

## **ChatGPT goes to operating room: Evaluating GPT-4 performance and future direction of surgical education and training in the era of large language models**

Namkee Oh, M.D., Gyu-Seong Choi, M.D., Ph.D., Woo Yong Lee, M.D., Ph.D.\*

Department of Surgery, Samsung Medical Center, Sungkyunkwan University School of Medicine, Seoul, Korea

### **\* Corresponding author**

Professor

Woo Yong Lee, M.D., Ph.D.

Department of Surgery, Samsung Medical Center, Sungkyunkwan University School of Medicine

Address: 81 Irwon-ro, Gangnam-gu, Seoul, Republic of Korea 06351

Tel: +82-2-3410-0261

Fax: +82-2-3410-6980

E-mail: [wooyong123.lee@samsung.com](mailto:wooyong123.lee@samsung.com)

## **Abstract**

### **Purpose**

This study aimed to assess the performance of ChatGPT, specifically the GPT-3.5 and GPT-4 models, on the Korean general surgery board exam questions and investigate the potential applications of large language models (LLM) for surgical education and training.

### **Method**

The dataset comprised 280 questions from the Korean general surgery board exams conducted between 2020 and 2022. Both GPT-3.5 and GPT-4 models were evaluated, and their performance was compared using the chi-square test.

### **Result**

GPT-3.5 achieved an overall accuracy of 46.8%, while GPT-4 demonstrated a significant improvement with an overall accuracy of 76.4%, indicating a notable difference in performance between the models ( $p < .001$ ). GPT-4 also exhibited consistent performance across all subspecialties, with accuracy rates ranging from 63.6% to 83.3%.

### **Conclusion**

ChatGPT, particularly GPT-4, demonstrates a remarkable ability to understand complex surgical clinical information, achieving an accuracy rate of 76.4% on the Korean general surgery board exam. As LLM technology continues to advance, its potential applications in surgical education, training, and continuous medical education (CME) are anticipated to enhance patient outcome and safety.

## Introduction

Recently, significant advancements in large language model (LLM) technology have revolutionized the field of artificial intelligence (AI), with ChatGPT released by OpenAI in November 2022 standing out as a prime example.[1] ChatGPT has exhibited exceptional performance in evaluating knowledge related to fields such as medicine, law, and management, which have traditionally been considered the domain of experts. Notably, the system achieved high accuracy on the USMLE, the Bar exam, and the Wharton MBA final exam, even without fine-tuning the pre-trained model.[2-5]

Surgical education and training demand a significant investment of time, with the process involving a combination of didactic learning, hands-on training, and supervised clinical experience.[6] During residency, surgical trainees work alongside experienced surgeons, gaining practical experience in patient care, surgery, and clinical decision-making. Additionally, trainees engage in a series of didactic courses and conferences covering the principles of surgery, medical knowledge, and surgical techniques. Due to the comprehensive nature of surgical education and training, it can take more than a decade to become a skilled and competent surgeon. Given the time-intensive nature of surgical education and training, it is important to explore how emerging technologies, such as AI and LLMs, can augment the learning process.[7]

This study aims to employ ChatGPT to evaluate the general surgery board exam in Korea and assess whether LLMs possess expert-level knowledge. Moreover, the study compared the performance of GPT-3.5 and GPT-4. By exploring how LLMs can be integrated into clinical education and practice, this study ultimately aims to contribute to the improvement of education and training for surgical residents and practicing surgeons.

## **Methods**

### **General surgery board exam of Korea**

The goal of surgical education and training is to develop the ability to actively evaluate the pathological conditions of surgical diseases and to acquire the surgical skills to treat traumatic, congenital, acquired, neoplastic, and infectious surgical diseases. To quantitatively evaluate this knowledge and skill set of surgical residents, a board certification exam is required after completion of their training, in order to become a board-certified general surgeon of Korea. The exam is composed of two stages: the first stage is a 200-question multiple-choice test, and those who pass the first stage are eligible to take the second stage. The second stage consists of questions based on high-resolution clinical images and surgical video clips. The questions are created and supervised by the Korean Surgical Society (KSS) and the Korean Academy of Medical Science (KAMS).

### **Dataset for model testing**

The actual board exam questions are held by KAMS, but due to limited access to the usage of these questions, we constructed our dataset by gathering questions recalled by examinees who took the actual exam. As the LLM cannot process visual information such as clinical images, radiology, and graphs, questions that included visual information were excluded from our dataset. All problems were manually inputted in their original Korean text. Finally, our dataset included a total of 280 questions from the first stage of the board exam in 2020, 2021, and 2022.

