

ChatGPT prompts for generating multiple-choice questions in medical education and evidence on their validity: a literature review

Yavuz Selim Kiyak^{1,*}, Emre Emekli²

¹Department of Medical Education and Informatics, Faculty of Medicine, Gazi University, Ankara 06500, Turkey

²Department of Radiology, Faculty of Medicine, Eskişehir Osmangazi University, Eskişehir 26040, Turkey

*Corresponding author. Department of Medical Education and Informatics, Gazi Üniversitesi Hastanesi, E Blok 9. Kat, 06500 Beşevler, Ankara, Turkey.

E-mail: yskiyak@gazi.edu.tr

Abstract

ChatGPT's role in creating multiple-choice questions (MCQs) is growing but the validity of these artificial-intelligence-generated questions is unclear. This literature review was conducted to address the urgent need for understanding the application of ChatGPT in generating MCQs for medical education. Following the database search and screening of 1920 studies, we found 23 relevant studies. We extracted the prompts for MCQ generation and assessed the validity evidence of MCQs. The findings showed that prompts varied, including referencing specific exam styles and adopting specific personas, which align with recommended prompt engineering tactics. The validity evidence covered various domains, showing mixed accuracy rates, with some studies indicating comparable quality to human-written questions, and others highlighting differences in difficulty and discrimination levels, alongside a significant reduction in question creation time. Despite its efficiency, we highlight the necessity of careful review and suggest a need for further research to optimize the use of ChatGPT in question generation.

Main messages

- Ensure high-quality outputs by utilizing well-designed prompts; medical educators should prioritize the use of detailed, clear ChatGPT prompts when generating MCQs.
- Avoid using ChatGPT-generated MCQs directly in examinations without thorough review to prevent inaccuracies and ensure relevance.
- Leverage ChatGPT's potential to streamline the test development process, enhancing efficiency without compromising quality.

Keywords: ChatGPT; artificial intelligence; automatic item generation; ChatGPT prompts; multiple-choice questions

Introduction

ChatGPT has emerged as a useful companion to humans with potential benefits in medical education [1]. Its adoption in medical education can be particularly significant, where medical educators are continually seeking innovative methods to enhance learning [2] and assessment. Among the potential innovations, the use of ChatGPT for generating multiple-choice questions (MCQs) has captured interest due to this format's ubiquity in written assessments. The format is ubiquitous because MCQs serve as a useful, which means efficient, scalable, and cost effective, materials in assessing medical students' and residents' knowledge, clinical reasoning, and problem-solving skills [3].

However, while ChatGPT's utility in creating educational content is promising, the validity of the MCQs generated through this artificial intelligence (AI)-driven approach remains unexplored. Validity, in the context of educational assessments, is crucial, ensuring that the questions accurately measure what they are intended to. The quality of the prompts given to ChatGPT directly influences the quality and relevance of generated

MCQs, thereby impacting the validity of MCQs. Despite its critical importance, there exists a noticeable gap in the literature regarding the validity of MCQs generated by ChatGPT in medical education.

This gap is significant as the integration of AI tools in educational settings accelerates. Cross-sectional studies showed that medical educators and students, both in undergraduate and postgraduate levels, use ChatGPT for various purposes [4–6] including for question generation [7]. Therefore, it necessitates an understanding of assessment validity of ChatGPT-generated MCQs. The previous reviews were too broad and have not fully explored and documented the prompts whether they generate high-quality MCQs [8].

Addressing this gap is crucial. It will provide medical educators with insights into the effective use of AI for MCQ generation. It also will contribute to the broader discourse on the integration of artificial intelligence in educational practices, particularly the field of medicine requiring high levels of accuracy. By identifying the strengths and limitations of ChatGPT-generated MCQs, this

Received: February 28, 2024. **Revised:** April 29, 2024. **Accepted:** May 23, 2024

© The Author(s) 2024. Published by Oxford University Press on behalf of Postgraduate Medical Journal.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

research can guide the development of optimal strategies for AI-assisted test development in medical education.

Therefore, this rapid review aims to answer the following research questions:

1. What prompts have been used for MCQ generation using ChatGPT in medical education?
2. What is the validity evidence of the MCQs generated using ChatGPT in medical education?

Materials and methods

We opted for conducting a rapid review due to the urgent need for information about the use of AI, specifically ChatGPT, in generating MCQs for medical education. Given the fast-paced developments in AI and its growing application in medical education settings, there was a pressing demand for immediate guidance and understanding that could inform educators and test developers. This approach allowed for a quicker synthesis of available evidence, providing early, actionable findings to stakeholders in a rapidly evolving field. We followed the practical guide for rapid reviews [9] to ensure a methodologically sound and efficient approach.

