

OPEN ACCESS

EDITED BY
Luke Andrew Woodham,
City St George's, University of London,
United Kingdom

REVIEWED BY Yuanda Zhu, Independent Researcher, Atlanta, GA, United States Tse-Yen Yang, China Medical University, Taiwan

*CORRESPONDENCE Ryunosuke Noda ☑ nodaryu00@gmail.com

RECEIVED 10 September 2025 ACCEPTED 30 October 2025 PUBLISHED 25 November 2025

CITATION

Noda R, Yuasa C, Kitano F, Ichikawa D and Shibagaki Y (2025) Performance of o1 pro and GPT-4 in Self-Assessment Questions for Nephrology Board Renewal. Front. Med. 12:1702668. doi: 10.3389/fmed.2025.1702668

COPYRIGHT

© 2025 Noda, Yuasa, Kitano, Ichikawa and Shibagaki. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Performance of o1 pro and GPT-4 in Self-Assessment Questions for Nephrology Board Renewal

Ryunosuke Noda*, Chiaki Yuasa, Fumiya Kitano, Daisuke Ichikawa and Yugo Shibagaki

Division of Nephrology and Hypertension, Department of Internal Medicine, St. Marianna University School of Medicine, Kawasaki, Japan

Background: Large language models (LLMs) are increasingly evaluated in medical education and clinical decision support, but their performance in highly specialized fields, such as nephrology, is not well established. We compared two advanced LLMs, GPT-4 and the newly released o1 pro, on comprehensive nephrology board renewal examinations.

Methods: We administered 209 Japanese Self-Assessment Questions for Nephrology Board Renewal from 2014 to 2023 to o1 pro and GPT-4 using ChatGPT pro. Each question, including images, was presented in separate chat sessions to prevent contextual carryover. Questions were classified by taxonomy (recall/interpretation/problem-solving), question type (general/clinical), image inclusion, and nephrology subspecialty. We calculated the proportion of correct answers and compared performances using chi-square or Fisher's exact tests.

Results: Overall, o1 pro scored 81.3% (170/209), significantly higher than GPT-4's 51.2% (107/209; p < 0.001). o1 pro exceeded the 60% passing criterion every year, while GPT-4 achieved this in only two out of the 10 years. Across taxonomy levels, question types, and the presence of images, o1 pro consistently outperformed GPT-4 (p < 0.05 for multiple comparisons). Performance differences were also significant in several nephrology subspecialties, such as chronic kidney disease, confirming o1 pro's broad superiority.

Conclusion: o1 pro significantly outperformed GPT-4 in a comprehensive nephrology board renewal examination, demonstrating advanced reasoning and integration of specialized knowledge. These findings highlight the potential of next-generation LLMs as valuable tools in nephrology, warranting further and careful validation.

KEYWORDS

large language models, ChatGPT, GPT-4, o1, o1 pro, nephrology

Introduction

Recent advances in artificial intelligence, particularly in large language models (LLMs), have dramatically improved their capabilities. By learning from vast amounts of data, LLMs have achieved an unprecedented level of language comprehension, including the ability to organize and summarize extensive information (1). Their accessibility—requiring no specialized programming skills and simulating the experience of interacting with an expert consultant—has made them increasingly appealing to clinicians. This progress has expanded

their potential use across healthcare, from medical education and clinical practice to research support (2–6). In nephrology, for example, LLMs are being investigated for dietary counseling in kidney disease (7), tailoring hemodialysis prescriptions (8), enhancing kidney transplant care (9), and supporting nephrology-specific literature retrieval (10).

Among LLMs, the Generative Pretrained Transformer (GPT) series from OpenAI has garnered particular attention (11). These models demonstrate strong adaptability to English and multilingual environments, excelling in various medical assessments—including national medical licensing examinations in Japan and the United States, and specialty exams (12–17). Notably, GPT-4, a flagship model in this series (18), reportedly achieved passing scores on the Polish nephrology board examinations (19), suggesting that LLMs might perform well even in highly specialized medical fields such as nephrology.

