The Educational Potential of ChatGPT: Assessing Large Language Model ChatGPT Responses to Common Patient Questions Regarding Total Knee Arthroplasty

Yasir AlShehri, Gerard Anthony Sheridan, Lisa Howard¹, Donald S Garbuz², Bassam A Masri³, Michael Neufeld³ ¹VGH, ²UNIVERSITY OF BRITISH COLUMBIA, ³University of British Columbia INTRODUCTION:

Artificial intelligence (AI) has gained significant popularity in the medical field and holds enormous potential to improve health outcomes. Machine learning protocols in Orthopaedics are becoming more widespread. ChatGPT, a large language model based on AI technology, has become mainstream in recent years. ChatGPT uses natural language processing and deep learning to process vast amounts of data and generate response with a natural conversational flow. Effective patient education has been found to decrease patient anxiety, postoperative complications, and promote a faster recovery. Currently, the evidence on the use of Large Language Models (LLMs) such as ChatGPT for patient education is limited. The purpose of our study was to assess ChatGPT's ability to accurately and adequately answer frequently asked patient questions related to TKA, based on expert opinion from specialized surgeons in the field of TKA.

METHODS:

Ten guestions were obtained from the frequently asked guestions (FAQ) patient education portal of The American Association of Hip and Knee Surgeons (AAHKS) website (https://hipknee.aahks.org/total-knee-replacement/) . The questions were selected and approved by the authors by consensus. Most questions required minor modifications to clarify the context to TKA (e.g changing "How long will it last?" to "How long will my knee replacement last?"). Each question was inputted into ChatGPT 3.5 on January 29, 2024, without additional prompts or limitations on response length (Figure 1). No alterations were made to the responses, and a new session was initiated for each question to remove context from previous interactions. The responses were evaluated by four subspecialty arthroplasty surgeons.

Quality analysis was performed using a grading system based on ones used in previously published literature. The ChatGPT responses, upon review by the surgeons, were graded from A to D (Table 1). Following individual grading of all responses by each surgeon, a consensus was attempted to be reached. Response readability was assessed using the Flesch-Kincaid Reading Ease Score (FRES) and Grade Level (FKGL), two widespread tools for assessing readability. The FRES provides a score from 0 to 100, with higher scores indicating easier readability. For example, scores of 70-80 correlate with a 7th-grade reading level, while scores of 30-50 reflect college-level difficulty. The FKGL correlates with the educational level required to comprehend the text, ranging from elementary to college levels in the US educational system.

RESULTS:

The ten responses were graded as follows (Table 2); A (Correct and sufficient response): 20% (n=2); B (Correct but insufficient response): 40% (n=4); C (Response containing correct and incorrect information): 30% (n=3); D (Incorrect response): 0% (n=0). One response (10%) was not given a concensus grade due to lack of agreement between the grading surgeons. The mean FRES was 25.3 [SD± 8.0], which corresponds to a college graduate level, and ranged from 10.4 to 41.1. The mean FKGL was 15.4 [SD±1.5], ranging from 12.9 to 18.5, also corresponding to a college graduate reading level (Table 3). The ten questions, response grades, and analyses are presented below.

DISCUSSION AND CONCLUSION:

This study aimed to assess the adequacy and accuracy of ChatGPT's responses to common TKA-related questions. The results show that it can provide responses with overall satisfactory accuracy to these questions. However, inaccuracies were encountered, and there was variability in the surgeons' grading of the responses. A critical distinction in our study is that accuracy does not presume adequacy. While ChatGPT may provide reasonable accuracy, it is often inadequate. ChatGPT in its current state does not incorporate intelligence but rather works on pattern recognition and machine learning platforms. As such, it can masquerade as sufficient but can often provide incorrect information. Patients need to be educated on its limitations as it becomes a commonly utilized tool.

In conclusion, ChatGPT provides accurate responses to common patient questions about TKA but lacks the depth of inperson discussions with surgeons. As the interface improves over time, outcomes will likely improve. According to the results of this study, ChatGPT can be a useful tool for patients. Patients should be encouraged to use such tools to enhance their understanding and facilitate informed consent discussions with their surgeon.



	Response grading system
٨	Corrort and sufficient response
8	Correct but insufficient response

Response Number	Response Quality Grade				
	Surgcon 41	Surgeon 12	Sergeon (3	Sergeon #4	Centerns
1	۸	в	A	A	А
2	в	с	с	Α.	NC
)	в	Α	8	8	B
4	8	в	с	с	c
5	в	в	в	A	в
6	с	с	с	с	с
т	c	c	c		c
8	^	A	A	A	A
9	^	в	в		в
10	п	в	в	A	в

Reporte Number	Front-Nitrodal Hearing Ears Sover 19800-91205 Shi guade Shila-51205 Shi guade Shila-51205 Shi guade Yana Alahi shi di Shi guadi Shila-5206 Chilege yashasino.	Floods Alsocial Grant Level 6-3: Kilostategarias 31-4: Elementary 6.5-9: Middle school 52: 1-0: High school 12:1-15: Collage 15:1-16: Collage 15:1-16: Collage
1	24.5	15.1
2	33.1	13.9
x.	355	14.7
+	16.9	16.0
5	41.1	12.9
6	342	15.3
7	10.4	18.5
8	243	17.3
9	30.8	14.8
10	22.5	15.9

Level (FRGL) for the responses.