

Human versus Artificial Intelligence: ChatGPT-4 Outperforming Bing, Bard, ChatGPT-3.5 and Humans in Clinical Chemistry Multiple-Choice Questions

Malik Sallam¹⁻³, Khaled Al-Salahat^{1,3}, Huda Eid³, Jan Egger⁴, Behrus Puladi⁵

¹Department of Pathology, Microbiology and Forensic Medicine, School of Medicine, The University of Jordan, Amman, Jordan; ²Department of Clinical Laboratories and Forensic Medicine, Jordan University Hospital, Amman, Jordan; ³Scientific Approaches to Fight Epidemics of Infectious Diseases (SAFE-ID) Research Group, The University of Jordan, Amman, Jordan; ⁴Institute for AI in Medicine (IKIM), University Medicine Essen (AöR), Essen, Germany; ⁵Institute of Medical Informatics, University Hospital RWTH Aachen, Aachen, Germany

Correspondence: Malik Sallam, Email malik.sallam@ju.edu.jo

Introduction: Artificial intelligence (AI) chatbots excel in language understanding and generation. These models can transform healthcare education and practice. However, it is important to assess the performance of such AI models in various topics to highlight its strengths and possible limitations. This study aimed to evaluate the performance of ChatGPT (GPT-3.5 and GPT-4), Bing, and Bard compared to human students at a postgraduate master's level in Medical Laboratory Sciences.

Methods: The study design was based on the METRICS checklist for the design and reporting of AI-based studies in healthcare. The study utilized a dataset of 60 Clinical Chemistry multiple-choice questions (MCQs) initially conceived for assessing 20 MSc students. The revised Bloom's taxonomy was used as the framework for classifying the MCQs into four cognitive categories: Remember, Understand, Analyze, and Apply. A modified version of the CLEAR tool was used for the assessment of the quality of AI-generated content, with Cohen's κ for inter-rater agreement.

Results: Compared to the mean students' score which was 0.68 ± 0.23 , GPT-4 scored 0.90 ± 0.30 , followed by Bing (0.77 ± 0.43), GPT-3.5 (0.73 ± 0.45), and Bard (0.67 ± 0.48). Statistically significant better performance was noted in lower cognitive domains (Remember and Understand) in GPT-3.5 ($P=0.041$), GPT-4 ($P=0.003$), and Bard ($P=0.017$) compared to the higher cognitive domains (Apply and Analyze). The CLEAR scores indicated that ChatGPT-4 performance was "Excellent" compared to the "Above average" performance of ChatGPT-3.5, Bing, and Bard.

Discussion: The findings indicated that ChatGPT-4 excelled in the Clinical Chemistry exam, while ChatGPT-3.5, Bing, and Bard were above average. Given that the MCQs were directed to postgraduate students with a high degree of specialization, the performance of these AI chatbots was remarkable. Due to the risk of academic dishonesty and possible dependence on these AI models, the appropriateness of MCQs as an assessment tool in higher education should be re-evaluated.

Keywords: AI in healthcare education, higher education, large language models, evaluation

Introduction

The domain of higher education is set for a new transformative era.^{1,2} This transformation will be driven by the infiltration of artificial intelligence (AI) into various academic aspects.³⁻⁷ Specifically, the incorporation of AI into higher education can help in enhancing personalized learning, supporting research, automating the grading, facilitating the human-computer interaction, time-saving assistance, and enhancing the students' satisfaction.⁸⁻¹²

Nevertheless, the AI utility in higher education does not only hold promising opportunities but also valid concerns, both of which warrant critical and robust examination.¹³⁻¹⁶ This research endeavor is necessary to guide the ethical, responsible, and productive use of AI to enhance higher education guided by a robust scientific evidence.^{14,16-18} The relevance of the quest to meticulously examine the benefits and challenges of AI in higher education is also important in light of the current evidence showing that a substantial number of university students are already using AI chatbots.¹⁹⁻²⁴

Despite the benefits of AI in higher education, it simultaneously raises valid concerns regarding academic integrity.^{25,26} The ease with which AI can perform complex tasks might inadvertently encourage academic dishonesty, potentially undermining the educational ethics.^{25,27} Furthermore, the reliance on AI for academic tasks could trigger a decline in critical thinking and personal development skills among students, both of which are essential outcomes to enable the graduates in achieving economic, technological, and social advancements.^{28,29} Beyond academic integrity and the fear of declining skills acquired by the students, there are broader societal concerns associated with the pervasive use of AI in higher education.^{30,31} One reasonable fear is job displacement, since AI models can efficiently perform more functions, including teaching and administrative roles traditionally held by humans, with subsequent apprehension about the potential loss of job opportunities within the educational sector and beyond.^{32,33} This could have far-reaching economic implications, affecting not only individuals but also the structure and funding of educational institutions.³⁴

Ultimately, notable capabilities of AI chatbots in understanding and engaging in helpful conversations may contribute to a paradigm shift in higher education.^{8,14,35} This AI-driven change could be a key moment in educational history, with impact surpassing the advent of the internet and the transition to online teaching.^{17,18,36} Therefore, the stakeholders in the academia must strike the right balance between embracing technological innovations while preserving the core values of education.^{37–39} Since the integration of AI into higher education appears inevitable, the academic organizations must adapt to this evolution.³ This adaptation involves the need to emphasize educational aspects such as self-reflection, critical thinking, problem-solving, and independent learning.⁴⁰ Consequently, the educational systems can benefit from AI as a tool to complement, rather than replace, human intellect and creativity.⁴¹

Related Work

In the quest of transition to the AI era in education, guidance by robust scientific evidence is crucial which is highlighted by several studies that addressed the role of generative AI models in healthcare education and practice.^{13,16,42–46} One of the primary steps in this process is to scientifically evaluate the performance of the commonly used and popular AI tools, such as ChatGPT (by OpenAI, San Francisco, CA) with recent introduction of SearchGPT as a prototype of new AI search features, Bing (by Microsoft Corporation, Redmond, WA) its Microsoft Copilot integration, and Bard (by Google, Mountain View, CA), with its latest model Gemini.¹⁶ Several recent studies explored the performance of AI-based models in multiple-choice questions (MCQs), particularly within a broad spectrum of healthcare fields as recently reviewed by Newton and Xiromeriti.⁴⁷ The observed variability in AI performance can be ascribed to several factors, including the different AI models tested, varying approaches to prompting, language variations, and the diversity of the topics tested, among others.^{47–49} Specifically, older or free versions of ChatGPT based on GPT-3 or GPT-3.5 generally underperformed, correctly answering about half of the exam questions, while ChatGPT-4 showed significant improvement, passing most examinations.⁵⁰ Additionally, the underperformance of several generative AI models in various languages compared to its performance in English was demonstrated in recent literature assessing the AI models in Arabic, Polish, Chinese, and Spanish.^{51–55} Thus, continued investigation into this research area is needed to elucidate the determinants of AI model performance across various dimensions which can guide improvements in AI algorithms. However, it is essential that such explorations are conducted utilizing a standardized, refined methodology.^{48,56}