The field of LLMs continues to evolve. In December 2024, OpenAI released the o1 series—models developed through advanced reinforcement learning that exhibited enhanced reasoning abilities and fewer instances of hallucinations (factually ungrounded responses) (20). In medical fields, o1 achieved over 90% accuracy on the benchmark consisting of Japanese national medical licensing examinations (21). Its enhanced version, o1 pro, reportedly performed even better in domains requiring rigorous scientific reasoning in mathematics, science, and coding (22). The emergence of o1 pro marks the next generation of LLMs, potentially surpassing GPT-4, thereby heightening expectations for specialized applications.

Despite these advances, it remains unclear whether LLMs can handle the rigor, complexity, and contextual demands of nephrology—a field requiring integrated clinical reasoning, specialized knowledge, and the interpretation of complex data, including imaging and pathology. The Self-Assessment Questions for Nephrology Board Renewal (SAQ-NBR), provided by the Japanese Society of Nephrology, offer a comprehensive set of Japanese-language multiple-choice nephrology questions (23). These encompass fundamental concepts and complex clinical scenarios, including image-based challenges (kidney biopsies, imaging findings), and cover a wide range of nephrology subspecialties. Passing the SAQ-NBR requires integrative knowledge and clinical reasoning, making it a strict benchmark for advanced nephrology competence.

This study compared the performance of GPT-4 and the newly released o1 pro on the SAQ-NBR, clarifying whether LLMs could meet the cognitive and domain-specific challenges posed by complex nephrology content. Through this comparison, we aimed to provide insights into their potential applications in nephrology education and clinical support.

Materials and methods

The Self-Assessment Questions for Nephrology Board Renewal

The Self-Assessment Questions for Nephrology Board Renewal (SAQ-NBR) are a series of multiple-choice questions administered annually by the Japanese Society of Nephrology (23, 24). These questions, presented in Japanese, serve as a reference tool for nephrologists seeking board renewal. The passing criterion is defined

as achieving ≥60% correct answers. Each year's examination comprises a range of clinical and basic science questions that collectively span the breadth of nephrology. A subset of the questions includes images—such as renal biopsy specimens and radiological images—to assess interpretive and diagnostic skills. In this study, we included a total of 209 SAQ-NBR items from examinations between 2014 and 2023, excluding a single question that had been officially withdrawn as invalid by the Society.

Classification by taxonomy, question type, image inclusion, and subspecialty

Following a previous report (16), we classified each question into four categories:

Taxonomy: Based on the question creation manual for the Japanese National Medical Examination created by the Japan Medical Association (25), questions were categorized into three cognitive levels: recall, interpretation, and problem-solving. This taxonomy reflects the escalating depth of cognitive processing required to arrive at the correct answer.

Question Type: Questions were divided into general (focusing on fundamental medical or nephrological knowledge) and clinical (requiring clinical decision-making or patient management strategies).

Image Inclusion: Questions were designated as either image-based (image questions), incorporating visual data such as histopathological or radiological findings, or text-only (non-image questions).

Subspecialty: Drawing on the classification scheme of the Japanese Society of Nephrology's case experience list (26), questions were assigned to one of several nephrology subspecialty areas: chronic kidney disease/end-stage kidney disease (CKD/ESKD), acute kidney injury (AKI), glomerular diseases, tubulointerstitial diseases, hypertension/vascular diseases, water/electrolyte/acid-base disorders, autosomal dominant polycystic kidney disease (ADPKD)/urology, or basic medicine.

LLM models (o1 pro and GPT-4)

The o1 pro model, released in December 2024, represents the latest generation of LLMs, purportedly offering superior reasoning capabilities (20). GPT-4, introduced in March 2023, has demonstrated high performance on various medical examinations and has been widely recognized for its medical reasoning prowess (3, 18). We accessed both o1 pro and GPT-4 through ChatGPT pro interface. All prompts were input in December 2024.

