



## ChatGPT Solving Complex Kidney Transplant Cases: A Comparative Study with Human Respondents

Journal:	<i>Clinical Transplantation</i>
Manuscript ID	CLTX-24-BCOM-0499
Wiley - Manuscript type:	Brief Communication
Date Submitted by the Author:	15-Jun-2024
Complete List of Authors:	Mankowski, Michal; New York University Department of Surgery Jaffe, Ian; New York University Department of Surgery Xu, Jingzhi; New York University Department of Surgery Bae, Sunjae ; New York University Department of Surgery; New York University Department of Population Health Oermann, Eric; New York University Department of Neurosurgery Aphinyanaphongs, Yindalon; New York University Department of Population Health; New York University Department of Medicine McAdams Demarco, Mara; New York University Department of Surgery; New York University Department of Population Health Lonze, Bonnie ; New York University Department of Surgery Orandi, Babak J.; New York University Department of Surgery; New York University Department of Medicine Stewart, Darren; New York University Department of Surgery Levan, Macey; New York University Department of Surgery; New York University Department of Population Health Massie, Allan B.; New York University Department of Surgery; New York University Department of Population Health Gentry, Sommer E.; New York University Department of Surgery; New York University Department of Population Health Segev, Dorry L.; New York University Department of Surgery; New York University Department of Population Health
Transplant Peer Review Network - Second Choice:	No referral
Transplant Peer Review Network - First Choice:	No referral
Discipline:	organ transplantation in general, kidney transplantation/nephrology
Keywords:	kidney (allograft) function / dysfunction, kidney failure / injury
Abstract:	<b>Introduction:</b> ChatGPT has shown the ability to answer clinical questions in general medicine but may be constrained by the specialized nature of kidney transplantation. Thus, it is important to explore how ChatGPT can be used in kidney transplantation and how its knowledge compares to human respondents. <b>Methods:</b> We prompted ChatGPT versions 3.5, 4, and 4 Visual (4V) with 12 multiple-choice questions related to six kidney transplant cases from the 2013-2015 American Society of Nephrology (ASN) fellowship program quizzes. We compared

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

	<p>the performance of ChatGPT with US nephrology fellowship program directors, nephrology fellows, and the audience of the ASN’s annual Kidney Week meeting. <b>Results:</b> Overall, ChatGPT 4V correctly answered 10 out of 12 questions, showing a performance level comparable to nephrology fellows (group majority correctly answered 9 of 12 questions) and training program directors (11 of 12). This surpassed ChatGPT 4 (7 of 12 correct) and 3.5 (5 of 12). All 3 ChatGPT versions failed to correctly answer questions where the consensus among human respondents was low. <b>Conclusion:</b> Each iterative version of ChatGPT performed better than the prior version, with version 4V achieving performance on par with nephrology fellows and training program directors. While it shows promise in understanding and answering kidney transplantation questions, Chat GPT should be seen as a complementary tool to human expertise rather than a replacement.</p>

**Brief Communication:****ChatGPT Solving Complex Kidney Transplant Cases: A Comparative Study with Human Respondents**

Michal A. Mankowski, PhD<sup>1</sup>, Ian S. Jaffe, BS<sup>1</sup>, Jingzhi Xu, PhD<sup>1</sup>, Sunjae Bae, MD, PhD,<sup>1,2</sup> Eric K. Oermann, MD,<sup>3</sup> Yindalon Aphinyanaphongs, MD, PhD,<sup>2,4</sup> Mara A. McAdams-DeMarco, PhD<sup>1,2</sup> Bonnie E. Lonze, MD, PhD,<sup>1</sup> Babak J. Orandi, MD, PhD<sup>1,4</sup>, Darren Stewart, MS<sup>1</sup>, Macey Levan, JD, PhD<sup>1,2</sup>, Allan Massie, PhD<sup>1,2</sup>, Sommer Gentry, PhD<sup>1,2</sup>, Dorry L. Segev, MD, PhD<sup>1,2</sup>

<sup>1</sup> Department of Surgery, NYU Grossman School of Medicine, New York, NY

<sup>2</sup> Department of Population Health, NYU Grossman School of Medicine, New York, NY

<sup>3</sup> Department of Neurosurgery, NYU Grossman School of Medicine, New York, NY

<sup>4</sup> Department of Medicine, NYU Grossman School of Medicine, New York, NY

**Running Title:** ChatGPT Kidney Transplant Cases

**Key Words:** ChatGPT; kidney transplantation; artificial intelligence; generative pre-trained transformer; quiz

**ORCID/X Handles:**

Michal Mankowski - @MichalMankowski

Ian S. Jaffe - 0000-0002-7309-308X

Jingzhi Xu -

Sunjae Bae - 0000-0003-0098-8816

Eric K. Oermann - 0000-0002-1876-5963

Yindalon Aphinyanaphongs - 0000-0001-8605-5392

Mara A. McAdams-DeMarco - 0000-0003-3013-925X

Bonnie E. Lonze - 0000-0002-0973-1657

Babak Orandi - 0000-0001-6026-7135

Darren Stewart - 0000-0002-6764-4842

Macey Levan - 0000-0002-4239-1252 / @DrMaceyLevan

Allan Massie - 0000-0002-5288-5125 / @AllanBMassie

Sommer Gentry - 0000-0003-4530-8917 / @shelikesmath

Dorry L. Segev - 0000-0002-1924-4801 / @Dorry\_Segev

**Corresponding author:**

Michal Mankowski

Center for Surgical and Transplant Applied Research

Department of Surgery

New York University Grossman School of Medicine

One Park Ave,

6<sup>th</sup> floor, Room 6-664

New York City, NY 10016

michal.mankowski@nyulangone.org

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Word count:** 1931/2000 (abstract 204/250)

For Review Only

M Mankowski, IS Jaffe, J Xu, S Bae, EK Oermann, Y Aphinyanaphongs, MA McAdams-DeMarco, BE Lonze, BJ Orandi, D Stewart, M Levan, A Massie, S Gentry, DL Segev

## ChatGPT Solving Complex Kidney Transplant Cases: A Comparative Study with Human Respondents

*Clin. Transpl.*

### Abbreviations

AI, artificial intelligence

GPT, generative pre-trained transformer

LLM, large language model

### Abstract

**Introduction:** ChatGPT has shown the ability to answer clinical questions in general medicine but may be constrained by the specialized nature of kidney transplantation. Thus, it is important to explore how ChatGPT can be used in kidney transplantation and how its knowledge compares to human respondents.