The use of multiple-choice questions (MCQs) have traditionally been fundamental as an objective approach in academic evaluation.⁵⁷ The versatility of MCQs is shown through the use of the Bloom's taxonomy and its subsequent revised framework.^{58,59} The Bloom's taxonomy can guide structuring MCQs to align with specific cognitive functions needed to provide correct answers.⁶⁰ This alignment is key in assessing the students' achievement of the intended learning outcomes.⁶¹ The taxonomy stratifies these cognitive functions into distinct categories. The lower cognitive levels encompass knowledge, which emphasizes "Recall", and comprehension, centered on "Understanding". Conversely, the higher cognitive functions include "Apply", key in problem-solving, and "Analyze", entailing the systematic breakdown of information.^{58,59}

In the context of assessing the performance of AI model performance in MCQs based on the Bloom's taxonomy, a pioneering study by Herrmann-Werner et al assessed ChatGPT-4 with 307 psychosomatic medicine MCQs.^{62,63} The study demonstrated ChatGPT-4 ability to pass the exam irrespective of the prompting method.⁶³ Notably, cognitive errors were more prevalent in "Remember" and "Understand" categories.⁶³ Another recent study demonstrated that ChatGPT-3.5 correctly answered 64 of 80

medical microbiology MCQs, albeit below student averages, with better performance in the “Remember” and “Understand” categories and more frequent errors in MCQ with longer choices in terms of word count.⁶⁴

In this study, the objective was to synthesize and expand upon recent research examining the performance of AI chatbots in various examinations.⁵⁰ This research was informed by seminal studies, such as the Kung et al evaluation of ChatGPT in the United States Medical Licensing Examination (USMLE),⁴⁹ and also it aimed to extend the evidence of AI chatbot performance in a topic rarely encountered in literature, namely the Clinical Chemistry at postgraduate level. The novel contribution of the current study lies in employing a standardized framework, termed “METRICS” for the design and reporting of AI assessment studies, coupled with an in-depth analysis of AI models’ rationale behind responses, using an evaluation tool specifically tailored for AI content evaluation referred to as “CLEAR”.^{48,56}

The study hypothesized that postgraduate students, particularly in the field of clinical chemistry, may demonstrate superior performance compared to AI models. We anticipate that this disparity may be especially evident in tasks requiring higher cognitive functions, such as “Apply” and “Analyze”. This study aimed to critically assess the current capabilities of AI in an academic setting and explore the differences of human versus artificial intelligence in complex problem-solving scenarios.

Material and Methods

Study Design

The study utilized the METRICS checklist for the design and reporting of AI studies in healthcare.⁴⁸ The basis of the study was a dataset of 60 MCQs, used in a Clinical Chemistry examination. This examination was part of the Medical Laboratory Sciences Clinical Chemistry course, tailored for Master of Science (MSc) students in Medical Laboratory Sciences at the School of Medicine, University of Jordan.

The specific exam in focus was conducted in-person and 20 students undertook the examination during the Autumn Semester of the 2019/2020 academic year. The students’ performance in each question was available for comparison with AI models.

The MCQs utilized in this exam were designed by the first author (M.S.), who is a Jordan Medical Council (JMC) certified consultant in Clinical Pathology. Additionally, the first author (M.S.) has been a dedicated instructor for this course since the Academic Year 2018/2019. The MCQs were original, ensuring there were no copyright concerns.

Ethical Considerations

This study was waived from IRB approval from the Deanship of Scientific Research at the University of Jordan (Reference number: 2/2024/19). Ethical clearance for this research was determined to be non-essential, given the nature of the data involved. The data utilized were entirely anonymized, ensuring no breach of confidentiality or personal privacy. Additionally, the university examination results, which formed part of our dataset, are publicly accessible and open for academic scrutiny. Moreover, the MCQs employed in the study were originally created by the first author. These questions are devoid of any copyright concerns, further reinforcing the ethical integrity of our research approach.

The nature of this study involved the analysis of anonymized data and publicly accessible university examination results. Given that the data were entirely anonymized and did not involve any direct interaction with participants or any intervention, informed consent was deemed unnecessary. This decision was confirmed by the waiver of IRB approval from the Deanship of Scientific Research at the University of Jordan (Reference number: 2/2024/19).

Prior to generative AI tools’ use in this study, a thorough review of the terms of use was conducted to comply with all applicable guidelines. Additionally, the use of generative AI tools in this study was done in agreement with the standard licensing agreements and no additional permissions were required.

MCQ Features and Indices of Human Students’ Performance

The indices of student performance included facility index defined as the proportion of students who correctly answered the MCQ divided by the total number of students ($n=20$). The paper-based exam in a classroom with answer sheets for answering was administered in December 2019 and was designed to last for a duration of 90 minutes, with an average

allocation of 1.5 minutes per question. Completing all questions was mandatory for candidates to be eligible for full marks. The scoring process was automated by the School of Medicine at the University of Jordan, utilizing an answer key sheet. The students were then divided into the upper group comprising the top 5 performing students, and the middle group comprising the middle 10 students and the lower group comprising the lower scoring 5 students. The “Discrimination Index” (DI) was then calculated based on the difference between the percent of correct responses in the upper group and the percent of correct responses in the lower group. This was followed by the calculation of the “Maximum Discrimination” based on the sum of the percent in the upper and lower groups marking the item correctly. Then, the Discrimination Efficiency (DE) of the MCQ was calculated as the ratio of DI to the Maximum Discrimination. The classification of the MCQs based on the revised Bloom’s taxonomy four cognitive levels “Remember”, “Understand”, “Apply”, and “Analyze” was based on a consensus between the first and second authors, both of which are certified Clinical Pathologists.

Models of AI Tested, Settings, Testing Time, and Duration

In this study, a detailed evaluation of four AI models was conducted, each selected for its relevance, popularity, and advanced capabilities in language processing as follows: First, ChatGPT-3.5 (OpenAI, San Francisco, CA):⁶⁵ This model is grounded in the GPT-3.5 architecture deployed using its default settings and was assessed as of its latest update at time of testing as of January 2022.