To prevent context learning, each question was presented in a new chat session. Exceptionally, when consecutive questions pertained to the same clinical case, the chat session was not refreshed, and the same session was used to input subsequent questions. For each prompt, we stated in Japanese: "We will now present a nephrology-related question. Please provide your answer and explanation." The full text of the question was then provided. For image-based questions, the corresponding image was captured as a PNG file using the standard Windows screenshot tool and input simultaneously with the question text. The responses from both models were recorded, and their correctness was adjudicated based on the official answers provided by the Japanese Society of Nephrology.

Statistical analysis

For both o1 pro and GPT-4, the overall proportion of correct answers for 10 years and the proportion of correct answers by year were calculated, and whether they met the pass criteria (\geq 60% correct) for each year was evaluated. Additionally, we calculated and compared correct answer proportions by taxonomy, question type, image inclusion, and nephrology subspecialty. Statistical analyses were performed using Python version 3.10.12. Differences in proportions were evaluated using chi-square or Fisher's exact tests, as appropriate. A p-value <0.05 was considered statistically significant.

Results

Overall and annual performance

Across the 209 SAQ-NBR questions between 2014 and 2023, we confirmed the absence of duplicate content. The distribution of the number of questions and their classifications by year is detailed in Supplementary Table S1. Overall, o1 pro achieved a proportion of correct answers of 81.3% (170/209), significantly surpassing GPT-4's 51.2% (107/209; *p* < 0.001; Table 1). When examined by year, o1 pro consistently maintained high accuracy (70-95%) and exceeded the ≥60% threshold in every examination year. In contrast, GPT-4 showed considerable variability (35-72%) and passed the 60% threshold in only two of the 10 years. Notably, in the 2016, 2018, and 2020 examinations, o1 pro recorded accuracy of 95.0, 95.0, and 84.2%, respectively, significantly outperforming GPT-4's 55.0, 55.0, and 31.6% for the corresponding years (p = 0.011, 0.011, 0.003). To mitigate potential data leakage from public SAQ-NBR material, we performed a sensitivity analysis restricted to the most recent exam years (2022-2023). The results (o1 pro 72.5% [29/40] vs. GPT-4 47.5% [19/40];

TABLE 1 The proportion of correct answers of o1 pro and GPT-4 by exam year on the Self-Assessment Questions for Nephrology Board Renewal.

Exam year	The proportion of correct answers		<i>p</i> -value
	o1 pro	GPT-4	
2014	22/25 (88.0%)	17/25 (68.0%)	0.172
2015	22/25 (88.0%)	18/25 (72.0%)	0.289
2016	19/20 (95.0%)	11/20 (55.0%)	0.011
2017	14/20 (70.0%)	7/20 (35.0%)	0.057
2018	19/20 (95.0%)	11/20 (55.0%)	0.011
2019	15/20 (75.0%)	10/20 (50.0%)	0.191
2020	16/19 (84.2%)	6/19 (31.6%)	0.003
2021	14/20 (70.0%)	8/20 (40.0%)	0.112
2022	14/20 (70.0%)	11/20 (55.0%)	0.513
2023	15/20 (75.0%)	8/20 (40.0%)	0.055
Overall	170/209 (81.3%)	107/209 (51.2%)	< 0.001

Performance of o1 pro and GPT-4 on Self-Assessment Questions for Nephrology Board Renewal. Overall performance and exam year breakdown are reported. Differences in performance between large language models were queried using chi-squared and Fisher's exact tests.

p = 0.022) confirmed that the between-model difference remained statistically significant on the newest items.

Category-specific performance

By taxonomy, o1 pro significantly outperformed GPT-4 across all cognitive levels (Table 2). For recall-level questions, o1 pro achieved an 83.3% accuracy (90/108) compared to GPT-4's 49.1% (53/108; p < 0.001). For interpretation-level questions, the respective rates were 75.0% (42/56) versus 50.0% (28/56; p = 0.011). At the most complex level, problem-solving questions, o1 pro maintained an advantage with 84.4% (38/45) versus 57.8% (26/45; p = 0.011).

Regarding question type, o1 pro also outpaced GPT-4. For general questions, o1 pro's correct rate was 83.8% (93/111), exceeding GPT-4's 49.5% (55/111; p < 0.001). For clinical questions, o1 pro similarly excelled, with an accuracy of 78.6% (77/98) compared to GPT-4's 53.1% (52/98; p < 0.001).