**Methods:** We prompted ChatGPT versions 3.5, 4, and 4 Visual (4V) with 12 multiple-choice questions related to six kidney transplant cases from the 2013-2015 American Society of Nephrology (ASN) fellowship program quizzes. We compared the performance of ChatGPT with US nephrology fellowship program directors, nephrology fellows, and the audience of the ASN's annual Kidney Week meeting.

**Results:** Overall, ChatGPT 4V correctly answered 10 out of 12 questions, showing a performance level comparable to nephrology fellows (group majority correctly answered 9 of 12 questions) and training program directors (11 of 12). This surpassed ChatGPT 4 (7 of 12 correct) and 3.5 (5 of 12). All 3 ChatGPT versions failed to correctly answer questions where the consensus among human respondents was low.

**Conclusion:** Each iterative version of ChatGPT performed better than the prior version, with version 4V achieving performance on par with nephrology fellows and training program directors. While it shows promise in understanding and answering kidney transplantation questions, Chat GPT should be seen as a complementary tool to human expertise rather than a replacement.

**Key Words:** ChatGPT; kidney transplantation; artificial intelligence; generative pre-trained transformer; quiz

### Off-Print Requests:

Michal Mankowski  
One Park Ave,  
6<sup>th</sup> floor, Room 6-664

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

New York City, NY 10016  
michal.mankowski@nyulangone.org

For Review Only

## 1. Introduction

ChatGPT has been increasingly utilized in medicine due to its effectiveness in understanding and generating human-like text. The platform is built on the GPT<sup>1</sup> (Generative Pre-trained Transformer) architecture, utilizing a transformer-based<sup>2</sup> model that combines self-attention mechanisms with feedforward layers, enabling it to capture linguistic patterns and contexts. Its versatility has been demonstrated in various settings, including answering medical board/license examination questions<sup>3-5</sup> and addressing patients' queries related to medications, treatment plans, and post-treatment care.<sup>6,7</sup> The use of generative AI tools like ChatGPT in medicine continues to grow, with its potential to improve patient communication,<sup>8,9</sup> medical education,<sup>10-12</sup> and clinical decision-making.<sup>13-15</sup> Moreover, newer versions of ChatGPT have extended its applications by incorporating interpretation of visual and audio materials.<sup>16</sup>

However, the generalizability of ChatGPT's knowledge about general medicine to kidney transplantation may be limited due to the unique complexities of transplant medicine. AI systems like ChatGPT may lack the specificity needed for transplant-related issues like immunosuppression management and graft rejection. The dynamic nature of transplant protocols further challenges the relevance of AI tools in this field. Therefore, this study aims to explore how ChatGPT can be specifically applied to clinical interpretation in kidney transplantation using complex clinical cases.

In this study, we evaluated the performance of ChatGPT in answering multiple-choice questions from the transplant section of the American Society of Nephrology’s annual fellows training quiz. We compared the answers of ChatGPT with answers provided by training program directors of US nephrology fellowship programs, nephrology fellows, and the audience of the annual Kidney Week meeting of the American Society of Nephrology. This study is the first to compare the performance of ChatGPT in solving kidney transplant cases with human respondents.

**2. Methods**

**2.1 American Society of Nephrology Quiz and Questionnaire**

The Nephrology Quiz and Questionnaire is an educational session at the annual Kidney Week meeting of the American Society of Nephrology. A panel of experts prepares various clinical cases, including kidney transplant cases, each with two multiple-choice, single best-answer questions. Before the live session, United States nephrology training program directors and nephrology fellows independently answer the questions through an internet-based questionnaire. During the session, audience members (which may include any attendee of Kidney Week, ranging from individuals with no clinical training to expert nephrologists) compare their knowledge and judgment on case-oriented questions with the training program directors and fellows.

**2.2 Cases**

In our study, we analyzed kidney transplantation cases from the American Society of Nephrology Quiz and Questionnaire from 2013,<sup>17</sup> 2014,<sup>18</sup> and 2015.<sup>19</sup> The quizzes



consisted of six cases (two per year) and 12 single-answer multiple-choice questions (two questions per case) in total. Each case has been summarized by ChatGPT 4, and these summaries alongside the related questions are presented in Table 1. Although small, this highly tailored dataset of high-quality questions is aligned with similar efforts to identify smaller, quality datasets to investigate specific aspects of LLM performance.<sup>20,21</sup>

### 2.3. Quiz Answers: ChatGPT versus Human Respondents

We used the NYU Langone Health's instance on Azure OpenAI Studio to access three ChatGPT versions: ChatGPT 3.5, ChatGPT 4, and ChatGPT 4 Vision (4V) and used default settings. Each version was prompted 10 times with each case description (including lab result tables and figures captions) along with the two associated multiple-choice questions; repeated prompting was used to mitigate random response selection given the moderate temperature settings in default ChatGPT model standard settings.<sup>22</sup> The last line of the prompt requested that ChatGPT select an answer to each question based on the case. ChatGPT 4V was also given access to the figures for each case (when applicable). As image interpretation is not a component of ChatGPT 3.5 or 4, we relied on the written image interpretations incorporated into the written case descriptions, which were present in all cases with images.