Second, ChatGPT-4 (OpenAI, San Francisco, CA):⁶⁵ An advancement in the Generative Pre-trained Transformer (GPT) series, with the most recent update from April 2023 at time of testing. Third, Bing Chat (GPT-4 Turbo):⁶⁶ This model uses the GPT-4 Turbo model. At the time of testing, the version was updated until April 2023 and we selected the more balanced conversation style. Fourth, Bard (Google, Mountain View, CA):⁶⁷ This Google AI GPT model was last updated on October 4, 2023, at time of testing.

The testing of these models was conducted over a concise period, spanning November 27 to November 28, 2023. Our methodological approach involved initiating interactions with GPT-3.5, GPT-4, and Bard using a single page. For Bing Chat, we used the “New Topic” option considering the limit of responses posed by this model (50 at maximum). Additionally, we opted not to use the “regenerate response” feature in ChatGPT and abstained from providing feedback in all models to avoid feedback bias.

Prompt and Language Specificity

In this study, we meticulously crafted the prompts used for interacting with the AI models to ensure clarity and consistency in the testing process. For ChatGPT-3.5, ChatGPT-4, and Bard, the following exact prompt was used:

For the following 60 Clinical Chemistry MCQs that will be provided one by one, please select the most appropriate answer for each MCQ, with an explanation for the rationale behind selecting this choice and excluding the other choices. Please note that only one choice is correct while the other four choices are incorrect. Please note that these questions were designed for masters students in medical laboratory sciences.

This was followed by prompting each MCQ one by one. For Bing, the following prompt was used for each MCQ:

For the following 60 Clinical Chemistry MCQs that will be provided one by one, please select the most appropriate answer for each MCQ, with an explanation for the rationale behind selecting this choice and excluding the other choices. Please note that only one choice is correct while the other four choices are incorrect. Please note that these questions were designed for masters students in medical laboratory sciences.

All MCQs were presented in English. This choice was based on the fact that English is the official language of instruction for the MSc program in Medical Laboratory Sciences at the University of Jordan.

AI Content Evaluation Approach and Individual Involvement in Evaluation

First, we objectively assessed the correctness of responses based on the key answers of the MCQs. Then, subjective evaluation of the AI generated content was based on a modified version of the CLEAR tool. This involved assessing the

content on three dimensions as follows: First, completeness of the generated response. Second, accuracy reflected by lack of false knowledge and the content being evidence-based. Third, appropriateness and relevance of content being easy to understand, well organized, and free from irrelevant content.⁵⁶ Each dimension was scored on a 5-point Likert scale ranging from 1 = poor, 2 = satisfactory, 3 = good, 4 = very good, to 5 = excellent. A list of the key points to be considered in the assessment was set beforehand to increase objectivity.

The content generated by the four models was evaluated by two raters independently; the first author (M.S.) a consultant in Clinical Pathology, and the second author (K.A.) a specialist in Clinical Pathology, both certified in Clinical Pathology from the Jordan Medical Council (JMC).

Data Source Transparency and Topic Range

The MCQs were totally conceived by the first author and sole instructor of the course. Sources of the material taught during the course were the following three textbooks: Tietz Textbook of Clinical Chemistry and Molecular Diagnostics; Clinical Chemistry: Principles, Techniques, and Correlations; and Henry's Clinical Diagnosis and Management by Laboratory Methods.^{68–70}

The scope of topics covered in the MCQs were as follows: Adrenal Function, Amino Acids and Proteins, Body Fluid Analysis, Clinical Enzymology, Electrolytes, Gastrointestinal Function, Gonadal Function, Liver Function, Nutrition Assessment, Pancreatic Function, Pituitary Function, Thyroid Gland, and Trace Elements.

Statistical and Data Analyses

The statistical analysis was conducted using IBM SPSS Statistics Version 26.0 (Armonk, NY: IBM Corp). The continuous variables were presented as means and standard deviations (SD), while categorical data were summarized as frequencies and percentages [N (%)]. To explore the associations between categorical variables, we employed the chi-squared test (χ^2), while to explore the associations between scale variables and categorical variables, non-parametric tests were utilized: the Mann–Whitney *U*-test (M-W) and the Kruskal Wallis test (K-W). The Kolmogorov–Smirnov test was employed to confirm the non-normality of the scale variables: facility index (FI, $P=0.042$), discriminative efficiency (DE, $P=0.011$), word count for both stem and choices ($P<0.001$ for both), average completeness, accuracy/evidence, appropriateness/relevance, and the mCLEAR scores ($P<0.001$ for the four scores). P values <0.050 were considered statistically significant. For multiple comparisons, post hoc analysis was conducted using the *M-W* test. To account type I error due to multiple comparisons, we adjusted the α level using the Bonferroni correction. Consequently, the adjusted α level for conducting pairwise comparisons between the four AI models was set at $P=0.0083$.

The MCQs were categorized based on the FI as “difficult” for an FI of 0.40 or less, “average” for an FI > 0.40 and ≤ 0.80 , and “easy” for an FI > 0.80 . Additionally, the DE was stratified into “poor discrimination” if the DE was between -1 to zero, “satisfactory discrimination” for DE > 0 to < 0.40 as satisfactory, and DE ≥ 0.4 indicating “good discrimination”.

The inter-rater agreement was assessed using Cohen's kappa (κ) values, which ranged from very good to excellent. For ChatGPT-3.5, the agreement was $\kappa=0.874$ for Completeness, $\kappa=0.921$ for Accuracy, and $\kappa=0.723$ for Relevance. ChatGPT-4 showed $\kappa=0.845$ for Completeness, a perfect $\kappa=1$ for Accuracy, and $\kappa=0.731$ for Relevance. Bing displayed κ values of 0.911 for Completeness, 0.871 for Accuracy, and 0.840 for Relevance. Lastly, Bard's agreement was $\kappa=0.903$ for Completeness, $\kappa=1$ for Accuracy, and $\kappa=0.693$ for Relevance. Finally, the overall modified CLEAR (mCLEAR) scores for AI content quality were averaged based on the scores of the two raters and categorized as: “Poor” (1–1.79), “Below average” (1.80–2.59), “Average” (2.60–3.39), “Above average” (3.40–4.19), and “Excellent” (4.20–5.00) similar to the previous approach in.⁷¹

Results

Overall Performance of the Tested AI Models Compared to the Human Students

The overall performance of the MSc students in the exam was reflected in the average score of 40.05 ± 7.23 (66.75%), with the range of scores of the students of 24–54 (40.00%–90.00%). The performance of the four AI models varied with

the best performance for ChatGPT-4 scoring 54/60 (0.90 ± 0.30), followed by Bing scoring 46/60 (0.77 ± 0.43), ChatGPT-3.5 scoring 44/60 (0.73 ± 0.45), and finally Bard scoring 40/60 (0.67 ± 0.48).