Analysis by image inclusion showed that o1 pro demonstrated superior performance in both non-image and image-based questions. For non-image questions, o1 pro achieved 82.2% (139/169) and GPT-4 achieved 53.3% (90/169; p < 0.001). For image-based questions, o1 pro's correct rate was 77.5% (31/40) compared to GPT-4's 42.5% (17/40; p = 0.003).

Subspecialty analysis revealed that o1 pro significantly outperformed GPT-4 in several key areas, including CKD/ESKD (75.6% [34/45] vs. 42.2% [19/45]; p = 0.003), glomerular diseases, tubulointerstitial diseases, and basic medicine. In the remaining domains (AKI, hypertension/vascular diseases, and ADPKD/urology), o1 pro's accuracy exceeded that of GPT-4, although these differences did not reach statistical significance. The water/electrolyte/acid-base disorders subspecialty was the sole exception, showing identical model accuracy of 73.9% [17/23]. A sub-analysis confirmed this lack of difference persisted across all subcategories. This held true for 5 recall, 9 interpretation, and 9 problem-solving questions, as well as 6 general and 17 clinical questions; all comparisons yielded p = 1.000 (Supplementary Table S2). Concordance was high, as 19/23 items had identical outcomes: 15 were correct by both models and 4 were incorrect by both.

In sum, o1 pro consistently surpassed GPT-4 across virtually all examined domains—overall performance, annual pass rates, taxonomy levels, question types, presence or absence of images, and multiple nephrology subspecialties.

Discussion

This study showed that the new-generation LLM, o1 pro significantly outperformed GPT-4 on the comprehensive Japanese nephrology board renewal questions. Not only did o1 pro achieve passing scores in every examination year, but it also excelled in higher-order cognitive tasks such as clinical reasoning and interpretation of medical images. These results suggest that next-generation LLMs could extend beyond simple knowledge retrieval, integrating specialized medical knowledge, visual data processing, and context-specific decision-making into more advanced reasoning capabilities.

Previous studies have demonstrated the utility of LLMs in medical contexts, including GPT-4's strong performance on general medical

TABLE 2 The proportion of correct answers of o1 pro and GPT-4 by four categories on the Self-Assessment Questions for Nephrology Board Renewal

Category	The proportion of correct answers		<i>p</i> -value		
	o1 pro	GPT-4			
Taxonomy					
Recall	90/108 (83.3%)	53/108 (49.1%)	< 0.001		
Interpretation	42/56 (75.0%)	28/56 (50.0%)	0.011		
Problem-Solving	38/45 (84.4%)	26/45 (57.8%)	0.011		
Question type					
General Questions	93/111 (83.8%)	55/111 (49.5%)	< 0.001		
Clinical Questions	77/98 (78.6%)	52/98 (53.1%)	< 0.001		
Image inclusion					
Non-Image Questions	139/169 (82.2%)	90/169 (53.3%)	< 0.001		
Image Questions	31/40 (77.5%)	17/40 (42.5%)	0.003		
Subspecialty					
CKD/ESKD	34/45 (75.6%)	19/45 (42.2%)	0.003		
AKI	9/10 (90.0%)	4/10 (40.0%)	0.061		
Glomerular Diseases	47/57 (82.5%)	29/57 (50.9%)	< 0.001		
Tubulointerstitial Diseases	17/19 (89.5%)	9/19 (47.4%)	0.015		
Hypertension/Vascular Diseases	13/18 (72.2%)	8/18 (44.4%)	0.176		
Water/Electrolytes/ Acid-Base Disorder	17/23 (73.9%)	17/23 (73.9%)	1.000		
ADPKD/Urology	10/12 (83.3%)	7/12 (58.3%)	0.369		
Basic Medicine	23/25 (92.0%)	14/25 (56.0%)	0.010		

The performance of o1 pro and GPT-4 is reported for each category of Self-Assessment Questions for Nephrology Board Renewal. Differences in performance between large language models were queried using chi-squared and Fisher's exact tests. CKD, chronic kidney disease; ESKD, end-stage kidney disease; AKI, acute kidney injury; ADPKD, autosomal dominant polycystic kidney disease.

