹⁰. This comprehensive, standardized assessment poses conceptually and
 75 linguistically challenging questions across medical domains, which makes it an
 76 excellent input for ChatGPT in clinical decision-making.

77 Give this background, this study aims to evaluate ChatGPT's performance on the
78 Chinese NMLE conducted within the Chinese context.

79 **Methods**

80 **Artificial Intelligence**

81 ChatGPT is an advanced language model that leverages self-attention mechanisms
82 and extensive training data to deliver natural language responses within
83 conversational settings. Its primary strengths encompass managing long-range
84 dependencies and producing coherent, contextually appropriate responses.
85 Nevertheless, it is essential to recognize that GPT-4 is a server-based language model
86 without internet browsing or search capabilities. Consequently, all generated
87 responses rely solely on the abstract associations between words, or "tokens," within
88 its neural network⁷.

89 **Input source**

90 The official website does not release the 2022 NMLE test questions. However, a
91 complete set of 600 questions, with a total value of 600 points, is available online
92 (Supplemental S1) and considered as original questions. These questions are divided
93 into four units, with each question worth one point.

94 The four units encompass the following areas: Unit 1 assesses medical knowledge,
95 policies, regulations, and preventive medicine; Unit 2 focuses on cardiovascular,
96 urinary, muscular, and endocrine systems; Unit 3 addresses digestive, respiratory, and

97 other related systems; while Unit 4 evaluates knowledge of female reproductive
98 systems, pediatric diseases, and mental and nervous systems.

99 All inputs provided to the GPT-4 model are valid samples that do not belong to the
100 training dataset, as the database has not been updated since September 2021,
101 predating the release of these questions, this was further confirmed by randomly spot
102 checking the inputs. To facilitate research efforts, the 600 questions have been
103 organized into distinct categories based on their question type and units.

104 1.Common Questions (n=340): These questions are distributed across all units,
105 including Unit 1 (n=108), Unit 2 (n=82), Unit 3 (n=79), and Unit 4 (n=71). They aim
106 to evaluate basic science knowledge in physiology, biochemistry, pathology, and
107 medical humanities. Each question has four choices, and the AI must select the single
108 correct answer. An example from Unit 1 is: "What type of hypoxia is likely to be
109 caused by long-term consumption of pickled foods? A. Hypoxia of blood type B.
110 Hypoxia of tissue type C. Circulatory hypoxia D. Anoxic hypoxia E. Hypoxia of
111 hypotonic type."

112 2.Case Analysis Questions (n=260): These questions are also distributed across all
113 units, including Unit 1 (n=42), Unit 2 (n=68), Unit 3 (n=71), and Unit 4 (n=79).
114 Employed in clinical medicine, it examines and evaluates patient cases through a
115 thorough review of medical history, symptoms and diagnostic findings to determine a
116 diagnosis and treatment plan. Each question has four choices, and the AI must select
117 the single correct answer. An example from Unit 1 is: "A 28-year-old male complains
118 of muscle and joint pain in his limbs three days after diving. He experienced

119 respiratory equipment failure during diving three days ago and immediately ascended
120 rapidly to the surface. Subsequently, he experienced symptoms such as dizziness,
121 orientation disorder, nausea, and vomiting. After rest and oxygen inhalation, the
122 symptoms improved, but he continued to experience persistent muscle spasms,
123 convulsions, and joint pain in his limbs. Therefore, what is the most likely cause of
124 the patient's pain? A. Chronic inflammation and cell infiltration B. Stress ulcers C.
125 Local tissue coagulative necrosis D. Increased carbon dioxide concentration in the
126 blood E. Gas embolism in the blood vessel lumen."

127 **Scoring**

128 We assembled a dataset of NMLE questions and their corresponding answers,
129 maintaining validity by cross-verification with senior medical professionals. This
130 dataset was used to evaluate ChatGPT's performance on the exam by comparing its
131 responses to the standard answers and calculating the scores it achieved. A high score
132 would indicate that ChatGPT effectively tackled this task.