The responses of ChatGPT 3.5, ChatGPT 4, and ChatGPT 4V were compared with the responses of the training program directors of US nephrology fellowship programs, nephrology fellows, and the audience of the Annual Kidney Week Meeting

of the American Society of Nephrology. The number of respondents, demographic information, and training details for the human groups was not reported in their initial publication, and thus are unavailable. We evaluated whether a given group (training program directors, fellows, audience, or ChatGPT version) answered a question correctly if the most frequently chosen answer by a given group was correct. In case two or more answers were most frequently chosen and had the same frequency, we deemed these answers incorrect. We also separately evaluated the reproducibility of correct responses from the ChatGPT models; responses were considered reproducibly correct when the ChatGPT model answered the question correctly in all 10 prompting sessions.

**3. Results**

Figure 1 summarizes the performance of each responder group in answering quiz questions. Across the 6 cases and 12 related questions, training program directors answered 11 of 12 questions correctly. ChatGPT 4V answered 10 of 12 questions correctly and fellows answered 9 of 12 questions correctly. Both the Kidney Week audience and ChatGPT 4 answered 7 of 12 questions correctly. ChatGPT 3.5 provided 5 of 12 correct answers.

In 2015, ChatGPT 4, ChatGPT 4V, the nephrology fellows, and the training program directors correctly answered both questions of case 1 (Q1A and Q1B), while

ChatGPT 3.5 failed to answer both. None of the groups answered Q2A correctly, while only the training program directors, the Kidney Week audience, and ChatGPT 4V answered Q2B for case 2 correctly.

In 2014, ChatGPT 3.5, ChatGPT 4V, and the training program directors answered all four questions for both cases correctly. The nephrology fellows and ChatGPT 4 each failed to answer one question, Q1B and Q2A, respectively, while the Kidney Week audience answered only two questions correctly across the two cases.

In 2013, only the training program directors and fellows had answered all questions correctly. ChatGPT 4V missed just Q1A, while ChatGPT 4 missed both Q1A and Q2B and ChatGPT 3.5 only answered Q2A correctly. The audience missed only one question (Q1A).

Table S1 shows the distribution of answers to multi-choice questions by the responders' group. Responses for ChatGPT models were considered reproducibly correct when the ChatGPT model not only answered the question correctly, but did so 100% of the time. The higher correct response rates seen for ChatGPT 4V (10 of 12 questions) and Chat GPT 4 (7 of 12) compared to ChatGPT 3.5 (5 of 12), corresponded to higher reproducibly correct response rates; ChatGPT 4V was reproducibly correct on 9 of 12 questions, ChatGPT 4 on 6 of 12 and ChatGPT 3.5 on 4 of 12.

For questions where all three ChatGPT versions failed, there was less consensus among human respondents. For instance, in the 2013 quiz Q1A, only 40% of fellows, 69% of training program directors, and 27% of the audience provided correct answers, while ChatGPT versions 3.5 and 4 did not provide correct answers in any of the 10 prompt sessions and ChatGPT 4V provided the correct answer in only a single session. Similarly, in the 2015 quiz Q2A, no group most frequently selected the correct response.

**4. Discussion**

The utilization of Large Language Models (LLMs), such as ChatGPT, within medicine<sup>23,24</sup> has introduced a novel approach to assimilating medical knowledge. Our findings indicate that various versions of ChatGPT successfully answered an array of case-based kidney transplant questions from the 2013-2015 American Society of Nephrology Quiz and Questionnaire. Interestingly, in the 2015 quiz, ChatGPT 4V answered 3 of the 4 questions correctly, matching the performance of nephrology training program directors. In the 2014 quiz, both ChatGPT 3.5 and 4V answered all questions correctly, paralleling the performance of the training program directors. Yet, both ChatGPT 3.5 and 4 underperformed in the 2013 quiz, answering fewer questions correctly than all human groups, while ChatGPT 4V performed on par with the Kidney Week audience, but fell short when compared to nephrology experts.

While earlier versions of ChatGPT were unable to outperform the human groups, specifically the nephrology fellow and training program director groups, ChatGPT 4V

1  
2  
3 demonstrated comparable performance to these knowledgeable groups. This suggests  
4  
5 an improvement in the capabilities of the GPT models to handle multiple-choice  
6  
7 reasoning tasks in highly specialized areas like transplant nephrology. They are  
8  
9 consistent with a recent study conducted on general nephrology test questions, where  
10  
11 ChatGPT 4 achieved an accuracy of 74%, slightly below the average 77% accuracy of  
12  
13 nephrology examinees.<sup>25</sup> Our study indicates a notable overall performance  
14  
15 improvement in ChatGPT 4V compared to ChatGPT 4, and in ChatGPT 4 compared  
16  
17 to ChatGPT 3.5, suggesting that LLM model performance may enhance with  
18  
19 subsequent generations and benefit from multimodal inputs such as images.  
20  
21  
22 Considering that these models were iteratively improved over the span of less than  
23  
24 two years, this improvement reflects significant change in a short period of time. The  
25  
26 swift progress in the AI field implies that these models may eventually exceed human  
27  
28 experts, especially in multiple-choice constrained reasoning tasks.  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Furthermore, there is burgeoning evidence that LLMs can offer beneficial responses to non-constrained questions specifically related to transplantation. For instance, ChatGPT responses to questions related to the treatment of kidney transplant recipients showed considerable knowledge of kidney transplantation, albeit with some inaccuracies and a lack of depth.<sup>26</sup> Regarding liver disease, ChatGPT has been evaluated in answering general questions on liver transplantation, where 70.6% of clinical experts found it accurate and comparable to practicing clinicians.<sup>27</sup> ChatGPT was also evaluated in generating research ideas in liver transplantation.<sup>28</sup>

There are several reasons why AI models like ChatGPT may provide incorrect answers when compared to humans. The medical knowledge of ChatGPT is derived from its training on a diverse range of internet text. Its performance is not a result of its understanding of medicine but its ability to generate plausible-sounding text based on the patterns it learned during training.<sup>29</sup> ChatGPT models do not have a conceptual understanding of the topics they are trained on. Instead, they generate responses based on patterns in the data they were trained on, generating text in a probabilistic manner. ChatGPTs tend to hallucinate,<sup>30,31</sup> and make up fictitious information. This lack of understanding can lead to mistakes or oversights that a human expert possibly would not make. Additionally, current ChatGPT models do not have the capability to autonomously gain new knowledge or learn from experience. Their knowledge is fixed at the time of training and is not being updated or expanded upon until the next software release.