Human Students' Performance Based on the Revised Bloom's Taxonomy

The MCQ metrics were derived from the performance of the 20 MSc students in the exam. The best performance was in the "Remember" category, followed by the "Apply" category, "Understand" category, while the worst performance was in the "Analyze" category; however, these differences lacked statistical significance (Table 1).

Performance of the AI Models Based on the MCQ Metrics

The performance of the four tested AI models was stratified based on the MCQ metrics. Significantly lower number of correct answers was seen in difficult MCQs in both Bing (44% correct answers for the difficult MCQs as opposed to 84% and 81% for the easy and average MCQs, respectively, $P=0.045$, $\chi^2=6.204$) and Bard (44% correct answers for the difficult MCQs as opposed to 90% and 59% for the easy and average MCQs, respectively, $P=0.027$, $\chi^2=7.213$) (Table 2), while the MCQ stem and choices word counts were not associated with AI models' performance.

Performance of the AI Models Based on the Revised Bloom's Taxonomy

Upon analyzing the AI models' performance in MCQs stratified per the four revised Bloom's categories, only ChatGPT-4 showed statistically significant better performance in the Remember and Understand categories compared to Apply and Analyze categories (Table 3).

On the other hand, ChatGPT-3.5, ChatGPT-4, and Bard showed statistically better performance in the lower cognitive MCQs compared to the higher cognitive MCQs (Figure 1). Specifically, GPT-3.5 correctly answered 81% of the lower cognitive MCQs compared to 56% correct answers in the higher cognitive MCQs ($P=0.041$, $\chi^2=4.156$); GPT-4 correctly answered 98% of the lower cognitive MCQs compared to 72% correct answers in the higher cognitive MCQs ($P=0.003$, $\chi^2=9.030$); and Bard correctly answered 76% of the lower cognitive MCQs compared to 44% correct answers in the higher cognitive MCQs ($P=0.017$, $\chi^2=5.714$).

Table 1 Multiple-Choice Questions (MCQs) Metrics Stratified by the Revised Bloom's Taxonomy as Derived from the Performance of 20 MSc Students

| Revised Bloom's taxonomy | Remember | Understand | Apply | Analyze | P value ^c |
|--|----------------------------|-------------------|-------------------|-------------------|----------------------|
| MCQ metric | Mean \pm SD ^b | Mean \pm SD | Mean \pm SD | Mean \pm SD | |
| Facility index | 0.74 \pm 0.22 | 0.61 \pm 0.28 | 0.71 \pm 0.21 | 0.6 \pm 0.17 | 0.180 |
| Discriminative efficiency | 0.24 \pm 0.25 | 0.24 \pm 0.27 | 0.17 \pm 0.43 | 0.43 \pm 0.41 | 0.482 |
| MCQ stem word count | 15.04 \pm 6.5 | 24 \pm 16.95 | 73.4 \pm 57.06 | 25.31 \pm 27.55 | 0.052 |
| MCQ choices word count | 13.5 \pm 9.29 | 24.83 \pm 17.76 | 22.2 \pm 8.07 | 29.54 \pm 28.64 | 0.153 |
| Revised Bloom's cognitive level ^a | Lower | | Higher | | P value ^d |
| Facility index | 0.68 \pm 0.25 | | 0.63 \pm 0.18 | | 0.225 |
| Discriminative efficiency | 0.24 \pm 0.25 | | 0.36 \pm 0.42 | | 0.205 |
| MCQ stem word count | 18.88 \pm 12.77 | | 38.67 \pm 42.34 | | 0.265 |
| MCQ choices word count | 18.36 \pm 14.54 | | 27.5 \pm 24.61 | | 0.268 |

Notes: ^aLower cognitive level includes "Remember" and "Understand" categories, while the higher cognitive level includes "Apply" and "Analyze" categories; ^bSD: Standard deviation; ^cCalculated using the Kruskal Wallis test; ^dCalculated using the Mann Whitney U-test.

Table 2 Artificial Intelligence (AI)-Based Model Performance Based on the Multiple-Choice Question (MCQ) Metrics

| AI model | Answer | FI ^a category | | | | DE ^c category | | | | MCQ stem word count | | MCQ choices word count | |
|----------|-----------|--------------------------|-----------|-----------|-------------------|--------------------------|--------------|-----------|-------------------|----------------------------|----------------------|------------------------|----------------------|
| | | Easy | Average | Difficult | P value, χ^2 | Poor | Satisfactory | Good | P value, χ^2 | Mean \pm SD ^d | P value ^e | Mean \pm SD | P value ^e |
| | | N ^b (%) | N (%) | N (%) | | N (%) | N (%) | N (%) | | | | | |
| GPT-3.5 | Correct | 15 (78.9) | 22 (68.8) | 7 (77.8) | 0.690, 0.741 | 10 (62.5) | 15 (65.2) | 19 (90.5) | 0.087, 4.891 | 23.8 \pm 28.49 | 0.063 | 18.2 \pm 16.07 | 0.055 |
| | Incorrect | 4 (21.1) | 10 (31.3) | 2 (22.2) | | 6 (37.5) | 8 (34.8) | 2 (9.5) | | 27.63 \pm 21.62 | | 29.06 \pm 22.41 | |
| GPT-4 | Correct | 19 (100) | 26 (81.3) | 9 (100) | 0.054, 5.833 | 14 (87.5) | 21 (91.3) | 19 (90.5) | 0.923, 0.160 | 24.33 \pm 27.49 | 0.315 | 20.5 \pm 18.67 | 0.339 |
| | Incorrect | 0 | 6 (18.8) | 0 | | 2 (12.5) | 2 (8.7) | 2 (9.5) | | 29.17 \pm 19.57 | | 26.5 \pm 16.49 | |
| Bing | Correct | 16 (84.2) | 26 (81.3) | 4 (44.4) | 0.045, 6.204 | 11 (68.8) | 18 (78.3) | 17 (81.0) | 0.667, 0.809 | 25.89 \pm 29.62 | 0.322 | 19.57 \pm 15.17 | 0.655 |
| | Incorrect | 3 (15.8) | 6 (18.8) | 5 (55.6) | | 5 (31.3) | 5 (21.7) | 4 (19.0) | | 21.29 \pm 13.53 | | 26.14 \pm 26.6 | |
| Bard | Correct | 17 (89.5) | 19 (59.4) | 4 (44.4) | 0.027, 7.213 | 9 (56.3) | 18 (78.3) | 13 (61.9) | 0.303, 2.387 | 26.9 \pm 31.18 | 0.660 | 17.63 \pm 14.08 | 0.114 |
| | Incorrect | 2 (10.5) | 13 (40.6) | 5 (55.6) | | 7 (43.8) | 5 (21.7) | 8 (38.1) | | 20.65 \pm 13.9 | | 28.05 \pm 23.89 | |

Notes: ^aFI: Facility index of the MCQ; ^bN: Number; ^cDE: Discriminative efficiency of the MCQ; ^dSD: Standard deviation; ^eCalculated using the Mann Whitney U-test.