knowledge exams and national licensing examinations (12–15). In nephrology, GPT-4 was able to pass most of the Polish nephrology specialty exams (19). Additionally, our earlier study showed that GPT-4 significantly outperformed GPT-3.5 on the SAQ-NBR and met passing standards in several examination years (16). However, this study is the first to demonstrate that o1 pro consistently surpasses GPT-4 on this rigorous, specialized assessment. Beyond knowledge accuracy, o1 pro maintained a robust performance in complex interpretive and clinical questions, suggesting a major improvement in integrating kidney disease pathophysiology. These advancements could result from expanded training data, architectural upgrades, reinforcement learning for factual consistency, and enhanced multimodal processing (20, 21). Further model updates may improve LLMs' performance even more in nephrology.

A particularly notable finding was o1 pro's superior performance on image-based questions, including radiological and pathological images. Interpreting such visual information requires both deep medical knowledge and clinical experience—points not often emphasized in prior LLM studies in nephrology (27). Our results may indicate o1 pro's potential as a multifaceted clinical tool in nephrology, where expertise with textual, numerical, and image-based data is essential. From an educational standpoint, these capabilities may enable more effective image-focused teaching and objective skill assessment. Based on previous studies demonstrating advancements in LLMs' performance for analyzing radiological and pathological images (28–31), LLMs may have the potential to 1 day assist pathologists and nephrologists, improving the accuracy and consistency of medical imaging assessments.

These findings indicate that LLMs are moving closer to practical utility in nephrology education and clinical decision-making. Traditionally, clinicians have relied on textbooks, literature, lectures, and clinical training—yet the exponential growth of medical information makes it challenging to stay current (32). LLMs like o1 pro may serve as on-demand knowledge resources, offering evidence summaries, restructured pathophysiological concepts, and assistance with image interpretation. A study showed that LLMs could improve clinicians' exam performance in nephrology (33). Given o1 pro's high accuracy in both basic knowledge recall and complex problemsolving, it holds promise as a comprehensive educational and clinical support system for a broad range of users.

Despite o1 pro's demonstrated advantages in numerous categories, its lack of performance improvement in the water/electrolytes/acidbase disorder subspecialty is a notable exception. This is particularly noteworthy given o1 pro's design for enhanced inference. This subspecialty often requires the strict application of well-established physiological principles and codified diagnostic algorithms, and clinical guidelines (34–36). It was therefore considered an area requiring multi-step logical reasoning where o1 pro was expected to excel. However, our results confirmed that this lack of advantage persisted even in the complex problem-solving and interpretation subcategories. This suggests that the model's reasoning capabilities in general tasks may not necessarily align with its reasoning capabilities in the nephrology domain.

This study had several limitations. First, this study was limited to Japanese-language multiple-choice questions. Although the SAQ-NBR serves as a rigorous benchmark, these findings cannot be generalized to other languages or other assessment formats. It should be noted that performance on multiple-choice questions does not necessarily translate to competence in clinical reasoning or real-world clinical decision-making, which generally require free-response answers. Second, the internal reasoning processes of both o1 pro and GPT-4 were undisclosed (18, 20). This "black box" nature is a major limitation, as it complicates efforts to pinpoint the exact mechanisms behind their performance differences. As this study can only confirm a performance difference rather than the mechanism of that difference, establishing trust and ensuring safety for future applications in clinical decision-making remains a substantial challenge. Future validation requires a deeper understanding of how these models arrive at their conclusions. Third, the SAQ-NBR items were publicly available on the internet, so there is a possibility that o1 pro and GPT-4 might have encountered these questions during training. Given o1 pro's knowledge cutoff in October 2023 and GPT-4's in September 2021, data leakage may have occurred. To clarify, we have listed the publication dates of these questions and answers in Supplementary Table S3. Finally, this study has limitations concerning the reproducibility and consistency of the LLM outputs. Variations in LLM outputs may arise from different prompts ("prompt engineering"). Furthermore, the inherent stochastic nature of LLMs, influenced by factors such as probabilistic sampling during response generation and the use of a temperature parameter to control randomness, can lead to different outputs even with the same input. Our study did not assess performance consistency by testing each question multiple times. Therefore, it is possible that the results could vary upon repeated trials.