133 **Encoding**

134 To better reflect the actual clinical situation, we modified the case analysis
135 questions to be open-ended. Questions were formatted by deleting all the choices and
136 adding a variable lead-in imperative or interrogative phrase, requiring ChatGPT to
137 provide a rationale for the answer choice. Examples include: "What could be the most
138 plausible explanation for the patient's nocturnal symptoms? Justify your answer for

139 each option," and "Which mechanism is most likely responsible for the most fitting
140 pharmacotherapy for this patient? Provide an explanation for its correctness."

141 However, a unique subset of questions could not be encoded in the same manner.
142 These questions required selecting one provided choice, so we transformed them into
143 a special form (n=3). For example, the original question, "Which can inhibit insulin
144 secretion? A. Increased free fatty acids in blood B. Increased gastric inhibitory peptide
145 secretion C. Sympathetic nerve excitation D. Growth hormone secretion increases"
146 was encoded as "Can an increase in free fatty acids in the blood, an increase in gastric
147 inhibitory peptide secretion, an increase in sympathetic nerve excitation, or an
148 increase in growth hormone secretion inhibit insulin secretion?" This encoding was
149 present only in Unit 1. To minimize memory retention bias, a new chat session was
150 initiated for each inquiry.

151 **Adjudication**

152 In our study, AI outputs from the two types of encoders were independently scored
153 for Accuracy and Concordance by two physicians blinded to each other's assessments.
154 Scoring was based on predefined criteria (**Supplemental S2**). To train the physician
155 adjudicators, who were not blinded to each other, a subset of 20 questions was used.
156 ChatGPT's responses were classified into three categories: accurate, inaccurate, and
157 indeterminate. Accurate responses indicated that ChatGPT provided the correct
158 answer, while inaccurate responses encompassed no answer, incorrect answers, or
159 multiple answers with incorrect options. Indeterminate responses implied that the AI

160 output did not provide a definitive answer selection or believed there was insufficient
161 information to do so. Concordance was defined as when ChatGPT's explanation
162 confirmed its provided answer, while discordant explanations contradicted the answer.

163 To minimize within-item anchoring bias, adjudicators first evaluated accuracy for
164 all items, followed by concordance. Two physicians were blinded to each other's
165 evaluations. In cases of discrepancy, a third physician adjudicator was consulted.
166 Ultimately, 17 items (2.7% of the dataset) required the intervention of a third
167 physician adjudicator. The interrater agreement between the physicians was assessed
168 using the Cohen kappa (κ) statistic for the questions (**Supplemental S3**).

169 A schematic overview of the study protocol is provided in **Fig 1**.

170 **Translation**

171 To evaluate if translating questions from Chinese to English could enhance
172 ChatGPT's performance, we utilized ChatGPT to translate unencoded case analysis
173 questions. We then assessed ChatGPT's performance on the translated exam by
174 comparing its responses to standard answers and calculating its scores. We compared
175 the scores obtained from the original questions to those from the translated questions
176 and employed the chi-square test to determine performance improvement.

177 **Result**

178 **ChatGPT passed Chinese NMLE with a high score.**

179 In the Chinese NMLE, 442 (73.67%) out of 600 items were correctly answered by
180 ChatGPT, which is significantly higher than the passing threshold (360) defined by
181 official agencies.

182 The score of each unit is shown in Fig 2. ChatGPT's performance varied across the
183 four units of questions, with a highest accuracy in Unit 4 (76.0%), followed by Unit 3
184 (74.7%), Unit 1 (74.0%) and Unit 2 (70.0%), while there was no statistically
185 difference among four units ($\chi^2 = 0.66$, $p = 0.883$).

186 **ChatGPT's performance declines when handling encoded questions**

187 Test questions were encoded as open-ended for case analysis questions, simulating
188 scenarios where a student poses a common medical question without answer choices
189 or a doctor diagnoses a patient based on multimodal clinical data (e.g., symptoms,
190 history, physical examination, laboratory values). The accuracy was 40.5%, 60.3%,
191 42.3%, and 34.2% for Units 1, 2, 3, and 4, respectively (**Fig 3A**). Compared to the
192 original questions, the accuracy of the encoded questions decreased by 40.5%, 9.7%,
193 32.4%, and 41.8% for Units 1, 2, 3, and 4, respectively (**Fig 3B**).