Search strategy

The literature search was designed to identify studies involving the use of ChatGPT for the generation of MCQs in medical education. Therefore, the search terms used were (“ChatGPT” OR “Chat GPT”) AND (“medicine” OR “medical”). The search was conducted without any time restrictions and up until February 15, 2024. The databases searched included PubMed and Web of Science, adhering to the recommendation of searching at least two electronic databases [9]. The search included all fields in the studies. The initial search was conducted by one researcher, yielding 1690 studies from PubMed and 1063 studies from Web of Science. After removing duplicates with the help of a reference management tool, 1920 articles remained for title and abstract screening. After the screening process, we also searched the reference lists of the eligible studies [9].

Inclusion and exclusion criteria

The selection criteria included any study that reported the use of a prompt for generating MCQs within a medical education context. We did not limit the type of publications (e.g. research articles, reviews, commentaries). Exclusion criteria were studies not published in English and those outside of medical education (other health professions education, such as, nursing, dentistry, and pharmacy).

Screening of the studies

To enhance the reliability, two reviewers involved in the screening process. The screening began with titles and abstracts together, resulting in 158 studies selected for full-text screening in terms of the criteria. Following the full-text screening, only 21 articles met the criteria. The search in the reference lists added two more studies suitable for full review, bringing the total to 23 studies. Since we invested considerable effort to achieve a shared understanding before screening, we encountered minimal number of discrepancies between the two reviewers. These were resolved through discussion.

Extraction of the data

Data extraction was carried out by the experienced researcher on ChatGPT for question generation (YSK), focusing on key information such as author(s), date of publication, publication type, ChatGPT version used, the prompts employed, the number of MCQs generated, any further modifications made to the questions, the basis of evaluation presented by the studies, and the outcomes. Although the guide recommended as a reasonable approach that a second reviewer should check a 10% random sample for accuracy [9], in this study, the second reviewer (EE) verified the extracted data from all 23 studies. This process identified and corrected five instances of inconsistency or error overlooked by the first reviewer in the initial review process.

Assessment of the results

Prompts were evaluated considering the prompt engineering strategies and tactics described in a guide provided by OpenAI [10]. Some of them are: Write clear instructions, provide/mention examples, instructing to answer using a reference text, and ask the model to adopt a persona. The complete list of these strategies and tactics is presented in Table 1.

Validity was assessed by considering five major sources of validity evidence in testing [11]: First, “content” examines the relationship between test content and the intended construct. Second, “response process” focuses on individual responses and their alignment with intended interpretations. Third, “internal structure” focuses on the reliability and statistical characteristics of items within the assessment, such as, item difficulty and discrimination, and functionality of distractor options. Fourth, “relations to other variables” involves external data analysis, such as correlations with independent measures. Fifth, “evidence based on consequences of testing” focuses on the impact of assessments on individuals, institutions, and society, considering both intended and unintended effects.

Results

What prompts have been used for MCQ generation using ChatGPT in medical education?

We presented the prompts in Table 2. Four studies did not provide the prompt they used [12–15] and five studies did not report the version of ChatGPT [12, 16–19].

Six studies asked ChatGPT to generate MCQs based on a known style, such as, American Board of Dermatology Applied Exam [20], The United States Medical Licensing Examination (USMLE) [12, 17, 21, 22], and The National Board of Medical Examiners (NBME) [19]. This strategy aligns with providing/mentioning examples as a prompt engineering tactic.

Three studies used another prompt engineering tactic by asking the model to adopt a persona. These prompts were “You are a developer of teaching materials ...” [23] and “You are developing a question bank for medical exams ...” [24, 25]. Three studies submitted text for ChatGPT to generate responses using the provided text as a reference [18, 20, 26], a tactic that is also recommended. Five studies were so kind in their prompts as they used “please” to ask the model to generate questions [22, 27–30], even if it is not a prompt engineering tactic.

Table 1. OpenAI guidance on prompt engineering.