It is important to consider the limitations of our study as well. We cannot exclude the possibility that the quiz questions and answers were themselves part of the training data for either or both versions of ChatGPT since these quizzes were published in 2013-2015. This would tend to bias our conclusions toward higher performance of ChatGPT, however we observed that none of these ChatGPT versions gave entirely correct answers, suggesting the quizzes themselves may not have been part of the training corpus. Another limitation is that the assessment of ChatGPT's performance was based solely on a set of questions from the American Society of Nephrology Quiz cases, which do not represent the full spectrum of clinical scenarios in kidney

transplantation. The study's format may have introduced a limitation as it required all case information to be converted into text format to be processed by both earlier ChatGPTs (3.5 and 4) in a comparable way. This constraint could have resulted in a potential loss of context or nuance in the presented medical cases. While this was mitigated by the fact that the images were also described in the case description and in the figure captions, this could provide an explanation for some of the improvement in performance observed in ChatGPT 4V.

In conclusion, while LLM models like ChatGPT have stirred enthusiasm for potential medical applications, including kidney transplantation, it is crucial to comprehend their limitations and utilize them as a complementary tool to human expertise rather than a replacement. Our study evaluated the performance of ChatGPT in solving clinical kidney transplant scenarios and compared the results with human respondents of varying levels of expertise. The advanced ChatGPT 4V model performed comparably to nephrology fellows and nephrology training program directors on this small sample of questions, correctly answering ten out of twelve questions related to kidney transplantation cases. Such performance supports further exploration in using LLMs to assist with clinical interpretation, answering patient questions, and other tasks specific to organ transplant.

**Acknowledgements and Funding:**

**Author Contributions:** Concept/design: MM, JX, SG, DL; Data collection: MM, JX; Data analysis/interpretation: all authors; Drafting article: all authors; Critical revision of article: all authors; Approval of article: all authors; Funding secured by: DLS.

**Funding:** This work was supported by grant number K24AI144954 (Segev) from the National Allergy and Infectious Disease (NIAID) and R01DK132395 (Massie) from the National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK). I.S. Jaffe was supported by the National Center for Advancing Translational Sciences (NCATS), National Institutes of Health, through grant award number UL1TR001445. M. McAdams-DeMarco was supported by K02AG076883 and R01AG077888 from the National Institute on Aging (NIA). The work described here is the responsibility of the authors alone and does not necessarily reflect the views or policies of the Department of Health and Human Services, nor does mention of trade names, commercial products, or organizations imply endorsement by the U.S. Government.

**Conflict of Interest Disclosures:**

M Mankowski: None, IS Jaffe: None, J Xu: None, None, S Bae: None, EK Oermann: None, Y Aphinyanaphongs: None, MA McAdams DeMarco: Dr. McAdams-DeMarco reports speaker honoraria from Chiesi Farmaceutici S.p.A. unrelated to this work, B Lonze: None, B Orandi: None, D Stewart: None, M Levan: None, A Massie: None, S Gentry: None, DL Segev: None



## References

1. Radford A, Narasimhan K, Salimans T, Sutskever I. Improving language understanding by generative pre-training. Published online 2018.
2. Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. *Advances in neural information processing systems*. 2017;30.
3. 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 Medical Education*. 2023;9(1):e48002. doi:10.2196/48002
4. Joly-Chevrier M, Nguyen AXL, Lesko-Krleza M, Lefrançois P. Performance of ChatGPT on a practice dermatology board certification examination. *Journal of cutaneous medicine and surgery*. 2023;27(4):407-409.
5. Beam K, Sharma P, Kumar B, et al. Performance of a large language model on practice questions for the neonatal board examination. *JAMA pediatrics*. 2023;177(9):977-979.
6. Harris E. Large language models answer medical questions accurately, but can't match clinicians' knowledge. *JAMA*. Published online 2023.
7. Haupt CE, Marks M. AI-generated medical advice—GPT and beyond. *Jama*. 2023;329(16):1349-1350.
8. Gordon EB, Towbin AJ, Wingrove P, et al. Enhancing patient communication with ChatGPT in radiology: evaluating the efficacy and readability of answers to common imaging-related questions. *Journal of the American College of Radiology*. 2024;21(2):353-359.
9. Nashwan AJ, Abujaber AA, Choudry H. Embracing the future of physician-patient communication: GPT-4 in gastroenterology. *Gastroenterology & Endoscopy*. 2023;1(3):132-135.
10. Lee H. The rise of ChatGPT: Exploring its potential in medical education. *Anatomical sciences education*. Published online 2023.
11. Khan RA, Jawaid M, Khan AR, Sajjad M. ChatGPT-Reshaping medical education and clinical management. *Pakistan Journal of Medical Sciences*. 2023;39(2):605.
12. Hswen Y, Abbasi J. AI will—and should—change medical school, says Harvard's dean for medical education. *JAMA*. Published online 2023.
13. Liu J, Wang C, Liu S. Utility of ChatGPT in clinical practice. *Journal of Medical Internet Research*. 2023;25:e48568.
14. Sorin V, Klang E, Sklair-Levy M, et al. Large language model (ChatGPT) as a support tool for breast tumor board. *NPJ Breast Cancer*. 2023;9(1):44.