### **Large language model and performance evaluation**

In this study, we utilized the ChatGPT generative pre-trained transformer (GPT) language model developed by OpenAI to evaluate its performance on a dataset of questions. We performed model testing using both GPT-3.5 and GPT-4, with the

former conducted from March 1st to March 3rd, 2023, and the latter scheduled for March 15th, 2023. To evaluate the model's performance, we manually entered the questions into the ChatGPT website and compared the answers provided by the model to those provided by examinees.

## Statistical analysis

This study compared the performance of the GPT-3.5 and GPT-4 models with the chi-square test. A p-value less than a 0.05 would indicate a statistically significant difference between the performance of the GPT-3.5 and GPT-4.

## Result

The dataset used for model evaluation consisted of a total of 280 questions, which were classified into subspecialties and listed in order of frequency as follows: endocrine (16.8%), breast (16.1%), lower gastrointestinal (LGI, 14.3%), upper gastrointestinal (UGI, 13.2%), general (13.2%), pediatric (6.4%), hepatobiliary and pancreas (HBP, 6.1%), vascular (6.1%), transplantation (4.0%), and trauma and critical care(4.0%). (Fig. 1)

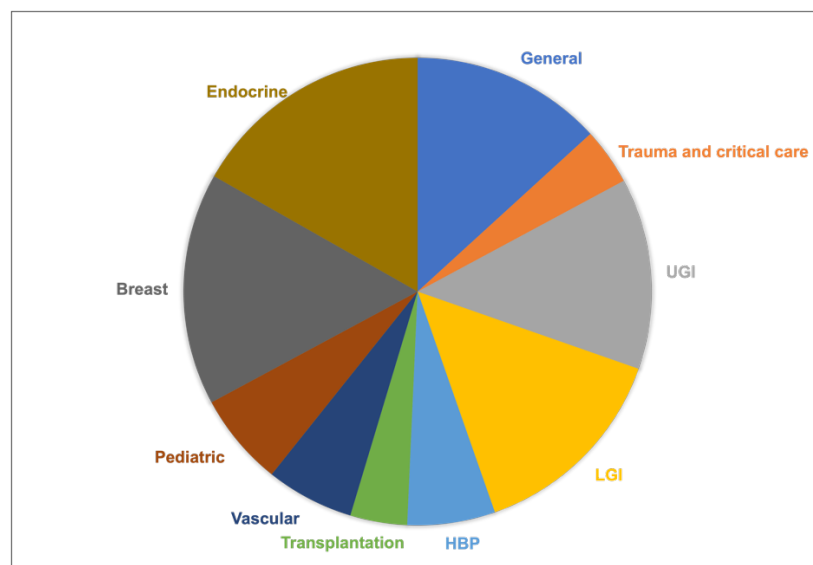


Figure 1. The dataset was composed of 280 questions, and it is classified into subspecialties in the field of general surgery.

A significant difference in performance was observed between the GPT-3.5 and GPT-4 models ( $p < .001$ ). The GPT-3.5 model achieved an overall accuracy of 46.8%, providing correct answers for 131 out of the 280 questions. In terms of individual subspecialties, the model's accuracy rates were as follows (sorted from highest to lowest): transplantation (72.7%), breast (62.2%), HBP (52.9%), general (48.6%), UGI (45.9%), trauma and critical care (45.5%), LGI (45.0%), endocrine (36.2%), pediatric (33.3%), and vascular (29.4%). In contrast, the GPT-4 model demonstrated a substantial improvement in overall accuracy, attaining a rate of 76.4% by providing correct answers for 214 out of the 280 questions. The accuracy rates for each subspecialty were as follows: pediatric (83.3%), breast (82.2%), UGI (81.1%), endocrine (78.7%), general (75.7%), transplantation (72.7%), LGI (72.5%), vascular (70.6%), HBP (64.7%), and trauma and critical care (63.6%). (Fig. 2)