Future research should include cross-specialty and multimodal evaluations and large-scale analyses using varied exam formats, moving beyond multiple-choice questions to incorporate more clinically relevant formats such as free-response questions (37). Comparisons with other models (e.g., Gemini, Claude, Llama) could clarify the relative advantages of different LLMs and contribute to the broader medical AI ecosystem. Furthermore, before these models can be safely integrated into daily clinical practice, critical challenges related to patient data privacy and security must be addressed, such as ensuring Health Insurance Portability and Accountability Act (HIPAA)-compliant environments when sensitive information is processed by LLMs (38). Prospective clinical trials are needed to establish whether LLM integration in clinical workflows improves diagnostic and therapeutic outcomes and, most importantly, benefits patient care.

Conclusion

This study suggested that o1 pro consistently surpassed GPT-4 in tackling diverse nephrology tasks, highlighting the potential of next-generation LLMs as valuable tools in nephrology. However, their utility and safety in actual clinical decision-making remain unknown, and further validation by prospective clinical trials is required.

Data availability statement

The original contributions presented in the study are included in the article/Supplementary material, further inquiries can be directed to the corresponding author.

Ethics statement

This study did not involve human participants, human tissue samples, or any experiments on humans. Therefore, informed consent was not applicable. We consulted with the Division of Graduate Student Affairs and Research Promotion, which serves as the Institutional Review Board (IRB) office at St. Marianna University Hospital. After careful review, it was concluded that IRB approval was not indicated or required for this study because the study did not involve human subjects and is outside the scope of research requiring ethical review.

References

- 1. Zhao WX, Zhou K, Li J, Tang T, Wang X, Hou Y, et al. (2023). A survey of large language models. Available online at: https://arxiv.org/abs/2303.18223. (Accessed Mach 31, 2023)
- 2. Thirunavukarasu AJ, Ting DSJ, Elangovan K, Gutierrez L, Tan TF, Ting DSW. Large language models in medicine. *Nat Med.* (2023) 29:1930–40. doi: 10.1038/s41591-023-02448-8

Author contributions

RN: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Writing – original draft, Writing – review & editing. CY: Data curation, Investigation, Writing – review & editing. FK: Data curation, Writing – review & editing. DI: Supervision, Writing – review & editing. YS: Supervision, Writing – review & editing.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative Al statement

The authors declare that Gen AI was used in the creation of this manuscript. During the preparation of this manuscript, the authors used ChatGPT to improve language and readability. Following the use of this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

Any alternative text (alt text) provided alongside figures in this article has been generated by Frontiers with the support of artificial intelligence and reasonable efforts have been made to ensure accuracy, including review by the authors wherever possible. If you identify any issues, please contact us.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary material for this article can be found online at: https://www.frontiersin.org/articles/10.3389/fmed.2025.1702668/full#supplementary-material

- 3. Lee P, Bubeck S, Petro J. Benefits, limits, and risks of GPT-4 as an AI Chatbot for medicine. N Engl J Med. (2023) 388:1233–9. doi: 10.1056/NEJMsr2214184
- 4. Clusmann J, Kolbinger FR, Muti HS, Carrero ZI, Eckardt J-N, Laleh NG, et al. The future landscape of large language models in medicine. $Commun\ Med$. (2023) 3:1–8. doi: 10.1038/s43856-023-00370-1