15. Rao A, Kim J, Kamineni M, et al. Evaluating GPT as an adjunct for radiologic decision making: GPT-4 versus GPT-3.5 in a breast imaging pilot. *Journal of the American College of Radiology*. 2023;20(10):990-997.

16. Liu H, Li C, Wu Q, Lee YJ. Visual instruction tuning. *Advances in neural information processing systems*. 2024;36.

17. Josephson MA, Perazella MA, Choi MJ. American society of nephrology quiz and questionnaire 2013: Transplantation. *Clinical Journal of the American Society of Nephrology*. 2014;9(7):1319-1327.

18. Josephson MA, Perazella MA, Choi MJ. American society of Nephrology Quiz and Questionnaire 2014: transplantation. *Clinical Journal of the American Society of Nephrology*. 2015;10(5):903-909.

19. Josephson MA, Perazella MA, Choi MJ. American society of nephrology quiz and questionnaire 2015: Transplantation. *Clinical Journal of the American Society of Nephrology*. 2016;11(6):1114-1122.

20. Jimenez CE, Yang J, Wettig A, et al. SWE-bench: Can Language Models Resolve Real-World GitHub Issues? Published online April 5, 2024. doi:10.48550/arXiv.2310.06770

21. Rein D, Hou BL, Stickland AC, et al. GPQA: A Graduate-Level Google-Proof Q&A Benchmark. Published online November 20, 2023. doi:10.48550/arXiv.2311.12022

22. Davis J, Van Bulck L, Durieux BN, Lindvall C. The temperature feature of ChatGPT: Modifying creativity for clinical research. *JMIR Human Factors*. 2024;11(1):e53559.

23. Lee P, Bubeck S, Petro J. Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine. Drazen JM, Kohane IS, Leong TY, eds. *N Engl J Med*. 2023;388(13):1233-1239. doi:10.1056/NEJMSr2214184

24. Shah NH, Entwistle D, Pfeffer MA. Creation and adoption of large language models in medicine. *Jama*. 2023;330(9):866-869.

25. Miao J, Thongprayoon C, Valencia OAG, et al. Performance of ChatGPT on nephrology test questions. *Clinical Journal of the American Society of Nephrology*. Published online 2023;10-2215.

26. Rawashdeh B, Kim J, AlRyalat SA, Prasad R, Cooper M. ChatGPT and artificial intelligence in transplantation research: is it always correct? *Cureus*. 2023;15(7).

27. Endo Y, Sasaki K, Moazzam Z, et al. Quality of ChatGPT responses to questions related to liver transplantation. *Journal of Gastrointestinal Surgery*. 2023;27(8):1716-1719.

28. Akabane M, Iwadoh K, Melcher ML, Sasaki K. Exploring the potential of ChatGPT in generating unknown clinical questions about liver transplantation: A feasibility study. *Liver Transplantation*. 2024;30(2):229-234.

29. Bender EM, Gebru T, McMillan-Major A, Shmitchell S. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? □. In: *Proceedings of the 2021 ACM*

*Conference on Fairness, Accountability, and Transparency*. FAccT '21. Association for Computing Machinery; 2021:610-623. doi:10.1145/3442188.3445922

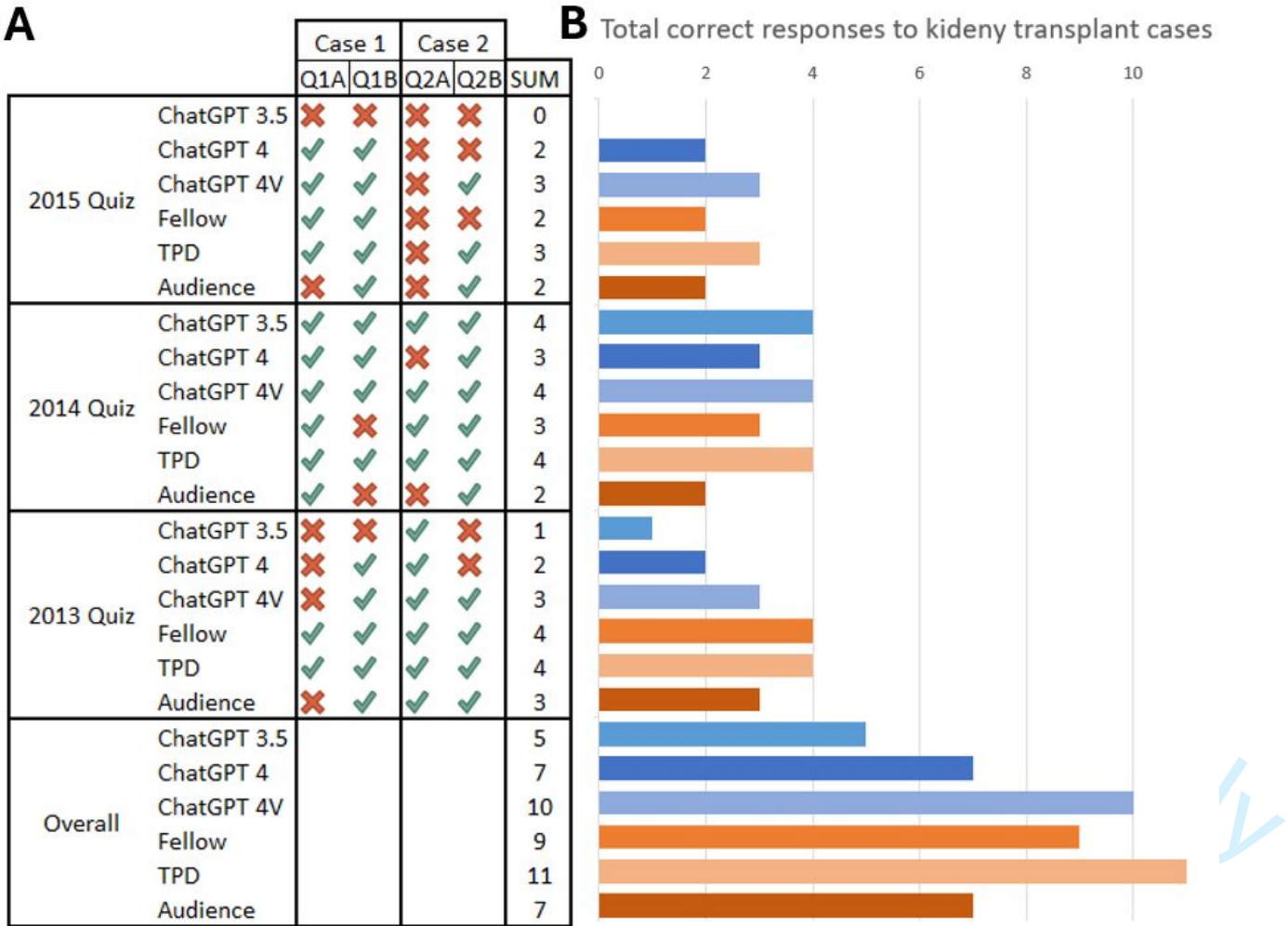
30. Hanna E, Levic A. Comparative Analysis of Language Models: hallucinations in ChatGPT: Prompt Study. Published online 2023.
31. Azamfirei R, Kudchadkar SR, Fackler J. Large language models and the perils of their hallucinations. *Critical Care*. 2023;27(1):120.