Strategy	Tactics	When/Why to use
Write clear instructions	Include details in your prompt	When specificity is needed for relevance.
	Ask the model to adopt a persona	To guide the model's tone or perspective.
	Use delimiters to clearly indicate parts	For complex inputs requiring clear separation.
	Specify steps to complete a task	For tasks that can be broken down into sequential actions.
	Provide examples	To model a desired format or approach.
Provide reference Text	Specify desired output length	When output needs to be within certain length constraints.
	Instruct to answer using a reference text	To base responses on provided material for accuracy.
	Instruct to answer with citations from reference text	When referencing specific information directly from provided texts.
Split complex tasks into simpler subtasks	Use intent classification	To identify and address specific parts of a complex query.
	Summarize/filter previous dialog Summarize long documents piecewise	In long dialogs, to maintain relevance and context. When dealing with content that exceeds model's context window (allowed input load).
Give time to "think"	Instruct model to work out its own solution before concluding	For reasoning tasks requiring step-by-step processing.
	Use inner monolog or sequence of queries	To internally process reasoning before providing an answer.
	Ask if anything was missed previously	To ensure comprehensiveness in tasks requiring thorough exploration.
Use external tools	Implement embeddings-based search	For enhancing model's knowledge with external information.
	Use code execution for calculations or API calls	For tasks requiring precise calculations or external data.
	Give access to specific functions	To benefit from specialized functionalities in specific contexts.
Test changes systematically	Evaluate outputs with reference to gold-standard answers	Testing is always helpful.

What is the validity evidence for the MCQs generated using ChatGPT in medical education?

Content

The generated questions were on various domains and diseases: Physiology [16, 31], dermatology [20], anatomy [13, 14, 30], immunology [32], internal medicine [26], surgery [26] and neurosurgery [14], anesthesia [15], clinical pharmacology [25], carpal tunnel syndrome [12], syndrome of inappropriate antidiuretic hormone secretion [21], diabetes [28], hyperlipidemia [29], anterior cutaneous nerve entrapment syndrome [23], urinary tract infection [33], hypertension [24, 25], chondroid tumors [18], and reproductive system [19].

The rates of MCQs that have inaccuracies based on expert reviews were 60% [20], 21% and 37% [34], 1% and 15% [27], and 16% [31]. One study reported that all MCQs were considered acceptable in an expert panel [25]. However, each study used different criteria with small number of experts. One study found that a large proportion of the ChatGPT-generated questions were nearly identical [34].

Response process

All studies generated single best answer MCQs in English, some with five options, and others with four options.

Internal structure

Three studies conducted item analyses [19, 25, 31]. In one study, two case-based MCQs produced ideal discrimination levels (above 0.30), which are 0.41 and 0.39 [25], in an exam consists of 25 MCQs (23 were human-written). Difficulty levels were 0.78 and 0.58 [25], which are easy and moderate, respectively. In another study, the average of discrimination scores of 21 MCQs was 0.24, which is

acceptable (above 0.20) but some particular MCQs had unacceptable discrimination levels [31]. Average of difficulty levels of 21 MCQs was 0.62 [31]. In a study that carried out an editing process after generating MCQs, 29 MCQs' discrimination level was 0.23, which is not ideal but acceptable, and average difficulty level was 0.71 [19].

Relations to other variables

Some studies compared different aspects of ChatGPT-generated and human-written MCQs. In one study based only on expert reviews, "appropriateness, clarity and specificity, discriminative power, and suitability for medical school exams" of MCQs generated by using ChatGPT were comparable to human-written MCQs but inferior in the relevance criterion [26]. In one study based on administering MCQs in exams, human-written MCQs had significantly higher levels of discrimination but similar level of difficulty [31]. In another study, ChatGPT-based MCQs had similar levels of discrimination and difficulty with human-written MCQs but this study carried out further modifications on the MCQs before administering them [19].

Evidence based on consequences of testing

Although several studies mentioned that ChatGPT can expedite item writing process, one study reported that ChatGPT reduced MCQ creation time from an estimated 30–60 minutes per MCQ to 5–15 minutes but these were case-based MCQs [19]. Another study reported that ChatGPT generated 50 MCQs (not case-based) within 20 minutes and 25 seconds, whereas humans spent 211 minutes and 33 seconds for writing 50 MCQs [26].

Table 2. ChatGPT prompts for generating multiple-choice questions and their characteristics published in the literature.