194 These findings demonstrate that while ChatGPT's ability to answer questions in
195 Chinese as applied to common medical situations is commendable, there is still room
196 for improvement. During the adjudication stage, physician agreement was good for
197 open-ended questions (with a κ range from 0.83 to 1.00).

198 *ChatGPT demonstrates high internal concordance*

199 Concordance is a measure of the agreement or similarity between the option
200 selected by AI and its subsequent explanation. The results indicated that ChatGPT
201 maintained a >75% concordance across all questions, and this high level of
202 concordance was consistent across all four units (**Fig 4**). Furthermore, we examined

203 the concordance difference between correct and incorrect answers, discovering that
204 concordance was perfect and significantly higher among accurate responses compared
205 with inaccurate ones (85% vs. 59.5%, $p < 0.005$) (**Fig 4**).

206 These findings suggest that ChatGPT exhibits a high level of answer-explanation
207 concordance in Chinese, which can be attributed to the strong internal consistency of
208 its probabilistic language model.

209 **Translating the input into English may not improve ChatGPT's performance.**

210 After translating the original case analysis questions in Chinese into English to
211 explore a potential way to improve ChatGPT's performance, the improvement for
212 Units 1, 2, 3, and 4 was minimal, with only one point gained in each unit. The total
213 number of correct answers increased from 256 to 260. The accuracy improvement for
214 translated case analysis questions was subtle ($\chi^2 = 0.1206$, $P = 0.728$). This suggests
215 that ChatGPT's performance when facing questions in Chinese may not be improved
216 by translating them into English, and solely building a database in English while
217 translating other languages into English may not be an effective approach.

218 **Discussion**

219 In present study, we firstly investigated ChatGPT's performance on the Chinese
220 NMLE. Our findings can be summarized under two major themes: (1) ChatGPT's
221 score is satisfactory but requires improvement when addressing questions posed in the
222 Chinese language; and (2) Translation into English showed slight performance
223 improvement. This study provides new evidence for the ability of ChatGPT in

224 medical education and clinical decision-making within the Chinese context, offering
225 valuable insights into the applicability of AI language models for non-English medical
226 education settings and laying the groundwork for future research in this area.

227 **ChatGPT's performance in the Chinese NMLE is acceptable, yet further**
228 **improvement**

229 In the Chinese NMLE, ChatGPT achieved a score of 442 (73.67%), exceeding the
230 passing requirement of 360 points for the Chinese language. In the 2022 NLME, the
231 average score of 65 medical students was 412.7 (68.7%), with a minimum score of
232 295 (49.2%) and a maximum score of 474 (79.0%). According to the statistics, the
233 national pass rate for the exam in 2022 was 55%. When compared to medical students
234 who have undergone a traditional 5-year medical education and a one-year internship,
235 ChatGPT's performance is currently satisfactory, however, there is still potential for
236 improvement. Several underlying reasons may be responsible. 1) Limitations in
237 training data: If ChatGPT's training data contains less information about the Chinese
238 medical field, its performance when handling Chinese medical questions could be
239 impacted, resulting in a lower accuracy rate for such queries. 2) Knowledge updates:
240 With a knowledge cutoff date in September 2021, the most recent developments in the
241 Chinese medical field may not have been adequately learned by the model, affecting
242 its accuracy when answering Chinese medical questions.

243 **ChatGPT's accuracy can be improved by addressing data limitations, refining its**
244 **architecture, and using domain-specific knowledge**

245 Moreover, we observed that outputs with high accuracy exhibited high concordance,

246 while lower accuracy was associated with reduced concordance. Consequently, we
247 speculate that ChatGPT's inaccurate responses primarily arise from missing
248 information, leading to indecision in the AI rather than adherence to an incorrect
249 answer. Language models like ChatGPT are built on vast amounts of text, and their
250 accuracy depends on the quality and diversity of their training data¹¹. When the model
251 encounters scenarios with limited or underrepresented data, its performance may
252 suffer, leading to indecision or inaccurate responses. To address this issue, one could
253 consider expanding the training data to cover a broader range of contexts or refining
254 the model's architecture to handle uncertainty more effectively. Additionally,
255 incorporating domain-specific knowledge and data sources can help improve the
256 model's performance in specialized areas.