For Review Only

**Table 1. Summaries of the 2015, 2014, and 2013 kidney transplant cases from the American Society of Nephrology Quiz and related questions. Correct answers are bolded.**

Year	Case	Case Summary (by ChatGPT 4)	Lab results	Image	Question
2015	1	A 53-year-old woman with a history of ESRD, aortic stenosis, cancer, COPD, genitourinary reflux disease, and recent kidney transplant presents to the ER with shortness of breath, dyspnea, pleuritic chest pain, and cough. CT and CXR reveal dense consolidation in the left lower lobe, multiple areas of nodular opacities in both lungs, and an abnormal aortic valve, including thickened prosthesis leaflets indicating severe prosthetic aortic valve stenosis.	Laboratory results: Complete blood count	No	1A: Which of the following is the most likely cause of her chest pain and CXR findings? (A) Lung cancer, (B) Recurrent breast cancer, <b>(C) Pneumonia</b> , (D) Sarcoidosis
					1B: With the addition of voriconazole, which of the following changes would develop in her immunosuppression levels? <b>(A) Tacrolimus trough increased</b> , (B) Tacrolimus trough decreased, (C) MMF level increased, (D) MMF level decreased
	2	A 63-year-old woman with a history of non-alcoholic steatohepatitis, type 2 diabetes, hypertension, and end-stage renal disease treated with a kidney transplant, is admitted to the hospital for liver failure and volume overload and is being considered for a combined liver-kidney transplant. Her condition deteriorates with declining kidney function, rising serum creatinine, increased fatigue, decreased urine output, and the presence of numerous muddy brown granular casts in her urine sediment.	Laboratory results: Complete blood count	No	2A: What is the most likely underlying cause of the patient's decline in kidney function? (A) Hepatorenal syndrome, (B) Antibody-mediated rejection, (C) Acute tubular necrosis (ATN), (D) Recurrent diabetes mellitus, <b>(E) Impossible to predict</b>
					2B: Which of the following is likely contributing most to her massive volume overload? (A) Hepatorenal syndrome, (B) ATN, (C) Nephrotic syndrome, <b>(D) Heart failure</b>
2014	1	A 50-year-old man with a history of ESRD due to diabetes and hypertension, who had received a deceased donor kidney transplant, presented 9 months after operation with fever, fatigue, lightheadedness, loose stool and right upper quadrant pain. His medication list was extensive, including tacrolimus, mycophenolate mofetil, prednisone and trimethoprim sulfamethoxazole, and his lab results revealed anemia and underproduction of cells, with no evidence of malignancy, ulceration, or bleeding from prior colonoscopy and esophagogastroduodenoscopy.	(1) Laboratory findings, (2) Anemia laboratory evaluation	Post-transplant hematocrit	1A: Testing for which of the following may be most helpful in evaluating the underproduction anemia? (A) Cytomegalovirus (CMV), (B) BK virus, <b>(C) Parvovirus B19</b> , (D) Clostridium difficile, (E) JC virus
	2	A 52-year-old woman who received a kidney transplant was readmitted to the hospital with an altered mental status following a time period of post-transplant complications including pancreatitis, drain placement, and rising creatinine	Cerebrospinal fluid results	Brain imaging studies. (A) Head CT. (B) Brain MRI.	1B: Which of the following is most effective in treating Parvovirus B19? <b>(A) Intravenous Ig (IVIG)</b> , (B) Cidofovir, (C) Valganciclovir, (D) Acyclovir, (E) Immunosuppression reduction 2A: What is the most likely diagnosis? <b>(A) Posterior reversible encephalopathy syndrome (PRES)</b> , (B) CMV encephalitis, (C) Progressive multifocal leukoencephalopathy (PML), (D) Herpes simplex virus (HSV) encephalitis, (E) JC virus encephalitis,