- 5. Moulaei K, Yadegari A, Baharestani M, Farzanbakhsh S, Sabet B, Reza Afrash M. Generative artificial intelligence in healthcare: a scoping review on benefits, challenges and applications. *Int J Med Inform.* (2024) 188:105474. doi: 10.1016/j.ijmedinf.2024.105474
- 6. Sallam M. ChatGPT utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. *Healthcare*. (2023) 11:887. doi: 10.3390/healthcare11060887
- 7. Qarajeh A, Tangpanithandee S, Thongprayoon C, Suppadungsuk S, Krisanapan P, Aiumtrakul N, et al. AI-powered renal diet support: performance of ChatGPT, bard AI, and Bing chat. *Clin Pract.* (2023) 13:1160–72. doi: 10.3390/clinpract13050104
- 8. Hueso M, Álvarez R, Marí D, Ribas-Ripoll V, Lekadir K, Vellido A. Is generative artificial intelligence the next step toward a personalized hemodialysis? *Rev Investig Clin.* (2023) 75:309–17. doi: 10.24875/RIC.23000162
- 9. Garcia Valencia OA, Thongprayoon C, Jadlowiec CC, Mao SA, Miao J, Cheungpasitporn W. Enhancing kidney transplant care through the integration of Chatbot. *Healthcare*. (2023) 11:2518. doi: 10.3390/healthcare11182518
- 10. Aiumtrakul N, Thongprayoon C, Suppadungsuk S, Krisanapan P, Miao J, Qureshi F, et al. Navigating the landscape of personalized medicine: the relevance of ChatGPT, BingChat, and bard AI in nephrology literature searches. *J Pers Med.* (2023) 13:1457. doi: 10.3390/jpm13101457
- 11. ChatGPT (2022). Introducing ChatGPT: OpenAI. Available online at: https://openai.com/blog/chatgpt. (Accessed December 25 2024)
- 12. Kasai J, Kasai Y, Sakaguchi K, Yamada Y, Radev D. (2023) Evaluating GPT-4 and ChatGPT on Japanese medical licensing examinations. Available online at: https://arxiv.org/abs/2303.18027. (Accessed April 5 2024)
- 13. Takagi S, Watari T, Erabi A, Sakaguchi K. Performance of GPT-3.5 and GPT-4 on the Japanese medical licensing examination: comparison study. *JMIR Med Educ.* (2023) 9:e48002. doi: 10.2196/48002
- 14. Gilson A, Safranek CW, Huang T, Socrates V, Chi L, Taylor RA, et al. How does ChatGPT perform on the United States medical licensing examination? The implications of large language models for medical education and knowledge assessment. *JMIR Med Educ.* (2023) 9:e45312. doi: 10.2196/45312
- 15. Kung TH, Cheatham M, Medenilla A, Sillos C, De Leon L, Elepaño C, et al. Performance of ChatGPT on USMLE: potential for AI-assisted medical education using large language models. *PLOS Digit Health*. (2023) 2:e0000198. doi: 10.1371/journal.pdig.0000198
- 16. Noda R, Izaki Y, Kitano F, Komatsu J, Ichikawa D, Shibagaki Y. Performance of ChatGPT and bard in self-assessment questions for nephrology board renewal. *Clin Exp Nephrol.* (2024) 28:465–9. doi: 10.1007/s10157-023-02451-w
- 17. Miao J, Thongprayoon C, Garcia Valencia OA, Krisanapan P, Sheikh MS, Davis PW, et al. Performance of ChatGPT on nephrology test questions. *Clin J Am Soc Nephrol.* (2023) 19:35–43. doi: 10.2215/CJN.000000000000330
- 18. OpenAI. (2023) GPT-4 technical report. Available online at: https://arxiv.org/abs/2303.08774. (Accessed March 4, 2024)
- 19. Nicikowski J, Szczepański M, Miedziaszczyk M, Kudliński B. The potential of ChatGPT in medicine: an example analysis of nephrology specialty exams in Poland. *Clin Kidney J.* (2024) 17:sfae193. doi: 10.1093/ckj/sfae193
- 20. OpenAI. (2024) OpenAI o1 system card. Available online at: https://arxiv.org/abs/2412.16720. (Accessed December 21, 2024)