Author(s) and date	Publication type	ChatGPT version	Prompt	Number of MCQs	Further modification	Evaluation based on	Outcome
Agarwal et al., January 2023	Research article	Unspecified	"Generate 5 difficult reasoning-based MCQs for MBBS undergraduates on [Topic]"	110	No	Expert review by two physiologists	The questions were "highly valid", "somewhat difficult or easy", and "somewhat required reasoning".
Ayub et al., August 2023	Research article	An app based on ChatGPT-3.5	"Create five American Board of Dermatology Applied Exam (ABD-AE)-style questions"	40	No	Expert review by two dermatologists	"... only 16 (40%) questions created using ChatGPT 3.5 were accurate and at an appropriate level of complexity for a trainee studying for ABD-AE."
Benitez et al., January 2024	Special article (Perspective)	Unspecified	Unspecified but it might include "USMLE-style"	4	No	Informal review	"the complexity of the questions appeared to be of a lower order"
Biswas, May 2023	Special article (Editorial)	Unspecified	"Provide some MCQs for USMLE step 1 exams."	5	No	None	None
Cheung, August 2023	Research article	ChatGPT-3.5	"Can you write a multiple-choice question based on the following criteria, with the reference I am providing for you and your medical knowledge?"	50	"Interaction is only allowed for clarification but not any modification."	Expert review by five experts (for a comparison to 50 human-written MCQs)	Appropriateness, clarity and specificity, discriminative power, and suitability for medical school exams of AI MCQs are comparable to human-written MCQs, AI was inferior only in the relevance criterion.
Divito et al., November 2023	Special article (New Wave)	ChatGPT-3.5	"For each learning issue, can you provide a USMLE style multiple choice (A-E) question that tests student knowledge on the topic? Provide the correct answer and a justification of the correct choice for each question."	10	No	None	None
Doggett et al., February 2024	Special article (Commentary)	ChatGPT-3.5	"Create a 5 answer Single Best Answer question at the level of a [final year/preclinical] medical student on the topic of [topic] where only one answer is correct. Give the answer to this question"	100 pre-clinical, 100 final-year	No	Expert review by six experts (two for preclinical, four for final-year)	"Of 100 preclinical questions, 21 had factual inaccuracies. Of 100 final-year questions, 37 had factual inaccuracies. ... a large proportion of the regenerated questions were near identical. ... Questions were also generally first-order (recall) ..."
E et al., October 2023	Research article	ChatGPT-4	"Could you please generate Multiple Choice Questions (MCQs) on Internal Medicine for medical students? They should be of the same difficulty level as the examples I've provided. Please start numbering from 1, and label the choices from a-d, marking the correct answer with an asterisk. Please write 5 as knowledge questions and 5 with 'clinical history' questions"	210	Two different prompts were used, the description of the process is unclear.	Expert review by five specialists (the number of questions reviewed by each expert is unclear)	"Only 1 question (0.5%) was defined as false; ... approximately 15% of the questions generated from the detailed prompt required some correction, primarily due to inaccuracies in content or faulty methodology."
Eysenbach, March 2023	Special article (Editorial)	ChatGPT-3.5	"Please generate a quiz that asks students to identify the symptoms of diabetes."	5	No	None	None

(Continued)

Table 2. Continued.

Author(s) and date	Publication type	ChatGPT version	Prompt	Number of MCQs	Further modification	Evaluation based on	Outcome
Han et al., October 2023	Research article	ChatGPT-3.5	"Please write multiple-choice questions with a vignette containing lab values"	3	Additional prompts "improve this question by adding more lab values in the stem" and "Can you change this question into one for which familial combined hyperlipidemia is the best answer?" were used.	None	None
Hirosawa and Shimizu, October 2023	Special article (Short communication)	ChatGPT-4	"You are a developer of teaching materials. Create a multiple-choice question with one correct answer regarding the symptoms and physical findings in a patient's medical history that would lead you to suspect anterior cutaneous nerve entrapment syndrome (ACNES). Do not generate the answer yet."	1	No	None	None
Ilgaz and Gelik, September 2023	Research article	ChatGPT-3.5	Unspecified	3	Yes but unspecified	Informal review	"... ChatGPT generated one incorrect question"
Indran et al., December 2023	Special article (Twelve Tips)	ChatGPT-3.5 and ChatGPT-4	"1. Generate a multiple-choice question requiring the learner to choose the single best answer out of 5 options. 2. Provide an answer and rationale. 3. The question should test multiple pharmacology-related concepts using a clinical scenario on nitrofurantoin. Since the prompt is too detailed, we could not fit in this cell. Please refer to the study [24].	2	No	None	None
Kiyak, October 2023	Special article	ChatGPT-3.5	Since the prompt is too detailed, we could not fit in this cell. Please refer to the study [25].	1	No	Informal review	"... it seems the MCQ is both plausible and well-constructed."
Kiyak et al., February 2024	Research article	ChatGPT-3.5	Since the prompt is too detailed, we could not fit in this cell. Please refer to the study [25].	10	No but the change in the patient names for cultural considerations.	Psychometric analysis (item difficulty and discrimination levels, and distractor functionality) in an administration together with 23 human-written questions	All MCQs were considered acceptable in the expert panel. Eight MCQs were not suitable to the context of the medical school. The MCQs had ideal levels of discrimination (0.41 and 0.39), which were above 0.30. Difficulty indices were 0.78 and 0.58. One MCQ had non-functional distractors.
Koga, May 2023	Special article (Letter to Editor)	ChatGPT-4	"Please generate multiple-choice questions formatted in the style of USMLE Step 1, each containing a clinical vignette of about 5 to 6 sentences"	1	No	None	None

(Continued)