Table 3 The Performance of the Four Artificial Intelligence (AI)-Based Models in the Clinical Chemistry Multiple-Choice Question (MCQs) Stratified per the Four Revised Bloom's Categories

| Revised Bloom's taxonomy | Answer | Remember | Understand | Apply | Analyze | P value, χ^2 |
|--------------------------|-----------|--------------------|------------|----------|-----------|-------------------|
| | | N ^a (%) | N (%) | N (%) | N (%) | |
| GPT-3.5 | Correct | 19 (79.2) | 15 (83.3) | 2 (40.0) | 8 (61.5) | 0.164, 5.104 |
| | Incorrect | 5 (20.8) | 3 (16.7) | 3 (60.0) | 5 (38.5) | |
| GPT-4 | Correct | 24 (100) | 17 (94.4) | 3 (60.0) | 10 (76.9) | 0.015, 10.532 |
| | Incorrect | 0 | 1 (5.6) | 2 (40.0) | 3 (23.1) | |
| Bing | Correct | 20 (83.3) | 15 (83.3) | 4 (80.0) | 7 (53.8) | 0.182, 4.859 |
| | Incorrect | 4 (16.7) | 3 (16.7) | 1 (20.0) | 6 (46.2) | |
| Bard | Correct | 18 (75.0) | 14 (77.8) | 3 (60.0) | 5 (38.5) | 0.090, 6.504 |
| | Incorrect | 6 (25.0) | 4 (22.2) | 2 (40.0) | 8 (61.5) | |

Notes: ^aN: Number.

Performance of the AI Models Based on the Modified CLEAR Tool

In our assessment of completeness, accuracy/evidence, and appropriateness/relevance, based on the modified CLEAR tool, ChatGPT-4 was the only model rated as “Excellent” across all categories. Bing achieved an “Excellent” rating solely in appropriateness/relevance. The other AI models were categorized as “Above average” in performance (Table 4). The statistical analysis revealed significant superiority of ChatGPT-4 compared to the other models in all CLEAR categories, with the exception of Bing where the difference was only significant in the completeness and the overall mCLEAR score (Table 4).

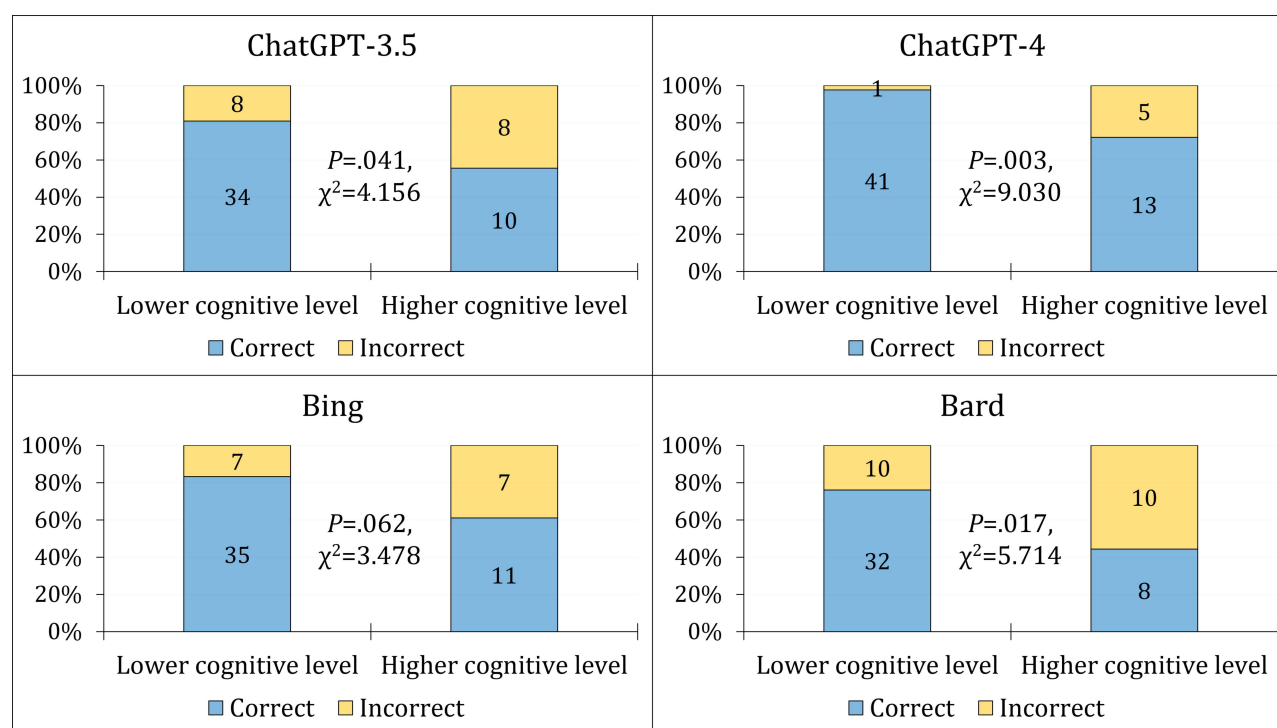
**Figure 1** The performance of the four artificial intelligence (AI)-based models in the MCQs stratified per the revised Bloom cognitive levels.