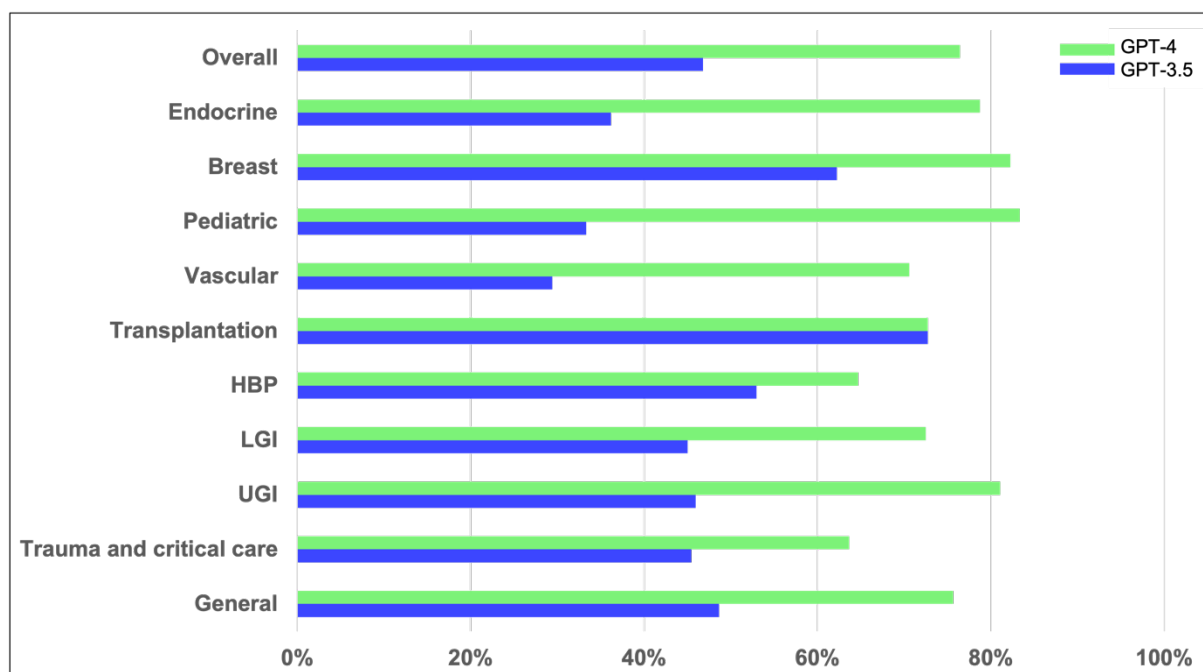


Figure 2. Comparison of the performance of GPT-4 and GPT3.5 with overall accuracy and accuracies according to its subspecialties.

## Discussion

The primary objective of this study was to conduct a quantitative assessment of ChatGPT's ability to comprehend complex surgical clinical information and to explore the potential implications of LLM technology for surgical education and training. Specifically, we tested the performance of ChatGPT using questions from the Korean general surgery board exam and observed that the model achieved an accuracy of 76.4% with GPT-4 and 46.8% with GPT-3.5. Remarkably, this accuracy was achieved without fine-tuning the model and by using prompts in Korean language exclusively, thus highlighting the significance of our findings.

In a study evaluating GPT-3.5's performance on the United States Medical Licensing Examination (USMLE), the model achieved accuracies of 41.2%, 49.5%, and 59.8% for Step 1, CK2, and 3, respectively, using multiple choice single answer questions with forced justification for selection.[2] Additionally, the model achieved accuracies of 64.4%, 57.8%, 44%, and 42% using the National Board of Medical Exam (NBME)-Free Step1 and 2, and AMBOSS-Step1 and Step2 question banks, respectively.[8] In another study using the American Academy of Ophthalmology's Basic and Clinical Science Course (BCSC) question bank, which is specific to the field of ophthalmology, the model showed an accuracy of 55.8%, while the Ophthalmic Knowledge Assessment Program (OKAP) exam an accuracy of 42.7%.[9] These results suggest that the performance of GPT 3.5 on the Korean general surgery board exam questions, overall accuracy of 46.8%, is comparable to the performance of previous studies.

The comparative analysis revealed a notable improvement in GPT 4's performance compared to GPT 3.5 model across all subspecialties. GPT 4 not only exhibited a higher overall accuracy rate but also demonstrated more consistent

performance in each subspecialty, with accuracy rates ranging from 63.6% to 83.3%. Although the current publicly available GPT 4 model on website can only process text-based information, the GPT 4 framework itself is capable of processing and analyzing visual information, including images and videos.[10] This capability raises the possibility that, in the future, the performance of GPT 4 could be evaluated on datasets containing clinical photos and surgical videos. Such advancements would further enhance the applicability of GPT 4 in medical and surgical fields, broadening its utility beyond text-based tasks and offering a more comprehensive understanding of complex clinical scenarios assisting professionals in their decision-making processes and contributing to improved patient care.