257 **ChatGPT performs best in English, with accuracy affected by translation issues**
258 **and data limitations in other languages**

259 ChatGPT's performance may differ across different language and this variation can
260 be attributed to factors such as the quality and quantity of training data available in
261 different language¹². Among these languages, a better performance of ChatGPT with
262 English task descriptions is reported¹². However, we observed a slight improvement in
263 ChatGPT's performance after translating questions into English, suggesting that
264 relying solely on English databases or translations might not be the an effective
265 approach in improving the abilities for medical tasks. Several potential reasons may
266 be responsible: 1) Translation limitations: When translating questions, some nuances
267 or specific terms may be lost or inaccurately translated, which could impact the AI's

268 understanding and subsequently its performance¹³. Additionally, some languages may
 269 have unique expressions or cultural context that are difficult to convey accurately in
 270 English, leading to potential misunderstandings or misinterpretations. 2) AI model's
 271 abilities: Language models like ChatGPT rely on the quality and quantity of training
 272 data available in each language. If the model has been trained extensively with
 273 English data, it may perform better when handling English text. For other languages,
 274 the AI model's performance could be affected by insufficient or lower-quality training
 275 data, leading to less accurate responses.

276 Regarding the best-performing languages, according to the answer of ChatGPT,
 277 English typically yields the highest accuracy since most training data is in English.
 278 Other widely spoken languages with a substantial number of online resources, such as
 279 Chinese, Spanish, French, and German, may also exhibit relatively better performance
 280 in terms of accuracy. However, these results may vary depending on the model and
 281 specific task. To assess the AI model's performance for a particular language, a
 282 targeted evaluation might be necessary.

283 **GPT-4 shows progress, but addressing healthcare standards, ethics, and culture** 284 **is crucial for AI integration in medicine**

285 ChatGPT-3.5 achieved near-passing threshold accuracy of 60% on the United
 286 States Medical Licensing Exam⁷. Furthermore, our previous study also showed a
 287 similar performance by ChatGPT-3.5 on the Clinical Medicine Entrance Examination
 288 for Chinese Postgraduates (scored 153.5/300, 51%) in Chinese language¹⁴. In the
 289 present study, a significant higher score was found by GPT-4 (scored 442/600, 73.6%).

290 This improvement may be attributed to differences in model sizes and training data.
 291 GPT-4's larger model size enables it to handle more complex tasks and generate more
 292 accurate responses due to its extensive training dataset, broader knowledge base, and
 293 improved contextual understanding¹⁵. On the other hand, Chinese medical licensing
 294 exams have many common-sense questions and fewer case analysis questions than
 295 United States Medical Licensing Exam, which may be another reasons for the
 296 relatively high pass rates.

297 Despite the promising potential of AI in medicine, it also faces several challenges.
 298 The development of standards for AI use in healthcare is still required^{16, 17},
 299 encompassing clinical care, quality, safety, malpractice, and communication
 300 guidelines. Moreover, the implementation of AI in healthcare necessitates a shift in
 301 medical culture, posing challenges for both medical education and practice. Ethical
 302 considerations, such as data privacy, informed consent, and bias prevention, must also
 303 be addressed to ensure that AI is employed ethically and for the benefit of patients.

304 **Limitations**

305 Several limitations should be noted. Firstly, the clinical tasks are highly
 306 complicated, the exams cannot fully stimulate the problems in clinical practices.
 307 Secondly, the limited input sample size may preclude us performing the depth and
 308 range of analyses, which potentially limiting the generalizability of findings. Thus,
 309 before large-scale applications of Large Language Model-based AI in medical
 310 education or clinical practice, their utility should be further studied in real-world

311 condition.

312 **Conclusion**

313 ChatGPT demonstrated impressive performance in the Chinese NMLE in Chinese
314 language, exceeding the passing threshold and exhibiting high internal consistency.
315 Nevertheless, its performance waned when faced with open-end encoded questions.
316 Translation into English did not substantially boost its performance. The findings
317 emphasize the ChatGPT's ability of comprehensible reasoning in medical education
318 and clinical decision-making in Chinese.