Table 4 Modified CLEAR Average Scores for the Four AI Models in Explaining the Rationale for Selecting Choices

| Assessment category | Mean±SD ^b | Rank | P value ^c | Post hoc test (Mann Whitney U-test) | | | | | |
|---|----------------------|---------------|----------------------|-------------------------------------|-----------------|-----------------|---------------|---------------|--------------|
| | | | | GPT-3.5 vs GPT-4 | GPT-3.5 vs Bing | GPT-3.5 vs Bard | GPT-4 vs Bing | GPT-4 vs Bard | Bing vs Bard |
| ChatGPT-3.5 completeness score | 4.03±1.26 | Above average | <0.001 | <0.001 | 0.343 | 0.249 | 0.001 | 0.003 | 0.745 |
| ChatGPT-4 completeness score | 4.73±0.77 | Excellent | | | | | | | |
| Bing completeness score | 4.14±1.34 | Above average | | | | | | | |
| Bard completeness score | 4.19±1.18 | Above average | | | | | | | |
| ChatGPT-3.5 accuracy/evidence score | 3.87±1.80 | Above average | 0.016 | 0.007 | 0.633 | 0.604 | 0.023 | 0.002 | 0.324 |
| ChatGPT-4 accuracy/evidence score | 4.6±1.21 | Excellent | | | | | | | |
| Bing accuracy/evidence score | 4.07±1.66 | Above average | | | | | | | |
| Bard accuracy/evidence score | 3.67±1.90 | Above average | | | | | | | |
| ChatGPT-3.5 appropriate/relevance score | 4.18±1.33 | Above average | 0.011 | 0.005 | 0.645 | 0.691 | 0.023 | 0.001 | 0.355 |
| ChatGPT-4 appropriate/relevance score | 4.76±0.76 | Excellent | | | | | | | |
| Bing appropriate/relevance score | 4.27±1.35 | Excellent | | | | | | | |
| Bard appropriate/relevance score | 4.15±1.22 | Above average | | | | | | | |
| ChatGPT-3.5 mCLEAR ^a score | 4.03±1.41 | Above average | <0.001 | <0.001 | 0.270 | 0.213 | 0.001 | 0.002 | 0.868 |
| ChatGPT-4 mCLEAR score | 4.70±0.90 | Excellent | | | | | | | |
| Bing mCLEAR score | 4.16±1.43 | Above average | | | | | | | |
| Bard mCLEAR score | 4.00±1.41 | Above average | | | | | | | |

Notes: ^amCLEAR: Modified CLEAR score based on the study by Sallam et al; ^bSD: Standard deviation; ^cCalculated using the Kruskal Wallis test. Significant P values are highlighted in bold style.

Discussion

The whole landscape of education, including higher education is set for a new era that can be described as a paradigm shift with the widespread and popularity of generative AI.^{13,72,73} In this study, a comparison between the human and AI abilities in a highly specialized field at a high level was undertaken. Specifically, the performance of MSc students in a Clinical Chemistry exam, with an average score of 40.05 ± 7.23 (66.75%), was used as a benchmark for comparison. Remarkably, ChatGPT-4 surpassed this human benchmark, achieving a score of 54/60 (90.00%). Bing followed with 46/60 (76.67%), outperforming both ChatGPT-3.5 (44/60, 73.33%) and Bard (40/60, 66.67%). Overall, the level of AI models' performance underlines the advancements in AI capabilities. Additionally, these results could pave the way for a broader scientific inquiry into both the potential role of AI in educational settings as well as the usefulness of the current assessment tools in higher education.

In this study, the initial central hypothesis assumed that the human students at a postgraduate level who undertook a specialized course in a highly specialized field, namely Clinical Chemistry, would show a superior performance compared to the tested AI models. The findings of this study showed that the AI models tested not only passed the exam but showed a noteworthy performance. For example, ChatGPT-4 score equaled the highest student score and thus would be rated as an "A" student. On the other hand, the performance of the AI models in this study was not entirely an unexpected finding. This comes in light of the recent evidence showing AI models' abilities to pass reputable exams in multiple languages such as the USMLE,⁴⁹ the German State Examination in Medicine,⁷⁴ the National Medical Licensing Examination in Japan,^{75,76} and the Brazilian National Examination for Medical Degree Revalidation.⁷⁷

From a broader perspective, a recent systematic review highlighted the abilities of ChatGPT as an example of LLMs in various exams.^{47,50} The review by Newton and Xiromeriti highlighted the capabilities of this popular AI model, with ChatGPT-3 outperforming human students in 11% of the included exams, with ChatGPT 4 achieving superior performance and outscoring the human performance in 35% of the included exams.⁴⁷ The current study findings were in line with the finding of better GPT-4 performance as opposed to the earlier and free GPT-3.5 version. Yet, the performance of ChatGPT-4 in comparison to the human students was noteworthy highlighting the refinements of LLMs over a short period of time.

In this study, analyzing the human students' performance based on the revised Bloom's taxonomy enabled elucidation of deeper insights into the assessment of cognitive aspects. The human students excelled in the "Remember" domain which is indicative of strong recalling and recognizing abilities. Additionally, the human students demonstrated a high performance in the "Understand" and "Apply" categories, with the lowest performance shown in the "Analyze" category. The lack of statistical significance in these differences suggest a balanced level of cognitive skills acquired among the students during the course despite the potential for improvement in higher-order cognitive skills entailing breakdown and organization of acquired knowledge.

On the other hand, the study findings revealed an interesting observation manifested in worse AI models' performance across the higher cognitive domains. This observation stands in contrast to the findings of Herrmann-Werner et al, which pioneered the use of the Bloom's taxonomy in AI model performance in MCQs.⁶² Herrmann-Werner et al demonstrated a lower level of ChatGPT performance in the lower cognitive skills in contrast to the findings of this study.⁶² To the contrary, a recent study that assessed ChatGPT-3 performance in medical microbiology MCQs showed a trend similar to our findings where the AI model performed at a higher level in the lower cognitive domains.⁶⁴ This divergence of findings suggests the need for more comprehensive studies to discern the abilities of AI models in different cognitive domains, which would be helpful to guide improvements in these models and to enhance their utility in higher education.

Upon examining the performance of AI models in this study based on the MCQ metrics (FI, DE, stem and choices word count), a significant drop in performance was noted in Bing and Bard for more difficult MCQs. This finding suggests that some AI models have yet to show evolution into the level where it can handle complex queries. The absence of a correlation between MCQ stem and choice word counts and AI performance indicates that the challenge was not related to the length of the queries but rather in the inherent complexity of the prompts.