The authors kindly recommend that the surgeon's society proactively adapts and utilizes these technological advancements to enhance patient safety and improve the quality of surgical care. In the context of surgical education, it is crucial to transition from the traditional rote learning approach to a method that emphasizes problem definition in specific clinical situations and the acquisition of relevant clinical information for problem resolution. Large Language Models (LLMs) serve as generative AI models, providing answers to given problems. Consequently, the quality of the answers relies on the questions posed.[11] Surgeons must conduct thorough history taking and physical examinations to accurately define the problems they face. By providing LLMs with comprehensive summaries of patients' chief complaints, present illnesses, and physical examinations, the models can offer valuable guidance on diagnostic tests and treatment options. Ultimately, it is essential for medical professionals to return to the fundamentals, maintaining close connections with patients and actively listening to their concerns.[12]

Moreover, active surgeons who completed their training over a decade ago



can benefit from using LLMs for continuous medical education (CME). Accessing new knowledge may be challenging for them due to the time elapsed since their training, potentially leading to outdated management practices. While numerous surgical societies offer CME programs, altering ingrained routines in clinical practice can be difficult. By maintaining an up-to-date LLM and integrating it into their decision-making processes, surgeons can consistently deliver the highest level of evidence-based care to their patients.[13]

In medicine, decision-making has a profound impact on patient safety, demanding a higher level of accuracy and a conservative approach to change compared to other fields. Although GPT-4 achieved a 76.4% accuracy rate on the Korean surgical board exam, it is not yet sufficient for immediate clinical application in patient care. However, it is noteworthy that a service released less than six months ago exhibits such remarkable performance. Furthermore, ChatGPT is only one example of LLMs. Recently, Microsoft released BioGPT, an LLM trained on PubMed literature, and Meta introduced LLaMA, an LLM with an accessible API for open innovation and fine-tuning.[14, 15] Based on these trends, we can anticipate future LLMs to be trained on an even larger and more diverse set of medical information, providing specialized knowledge in the medical field. Undoubtedly, the ultimate goal is to enhance the quality of care and improve patient outcomes and safety, and emphasizing this objective cannot be overstated.

The limitations of this study include the fact that the dataset was compiled using questions recalled by examinees, which may not accurately represent the full set of actual board exam questions due to restricted access. Another limitation is the exclusion of visual information. Since the models used in the study are unable to process visual information, such as clinical images, radiology, and graphs, questions

containing visual components were excluded from the dataset. As a result, we cannot determine whether ChatGPT would pass or fail the board exam based on these limitations. Despite these constraints, this study holds significance as it confirms the ability of LLMs to analyze surgical clinical information and make appropriate clinical decisions.

## Conclusion

ChatGPT, particularly GPT-4, demonstrates a remarkable ability to understand complex surgical clinical information, achieving an accuracy rate of 76.4% on the Korean general surgery board exam. As LLM technology continues to advance, its potential applications in surgical education, training, and continuous medical education (CME) are anticipated to enhance patient outcome and safety.

## Figure legend

1. OpenAI, <https://openai.com/blog/chatgpt>. 2022.
2. Kung, T.H., et al., *Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models*. PLOS Digital Health, 2023. **2**(2): p. e0000198.
3. Mbakwe, A.B., et al., *ChatGPT passing USMLE shines a spotlight on the flaws of medical education*. 2023, Public Library of Science San Francisco, CA USA. p. e0000205.
4. Bommarito II, M. and D.M. Katz, *GPT Takes the Bar Exam*. arXiv preprint arXiv:2212.14402, 2022.
5. Choi, J.H., et al., *Chatgpt goes to law school*. Available at SSRN, 2023.
6. Debas, H.T., et al., *American surgical association blue ribbon committee report on surgical education: 2004*. Annals of surgery, 2005. **241**(1): p. 1-8.
7. Wartman, S.A. and C.D. Combs, *Medical education must move from the information age to the age of artificial intelligence*. Academic Medicine, 2018. **93**(8): p. 1107-1109.
8. Gilson, A., 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 Medical Education, 2023. **9**(1): p. e45312.

9. Antaki, F., et al., *Evaluating the Performance of ChatGPT in Ophthalmology: An Analysis of its Successes and Shortcomings*. medRxiv, 2023: p. 2023.01.22.23284882.
10. OpenAI, *GPT-4 Technical Report*. 2023.
11. Radford, A., et al., *Improving language understanding by generative pre-training*. 2018.
12. Kapadia, M.R. and K. Kieran, *Being affable, available, and able is not enough: prioritizing surgeon-patient communication*. JAMA surgery, 2020. **155**(4): p. 277-278.
13. Han, E.-R., et al., *Medical education trends for future physicians in the era of advanced technology and artificial intelligence: an integrative review*. BMC medical education, 2019. **19**(1): p. 1-15.
14. Luo, R., et al., *BioGPT: generative pre-trained transformer for biomedical text generation and mining*. Briefings in Bioinformatics, 2022. **23**(6).
15. Touvron, H., et al., *Llama: Open and efficient foundation language models*. arXiv preprint arXiv:2302.13971, 2023.