319 **Data statement**

320 All data generated or analyzed during this study are included in this published article
321 [and its supplementary information files].

322 **Acknowledgments**

323 We acknowledge the ChatGPT for polishing our manuscript.

324 **Author contributions**

325 Xiao Liu was responsible for the entire project and revised the draft. Changchang
326 Fang, Jitao Ling, Jing Zhou, Xiaoling Liu, Yixuan Chen, Zhichen Zhu and Yuan Jiang
327 performed the study selection, data extraction, statistical analysis, and interpretation
328 of the data. Changchang Fang and Xiao Liu drafted the first version of the manuscript.
329 All authors participated in the interpretation of the results and prepared the final
330 version of the manuscript.

331 **Funding**

332 This work was supported in part by the National Natural Science Foundation

333 of China (X.L., 82100347; 82100869 to P.Y.), Basic and Applied Basic Research
334 Project of Guangzhou (202201011395 to X.L), Natural Science Foundation of
335 Guangdong Province (X.L., 2022A1515010582), Natural Science Foundation in
336 Jiangxi Province grant [No. 20224ACB216009 and No. 20212BAB216051 to J.Z., No.
337 20212BAB216047 and No. 202004BCJL23049 to P.Y.; No. 202004BCJL23049 to
338 P.Y.] Science and Technology Projects in Guangzhou (No. 202102010007)
339 None of the funding institutions had a role in design, methods, subject recruitment,
340 data collections, analysis, and preparation of the article. **Declarations**
341 Ethics approval This is a systematic review and meta-analysis. No ethical approval is
342 required.
343 **Conflict of interest**
344 All authors declare no competing interests.
345

346 Reference

- 347 1. Haleem A, Javaid M, Khan IH. Current status and applications of Artificial
348 Intelligence (AI) in medical field: An overview. *Current Medicine Research and*
349 *Practice*. 2019;9(6):231-237.
- 350 2. Haleem A, Vaishya R, Javaid M, et al. Artificial Intelligence (AI) applications in
351 orthopaedics: An innovative technology to embrace. *Journal of Clinical Orthopaedics*
352 *and Trauma*. 2019(0976-5662 (Print)).
- 353 3. Jha S, Topol EJ. Information and artificial intelligence. *Journal of the American*
354 *College of Radiology*. 2018;15(3):509-511.
- 355 4. Lupton ML. Some ethical and legal consequences of the application of artificial
356 intelligence in the field of medicine. 2018.
- 357 5. Murdoch TB, Detsky AS. The inevitable application of big data to health care.
358 *JAMA*. 2013(1538-3598 (Electronic)).
- 359 6. Misawa M, Kudo S-e, Mori Y, et al. Artificial intelligence-assisted polyp
360 detection for colonoscopy: initial experience. *Gastroenterology*.
361 2018;154(8):2027-2029.
- 362 7. Kung TH, Cheatham M, Medenilla A, et al. Performance of ChatGPT on USMLE:
363 Potential for AI-assisted medical education using large language models. *PLOS*
364 *digital health*. 2023;2(2):e0000198.
- 365 8. Bommarito J, Bommarito M, Katz DM, et al. GPT as Knowledge Worker: A
366 Zero-Shot Evaluation of (AI) CPA Capabilities. 2023.
- 367 9. Sun H. Are ChatGPT's knowledge and interpretation ability comparable to those
368 of medical students in Korea for taking a parasitology examination?: a descriptive
369 study. *Journal of Educational Evaluation for Health Professions*. 2023;20:1.
- 370 10. Xiancheng W. Experiences, challenges, and prospects of National Medical
371 Licensing Examination in China. *BMC Medical Education*. 2022;22(1):349.
- 372 11. Almazyad M, Aljofan F, Abouammoh NA, et al. Enhancing Expert Panel
373 Discussions in Pediatric Palliative Care: Innovative Scenario Development and
374 Summarization With ChatGPT-4. *Cureus*. 2023;15(4).
- 375 12. Lai VD, Ngo NT, Veyseh APB, et al. ChatGPT Beyond English: Towards a
376 Comprehensive Evaluation of Large Language Models in Multilingual Learning.
377 arXiv preprint arXiv:230405613. 2023.
- 378 13. Peng K, Ding L, Zhong Q, et al. Towards making the most of chatgpt for machine
379 translation. arXiv preprint arXiv:230313780. 2023.
- 380 14. Liu X, Fang C, Wang J. Performance of ChatGPT on Clinical Medicine Entrance
381 Examination for Chinese Postgraduate in Chinese. *medRxiv*. 2023:2023.2004.
382 2012.23288452.
- 383 15. Butler S. GPT 3.5 vs GPT 4: What's Difference Available:
384 <https://www.howtogeek.com/882274/gpt-3-5-vs-gpt-4/>. Accessed MAR 31, 2023.
- 385 16. F. D-V. Considerations for the Practical Impact of AI in Healthcare Food and
386 Drug Administration. 2023.
- 387 17. Zweig M EBRH. How should the FDA approach the regulation of AI and
388 machine learning in healthcare?