In this study, the use of the validated CLEAR tool for assessment of the quality of AI generated content presented a robust approach.⁵⁶ The rating of ChatGPT-4 as “Excellent” across all categories of completeness, accuracy/evidence, and appropriateness/relevance serves as a clear demonstration of its superiority. The Bing’s —which uses similar GPT-4 architecture— rating as “Excellent” in appropriateness/relevance was a noteworthy finding; nevertheless, the performance of this Microsoft AI model did not match ChatGPT-4 in terms of completeness and accuracy. The other AI models in this study were rated as “Above average” based on the modified CLEAR tool. This result, albeit lower than ChatGPT-4, still showed the huge potential of these freely available models, but with an evident room for improvement. The significant superiority of ChatGPT-4 over the other AI models tested in this study highlights the swift evolution of AI capabilities.⁷⁸

In the field of higher education, the implications of the study findings can be profound. The noteworthy capabilities of AI models, especially those shown by ChatGPT-4, to outperform humans at a postgraduate level could serve as a red flag necessitating the re-evaluation of traditional assessment approaches currently utilized for evaluation of students’ achievement of learning outcomes.^{72,79} Additionally, the study findings highlighted the current possible AI limitations in addressing higher-order cognitive tasks, which shows the unique value of human critical thinking and analytical skills.⁸⁰ Nevertheless, more studies are needed to confirm this finding based on a recent evidence showing the satisfactory performance of ChatGPT in tasks requiring higher-order thinking specifically in the field of medical biochemistry as shown by Ghosh and Bir.⁸¹

Future research could focus on investigating the feasibility of integrating AI into higher education frameworks in terms of utilizing an approach that could augment the human learning (eg, through enhancing personalized learning experience and providing instantaneous feedback) without compromising the development of critical thinking and analytical skills.^{5,9,72,82–84} Additionally, the ethical considerations of academic integrity should be considered in light of opportunities of academic dishonesty posed by AI models in educational settings.^{85–87} This issue also extends to warrant a thorough investigation into the implications of possible decline in students’ analytical and critical thinking skills and prioritizing the human needs and value.^{29,88,89}

Finally, while the current study can provide valuable insights into the performance of AI models compared to human students in the context of Clinical Chemistry topic, several limitations should be considered when interpreting the results. Future research in this area would benefit from addressing these limitations that included: First, this study employed a limited dataset of 60 MCQs. This limited number of MCQs inherently restricts the scope of performance evaluation. Second, the use of the CLEAR tool, albeit standardized, introduces a subjective element in evaluating the content generated by AI models. This subjectivity could lead to a potential bias in the assessment of AI responses if approached by different raters. Thus, the AI content evaluation was not entirely devoid of subjective judgment despite the use of key answers to reduce this subjectivity bias. Third, the exclusive concentration on Clinical Chemistry as a subject is both a strength and a limitation. While it allowed for a deep insight this specific health discipline, it limits the generalizability of the findings to other academic fields, since different subjects may present unique challenges that were not addressed in this study. Fourth, LLMs are evolving rapidly, and this study only provided a snapshot of AI models’ performance at a specific time point. Therefore, this study may not fully represent the potential improvements or advancements in AI capabilities that have occurred or may occur shortly after the study period. Fifth, the exam metrics, derived from the performance of a limited number of students (n=20), might have been influenced by various external factors. These include the format of the exam and its time limits and the specific cohort of students. Finally, the study results was based on prompting the AI-based models in English, which may also limit the generalizability of results based on varying levels of performance of AI models based on languages used.^{90,91}

Conclusion

The current study provided a comparative analysis of the human versus AI performance in a highly specialized academic context at the postgraduate level. The results could motivate future research to address the possible role of AI in higher education reaping its benefits while avoiding its limitations. The ideal approach would be to use the strengths of AI as a complement to the unique capabilities of human intellect. This can ensure the evolution of the educational process in an innovative way aiding in students’ intellectual development. Importantly, the study results call for a revision of the current assessment tools in higher education with a focus on improving the assessment of higher cognitive skills.

Data Sharing Statement

The data that support the findings of this study are available on request from the corresponding author (M.S.). The data are not publicly available due to the confidentiality of the questions created for an exam purpose.

Ethics Approval and Consent to Participate

This study was waived from IRB approval from the Deanship of Scientific Research at the University of Jordan (Reference number: 2/2024/19).

Acknowledgment

This paper has been uploaded to Research Square and medRxiv as a preprint which can be accessed through the following links:

Research Square: <https://doi.org/10.21203/rs.3.rs-3880412/v1>

medRxiv: <https://doi.org/10.1101/2024.01.08.24300995>

Author Contributions

All authors made a significant contribution to the work reported, whether that is in the conception, study design, execution, acquisition of data, analysis and interpretation, or in all these areas; took part in drafting, revising or critically reviewing the article; gave final approval of the version to be published; have agreed on the journal to which the article has been submitted; and agree to be accountable for all aspects of the work.

Funding

We declare that we received no funding nor financial support/grants by any institutional, private, or corporate entity.

Disclosure

The authors report no conflicts of interest in this work.

References

1. Chiu TKF. Future research recommendations for transforming higher education with generative AI. *Comp Educat*. 100197. doi:10.1016/j.caeai.2023.100197
2. Rawas S. ChatGPT: empowering lifelong learning in the digital age of higher education. *Educat Inform Technol*. 2023. doi:10.1007/s10639-023-12114-8
3. Rahiman HU, Kodikal R. Revolutionizing education: artificial intelligence empowered learning in higher education. *Cogent Educat*. 2024;11:2293431. doi:10.1080/2331186X.2023.2293431
4. Crompton H, Burke D. Artificial intelligence in higher education: the state of the field. *Int J Educa Technol High Educ*. 2023;20:22. doi:10.1186/s41239-023-00392-8
5. Karabacak M, Ozkara BB, Margetis K, Wintermark M, Bisdas S. The Advent of Generative Language Models in Medical Education. *JMIR Med Educ*. 2023;9:e48163. doi:10.2196/48163
6. Rodway P, Schepman A. The impact of adopting AI educational technologies on projected course satisfaction in university students. *Comput Educat*. 2023;5:100150. doi:10.1016/j.caeai.2023.100150
7. Giansanti D. The Chatbots Are Invading Us: a Map Point on the Evolution, Applications, Opportunities, and Emerging Problems in the Health Domain. *Life*. 2023;13:1130. doi:10.3390/life13051130
8. Dempere J, Modugu K, Hesham A, Ramasamy LK. The impact of ChatGPT on higher education. *Frontiers in Education*. 2023;8:1206936. doi:10.3389/educ.2023.1206936
9. Sallam M, Salim NA, Barakat M, Al-Tammemi AB. ChatGPT applications in medical, dental, pharmacy, and public health education: a descriptive study highlighting the advantages and limitations. *Narra J*. 2023;3:e103. doi:10.52225/narra.v3i1.103
10. Sáiz-Manzanares MC, Marticorena-Sánchez R, Martín-Antón LJ, González Díez I, Almeida L. Perceived satisfaction of university students with the use of chatbots as a tool for self-regulated learning. *Heliyon*. 2023;9:e12843. doi:10.1016/j.heliyon.2023.e12843
11. Labadze L, Grigolia M, Machaidze L. Role of AI chatbots in education: systematic literature review. *Int J Educa Technol High Educ*. 2023;20:56. doi:10.1186/s41239-023-00426-1
12. Imran M, Almusharraf N. Analyzing the role of ChatGPT as a writing assistant at higher education level: a systematic review of the literature. *Contemporary Educational Technology*. 2023;15:ep464. doi:10.30935/cedtech/13605
13. 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
14. Kooli C. Chatbots in Education and Research: a Critical Examination of Ethical Implications and Solutions. *Sustainability*. 2023;15:5614. doi:10.3390/su15075614