		levels. Despite a normal diagnostic work-up, she displays signs of encephalopathy, with her most likely diagnosis being Posterior Reversible Encephalopathy Syndrome (PRES), often associated with the use of tacrolimus.			2B: Which of our patient's medications has most commonly been associated with PRES? (A) Valganciclovir, (B) Famotidine, <b>(C) Tacrolimus</b> , (D) Simvastatin
2013	1	A 66-year-old Laotian man who had a history of end-stage renal disease, hyperlipidemia, and chronic hepatitis B, and had received a kidney transplant, presented with a nonproductive cough and right flank pain, three months post-transplant. A CT scan revealed a large mass encasing a rib, and treatment for tuberculosis was initiated.		Mass encasing a rib. A CT scan image.	1A: What is the most likely diagnosis in this patient? (A) Recurrent thymoma, (B) Post-transplant lymphoproliferative disorder (PTLD), (C) Malignant BK nodule, <b>(D) Extrapulmonary tuberculosis (TB)</b> , (E) Brown tumor 1B: How might TB treatment affect the tacrolimus level? <b>(A) Decrease the tacrolimus level, because rifampin is a CYP3A4 inducer</b> , (B) Increase the tacrolimus level, because INH is a CYP3A4 competitor, (C) Decrease the tacrolimus level, because ethambutol is a CYP3A4 inducer, (D) No marked change, because by using INH and ethambutol, they will cancel each other's effect, (E) Increase the tacrolimus level, because pyrazinamide is a CYP3A5 inducer
	2	A 30-year old man with ESRD and blood type AB, who had received a deceased donor kidney, developed deep vein thrombosis in his upper extremity 2.5 years post-transplant. Despite an initial delay in graft function, he developed complications such as acute T cell mediated rejection, chronic allograft arteriopathy, chronic transplant glomerulopathy, peritubular capillaritis, and moderate interstitial fibrosis and tubular atrophy.	(1) laboratory studies at the time of kidney biopsy and 19 days later (2) Selected additional laboratory studies	(1) Peritubular capillaritis consistent with AMR, (2) C4d staining, (3) Lymphocytes lifting up endothelial cells diagnostic of Banff IIA rejection.	2A: Which of the following is the most likely mechanism of anemia in this patient? (A) Iron deficiency anemia secondary to a bleed, (B) Decreased production anemia caused by bone marrow suppression, (C) Dilutional anemia caused by volume overload, <b>(D) Intravascular hemolysis</b> 2B: What is the most likely cause of the hemolytic anemia in this patient? The rejection episode (A) Solumedrol, (B) IVIG, <b>(C) Occult infection</b> , (D) Sirolimus toxicit



**Figure 1. Human and GPT model performance on the American Society of Nephrology Quiz transplantation questions.**

**A)** Summative performance on each question to the 2015, 2014, and 2013 kidney transplant cases from the American Society of Nephrology Quiz are shown. We assumed that a given group answered a question correctly if the most frequently chosen answer was correct. In case two or more answers were most frequently chosen, we deemed these answers incorrect. Responses for fellows, training program directors (TPD), and audience members at the Annual Kidney Week Meeting of the American Society of Nephrology were reported in Josephson et al., 2013, 2014, and 2015. Chat GPT 3.5, Chat GPT 4, and Chat GPT 4V were prompted 10 times with each case and associated questions. **B)** Total correct responses for each group by quiz year and for all three years combined.



## Supplementary Materials

**Table S1. Question-specific performance on American Society of Nephrology Quiz transplantation questions for human respondents and GPT models.**

Answers to 2015, 2014, and 2013 kidney transplant cases from the American Society of Nephrology Quiz transplantation questions are shown. Response frequencies for fellows, training program directors (TPD), and audience members at the Annual Kidney Week Meeting of the American Society of Nephrology were reported in Josephson et al., 2013, 2014, and 2015. Chat GPT 3.5, Chat GPT 4, and Chat GPT 4V were prompted 10 times with each case and associated questions. Gray columns indicate correct answers. The most frequent responses are bolded and underlined for each of the responder's groups.

		Case 1										Case 2									
		Q1A					Q1B					Q2A					Q2B				
		A	B	C	D	E	A	B	C	D	E	A	B	C	D	E	A	B	C	D	E
2015 Quiz	ChatGPT 3.5	<u>70%</u>	0%	30%	0%		0%	<u>80%</u>	10%	10%		0%	0%	<u>100%</u>	0%	0%	<u>60%</u>	0%	0%	40%	
	ChatGPT 4	20%	0%	<u>80%</u>	0%		<u>100%</u>	0%	0%	0%		0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	
	ChatGPT 4V	0%	0%	<u>100%</u>	0%		<u>100%</u>	0%	0%	0%		0%	0%	<u>100%</u>	0%	0%	20%	0%	0%	<u>80%</u>	
	Fellow	13%	16%	<u>68%</u>	3%		<u>73%</u>	16%	0%	11%		6%	3%	<u>82%</u>	0%	9%	33%	<u>37%</u>	3%	27%	
	TPD	0%	19%	<u>81%</u>	0%		<u>94%</u>	0%	0%	6%		0%	0%	<u>73%</u>	0%	27%	0%	27%	27%	<u>46%</u>	
	Audience	13%	29%	17%	41%		<u>76%</u>	18%	0%	6%		19%	3%	<u>69%</u>	0%	9%	7%	21%	8%	<u>64%</u>	
2014 Quiz	ChatGPT 3.5	0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	<u>100%</u>	0%	
	ChatGPT 4	0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	0%	30%	0%	<u>70%</u>	0%	0%	0%	0%	<u>100%</u>	0%	
	ChatGPT 4V	0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	<u>100%</u>	0%	
	Fellow	11%	4%	<u>84%</u>	1%	0%	28%	5%	6%	1%	<u>60%</u>	<u>84%</u>	4%	3%	0%	9%	6%	0%	<u>94%</u>	0%	
	TPD	4%	0%	<u>96%</u>	0%	0%	<u>90%</u>	5%	0%	0%	5%	<u>90%</u>	5%	0%	0%	5%	0%	0%	<u>100%</u>	0%	
	Audience	26%	8%	<u>63%</u>	2%	1%	20%	3%	6%	1%	<u>70%</u>	<u>24%</u>	<u>24%</u>	13%	22%	17%	11%	3%	<u>86%</u>	0%	
2013 Quiz	ChatGPT 3.5	0%	<u>80%</u>	20%	0%	0%	20%	<u>80%</u>	0%	0%	0%	30%	0%	0%	<u>70%</u>		0%	<u>50%</u>	0%	50%	0%
	ChatGPT 4	0%	<u>100%</u>	0%	0%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	0%	<u>100%</u>		0%	<u>100%</u>	0%	0%	0%
	ChatGPT 4V	0%	<u>90%</u>	0%	10%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	0%	<u>100%</u>		0%	0%	<u>100%</u>	0%	0%
	Fellow	13%	28%	7%	<u>40%</u>	12%	<u>60%</u>	19%	7%	7%	7%	13%	13%	13%	<u>61%</u>		20%	7%	<u>46%</u>	7%	20%
	TPD	10%	21%	0%	<u>69%</u>	0%	<u>78%</u>	0%	11%	11%	0%	5%	5%	0%	<u>90%</u>		17%	0%	<u>44%</u>	0%	39%
	Audience	11%	<u>40%</u>	15%	27%	7%	<u>56%</u>	21%	5%	12%	6%	1%	20%	3%	<u>76%</u>		9%	2%	<u>46%</u>	7%	36%