389 Available:
390 <https://rockhealth.com/how-should-the-fda-approach-the-regulation-of-ai-and-machine-learning-in-healthcare/>.
391

392

393 **Figure legends**

394 **Fig 1. Schematic of workflow for sourcing, encoding, and adjudicating results.**

395 The 600 questions were categorized into 4 units. The accuracy of the open-ended
396 encoded questions was evaluated, while the answer with forced justification encoded
397 questions were also assessed for the accuracy, concordance. The adjudication process
398 was carried out by two physicians, and in case of any discrepancies in the domains, a
399 third physician was consulted for adjudication. Additionally, any inappropriate output
400 was identified and required re-encoding.

401 **Fig 2. Score of ChatGPT4 on Chinese National Medical Licensing Examination** 402 **before encoding.**

403 ChatGPT's outputs from Units 1, 2, 3, and 4 were scored for each unit.

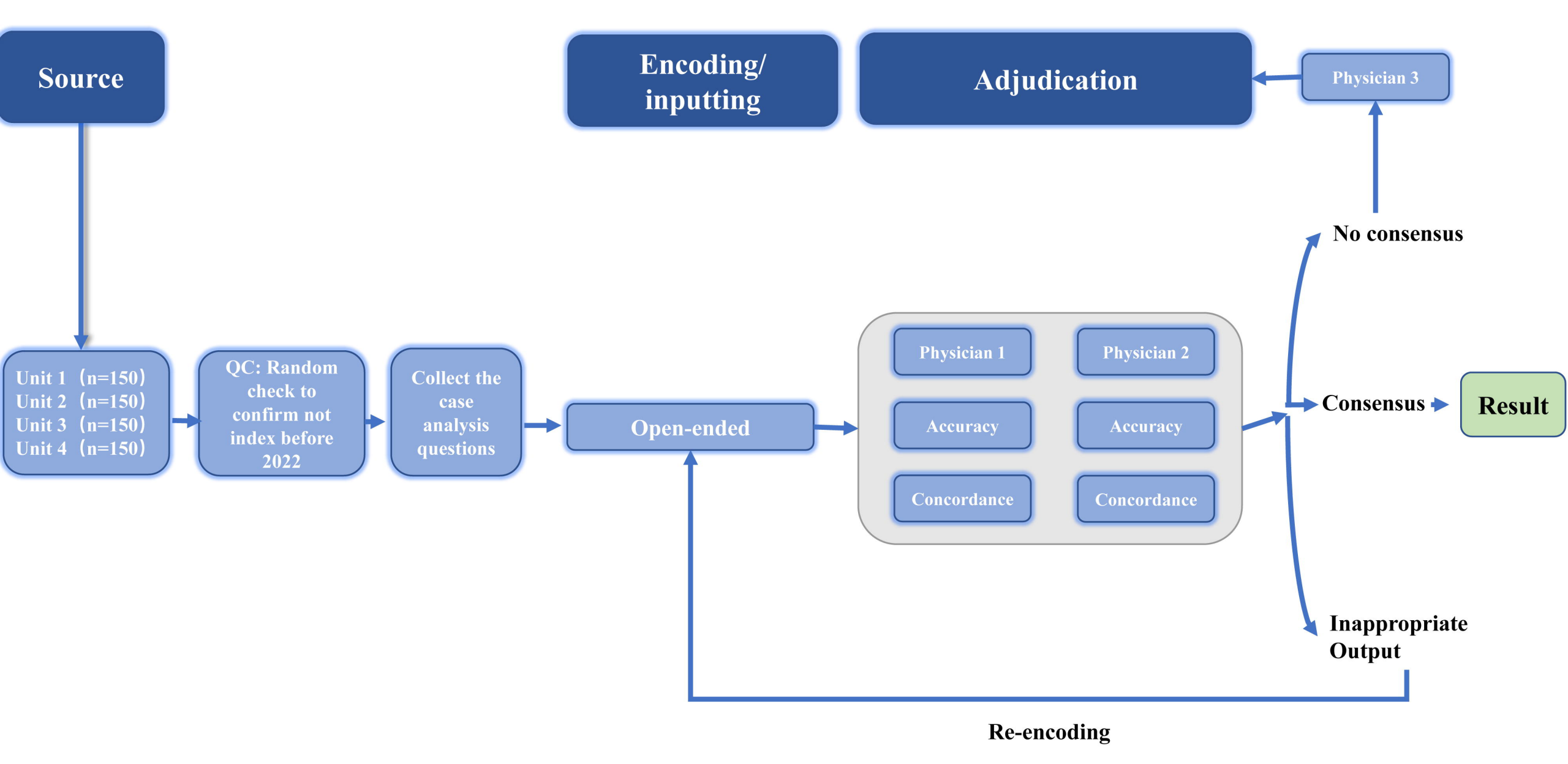
404 **Fig 3. Accuracy of ChatGPT4 on Chinese National Medical Licensing** 405 **Examination before encoding.**