- 21. Nori H, Usuyama N, King N, McKinney SM, Fernandes X, Zhang S, et al. (2024) From Medprompt to o1: exploration of run-time strategies for medical challenge problems and beyond. Available online at: https://arxiv.org/abs/2411.03590. (Accessed November 6, 2024)
- $22.\ ChatGPT\ (2024).\ Introducing\ ChatGPT\ pro.\ Available\ online\ at: https://openai.\ com/index/introducing-chatgpt-pro/.\ (Accessed\ December\ 25,\ 2024)$
- 23. Japanese Society of Nephrology. (2024). Self-assessment questions for nephrology board renewal. Available online at: https://jsn.or.jp/medic/specialistsystem/question-unitupdate.php. (Accessed December 25, 2024).
- 24. Japanese Society of Nephrology. (2024). Overview of the JSN. Available online at: https://jsn.or.jp/en/about-jsn/overview-of-the-jsn/. (Accessed December 25, 2024).
 - 25. Uemura K. Exam preparation and taxonomy. Med Educ. (1982) 13:315-20.
- 26. Japanese Society of Nephrology. (2024). List of nephrologist experienced cases. Available online at: https://jsn.or.jp/education-specialist-committee/file-02_20210829. pdf. (Accessed December 25, 2024).
- 27. Koga S. Advancing large language models in nephrology: bridging the gap in image interpretation. *Clin Exp Nephrol.* (2024) 29:128–9. doi: 10.1007/s10157-024-02581-9
- 28. Zhang S, Xu Y, Usuyama N, Xu H, Bagga J, Tinn R, et al. A multimodal biomedical foundation model trained from fifteen million image–text pairs. *NEJM AI*. (2025) 2:AIoa2400640. doi: 10.1056/AIoa2400640
- 29. Tanno R, Barrett DGT, Sellergren A, Ghaisas S, Dathathri S, See A, et al. Collaboration between clinicians and vision–language models in radiology report generation. *Nat Med.* (2024) 31:599–608. doi: 10.1038/s41591-024-03302-1
- 30. Ferber D, Wölflein G, Wiest IC, Ligero M, Sainath S, Ghaffari Laleh N, et al. Incontext learning enables multimodal large language models to classify cancer pathology images. *Nat Commun.* (2024) 15:10104. doi: 10.1038/s41467-024-51465-9
- 31. Lu MY, Chen B, Williamson DFK, Chen RJ, Liang I, Ding T, et al. A visual-language foundation model for computational pathology. *Nat Med.* (2024) 30:863–74. doi: 10.1038/s41591-024-02856-4
- 32. Densen P. Challenges and opportunities facing medical education. *Trans Am Clin Climatol Assoc.* (2011) 122:48–58. Available online at: https://pmc.ncbi.nlm.nih.gov/articles/PMC3116346/
- 33. Noda R, Tanabe K, Ichikawa D, Shibagaki Y. (2024) Chatgpt's performance in supporting physician decision-making in nephrology multiple-choice questions. Available online at: https://www.researchsquare.com/article/rs-4947755/v1. (Accessed May 1, 2025)
- 34. Verbalis JG, Goldsmith SR, Greenberg A, Korzelius C, Schrier RW, Sterns RH, et al. Diagnosis, evaluation, and treatment of hyponatremia: expert panel recommendations. *Am J Med.* (2013) 126:S1–S42. doi: 10.1016/j.amjmed.2013.07.006
- 35. Spasovski G, Vanholder R, Allolio B, Annane D, Ball S, Bichet D, et al. Clinical practice guideline on diagnosis and treatment of hyponatraemia. *Nephrol Dial Transplant*. (2014) 29:i1–i39. doi: 10.1530/EJE-13-1020
- 36. Adrogué HJ, Gennari FJ, Galla JH, Madias NE. Assessing acid-base disorders. *Kidney Int.* (2009) 76:1239–47. doi: 10.1038/ki.2009.359
- 37. Johri S, Jeong J, Tran BA, Schlessinger DI, Wongvibulsin S, Barnes LA, et al. An evaluation framework for clinical use of large language models in patient interaction tasks. *Nat Med.* (2025) 31:77–86. doi: 10.1038/s41591-024-03328-5
- 38. Marks M, Haupt CE. AI Chatbots, health privacy, and challenges to HIPAA compliance. *JAMA*. (2023) 330:309–10. doi: 10.1001/jama.2023.9458