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Review Only



## Supplementary Materials

Supplement to M Mankowski, IS Jaffe, J Xu, et al. ChatGPT Solving Complex Kidney Transplant Cases: A Comparative Study with Human Respondents.

### **Table of Contents**

#### ***Supplemental Tables***

Table S1.....page 2

For Review Only

**Table S1. Question-specific performance on American Society of Nephrology Quiz transplantation questions for human respondents and GPT models.** Answers to 2015, 2014, and 2013 kidney transplant cases from the American Society of Nephrology Quiz transplantation questions are shown. Response frequencies for fellows, training program directors (TPD), and audience members at the Annual Kidney Week Meeting of the American Society of Nephrology were reported in Josephson et al., 2013, 2014, and 2015. Chat GPT 3.5, Chat GPT 4, and Chat GPT 4V were prompted 10 times with each case and associated questions. Gray columns indicate correct answers. The most frequent responses are bolded and underlined for each of the responder’s groups.

		Case 1										Case 2									
		Q1A					Q1B					Q2A					Q2B				
		A	B	C	D	E	A	B	C	D	E	A	B	C	D	E	A	B	C	D	E
2015 Quiz	ChatGPT 3.5	<u>70%</u>	0%	30%	0%		0%	<u>80%</u>	10%	10%		0%	0%	<u>100%</u>	0%	0%	<u>60%</u>	0%	0%	40%	
	ChatGPT 4	20%	0%	<u>80%</u>	0%		<u>100%</u>	0%	0%	0%		0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	
	ChatGPT 4V	0%	0%	<u>100%</u>	0%		<u>100%</u>	0%	0%	0%		0%	0%	<u>100%</u>	0%	0%	20%	0%	0%	<u>80%</u>	
	Fellow	13%	16%	<u>68%</u>	3%		<u>73%</u>	16%	0%	11%		6%	3%	<u>82%</u>	0%	9%	33%	<u>37%</u>	3%	27%	
	TPD	0%	19%	<u>81%</u>	0%		<u>94%</u>	0%	0%	6%		0%	0%	<u>73%</u>	0%	27%	0%	27%	27%	<u>46%</u>	
	Audience	13%	29%	17%	41%		<u>76%</u>	18%	0%	6%		19%	3%	<u>69%</u>	0%	9%	7%	21%	8%	<u>64%</u>	
2014 Quiz	ChatGPT 3.5	0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	<u>100%</u>	0%	
	ChatGPT 4	0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	0%	30%	0%	<u>70%</u>	0%	0%	0%	0%	<u>100%</u>	0%	
	ChatGPT 4V	0%	0%	<u>100%</u>	0%	0%	<u>100%</u>	0%	0%	0%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	<u>100%</u>	0%	
	Fellow	11%	4%	<u>84%</u>	1%	0%	28%	5%	6%	1%	<u>60%</u>	<u>84%</u>	4%	3%	0%	9%	6%	0%	<u>94%</u>	0%	
	TPD	4%	0%	<u>96%</u>	0%	0%	<u>90%</u>	5%	0%	0%	5%	<u>90%</u>	5%	0%	0%	5%	0%	0%	<u>100%</u>	0%	
	Audience	26%	8%	<u>63%</u>	2%	1%	20%	3%	6%	1%	<u>70%</u>	<u>24%</u>	<u>24%</u>	13%	22%	17%	11%	3%	<u>86%</u>	0%	
2013 Quiz	ChatGPT 3.5	0%	<u>80%</u>	20%	0%	0%	20%	<u>80%</u>	0%	0%	0%	30%	0%	0%	<u>70%</u>		0%	<u>50%</u>	0%	50%	0%
	ChatGPT 4	0%	<u>100%</u>	0%	0%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	0%	<u>100%</u>		0%	<u>100%</u>	0%	0%	0%
	ChatGPT 4V	0%	<u>90%</u>	0%	10%	0%	<u>100%</u>	0%	0%	0%	0%	0%	0%	0%	<u>100%</u>		0%	0%	<u>100%</u>	0%	0%
	Fellow	13%	28%	7%	<u>40%</u>	12%	<u>60%</u>	19%	7%	7%	7%	13%	13%	13%	<u>61%</u>		20%	7%	<u>46%</u>	7%	20%
	TPD	10%	21%	0%	<u>69%</u>	0%	<u>78%</u>	0%	11%	11%	0%	5%	5%	0%	<u>90%</u>		17%	0%	<u>44%</u>	0%	39%
	Audience	11%	<u>40%</u>	15%	27%	7%	<u>56%</u>	21%	5%	12%	6%	1%	20%	3%	<u>76%</u>		9%	2%	<u>46%</u>	7%	36%

For Review Only