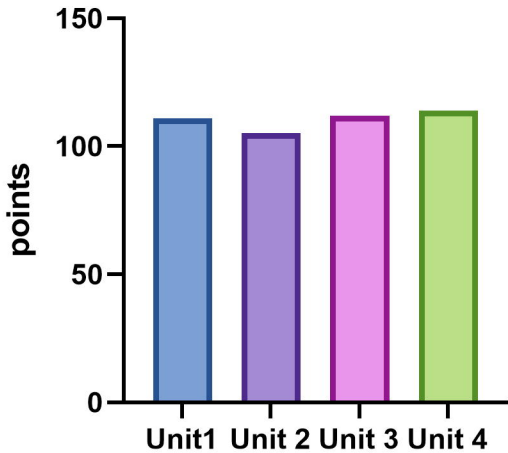
406 ChatGPT's outputs for Units 1, 2, 3, and 4 were evaluated as accurate, inaccurate, or
407 indeterminate using the scoring system outlined in S2 Data after encoding. (A)
408 Assessment of accuracy for open-ended question encodings. (B) Reduced accuracy of
409 encoded questions across Units 1, 2, 3, and 4.

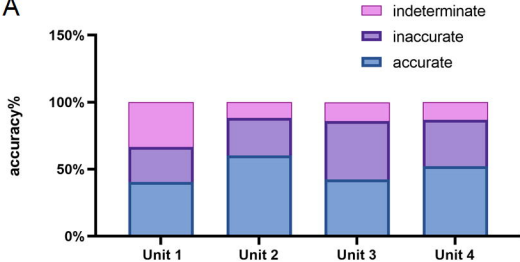
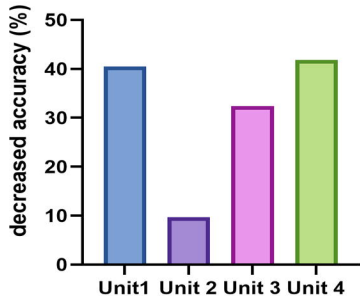
410 **Fig 4. Concordance of ChatGPT4 on Chinese National Medical Licensing**

411 **Examination after encoding.**

412 For Units 1, 2, 3, and 4 after encoding, ChatGPT's outputs were evaluated as
413 concordant or discordant, based on the scoring system detailed in S2 Data. This figure
414 demonstrates the concordance rates stratified between accurate, inaccurate, and
415 indeterminate outputs, across all case analysis questions.





A**B**

concordance
%