15. Ray PP. ChatGPT: a comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Inter Thin Cyber-Physi Syst.* **2023**;3:121–154. doi:10.1016/j.iotcps.2023.04.003
16. Sallam M, Al-Farajat A, Egger J. Envisioning the Future of ChatGPT in Healthcare: insights and Recommendations from a Systematic Identification of Influential Research and a Call for Papers. *Jordan Medical Journal.* **2024**;58. doi:10.35516/jmj.v58i1.2285
17. Grassini S. Shaping the Future of Education: exploring the Potential and Consequences of AI and ChatGPT in Educational Settings. *Education Sciences.* **2023**;13:692. doi:10.3390/educsci13070692
18. Kamalov F, Santandreu Calonge D, Gurrib I. New Era of Artificial Intelligence in Education: towards a Sustainable Multifaceted Revolution. *Sustainability.* **2023**;15:12451. doi:10.3390/su151612451
19. von Garrel J, Mayer J. Artificial Intelligence in studies—use of ChatGPT and AI-based tools among students in Germany. *Humanit Soc Sci Commun.* **2023**;10:799. doi:10.1057/s41599-023-02304-7
20. Sallam M, Salim NA, Barakat M, et al. Assessing Health Students' Attitudes and Usage of ChatGPT in Jordan: validation Study. *JMIR Med Educ.* **2023**;9:e48254. doi:10.2196/48254
21. Malik AR, Pratiwi Y, Andajani K, Numertayasa IW, Suharti S, Darwis A. Exploring Artificial Intelligence in Academic Essay: higher Education Student's Perspective. *Interl J Educat Re Open.* **2023**;5:100296. doi:10.1016/j.ijedro.2023.100296
22. Rodríguez JMR, Montoya MSR, Fernández MB, Lara FL. Use of ChatGPT at university as a tool for complex thinking: students' perceived usefulness. *NAER.* **2023**;12:323–339. doi:10.7821/naer.2023.7.1458
23. Abdaljalael M, Barakat M, Alsanafi M, et al. A multinational study on the factors influencing university students' attitudes and usage of ChatGPT. *Sci Rep.* **2024**;14. doi:10.1038/s41598-024-52549-8
24. Sallam M, Elsayed W, Al-Shorbagy M, et al. ChatGPT usage and attitudes are driven by perceptions of usefulness, ease of use, risks, and psycho-social impact: a study among university students in the UAE. *Frontiers in Education.* **2024**;9:1414758. doi:10.3389/educ.2024.1414758
25. Cotton DRE, Cotton PA, Shipway JR. Chatting and cheating: ensuring academic integrity in the era of ChatGPT. *In Educ Teach Inter.* **2023**:1–12. doi:10.1080/14703297.2023.2190148
26. Bin-Nashwan SA, Sadallah M, Bouteraa M. Use of ChatGPT in academia: academic integrity hangs in the balance. *Technol Soc.* **2023**;75:102370. doi:10.1016/j.techsoc.2023.102370
27. Birks D, Clare J. Linking artificial intelligence facilitated academic misconduct to existing prevention frameworks. *Interl J Educat Integrity.* **2023**;19:20. doi:10.1007/s40979-023-00142-3
28. Hasanein AM, Sobaih AEE. Drivers and Consequences of ChatGPT Use in Higher Education: key Stakeholder Perspectives. *Eur J Investig Health Psychol Educ.* **2023**;13:2599–2614. doi:10.3390/ejihpe13110181
29. Ahmad SF, Han H, Alam MM, et al. Impact of artificial intelligence on human loss in decision making, laziness and safety in education. *Humanit Soc Sci Commun.* **2023**;10:311. doi:10.1057/s41599-023-01787-8
30. Moya B, Eaton S, Pethrick H, et al. Academic Integrity and Artificial Intelligence in Higher Education (HE) Contexts: a Rapid Scoping Review. *Canadian Perspect Acad Integr.* **2024**;7. doi:10.55016/ojs/cpai.v7i3.78123
31. Unogwu OJ, Doshi R, Hiran KK, Mijwil MM, Catherine AT, Abotaleb M. Exploring the Implications of Emerging Artificial Intelligence Technologies at Edge Computing in Higher Education. In: Doshi R, Dadhich M, Poddar S, Hiran KK, editors. *Integrating Generative AI in Education to Achieve Sustainable Development Goals.* Hershey, PA, USA: IGI Global; **2024**:169–182. doi:10.4018/979-8-3693-2440-0.ch009
32. Wong WKO. The sudden disruptive rise of generative artificial intelligence? An evaluation of their impact on higher education and the global workplace. *J Open Innovati.* **2024**;10:100278. doi:10.1016/j.joitmc.2024.100278
33. Sallam M, Al-Mahzoum K, Almutairi Y, et al. Anxiety among Medical Students Regarding Generative Artificial Intelligence Models: a Pilot Descriptive Study. *Preprints.* **2024**. doi:10.20944/preprints202408.1215.v1
34. Ali O, Murray PA, Momin M, Dwivedi YK, Malik T. The effects of artificial intelligence applications in educational settings: challenges and strategies. *Technol Forec Social Change.* **2024**;199:123076. doi:10.1016/j.techfore.2023.123076
35. George B, Wooden O. Managing the Strategic Transformation of Higher Education through Artificial Intelligence. *Administrative Sciences.* **2023**;13(196). doi:10.3390/admsci13090196
36. Roll I, Wylie R. Evolution and Revolution in Artificial Intelligence in Education. *Inter J Artificial Intell Educ.* **2016**;26:582–599. doi:10.1007/s40593-016-0110-3
37. Chan CKY. A comprehensive AI policy education framework for university teaching and learning. *Int J Educa Technol High Educ.* **2023**;20:38. doi:10.1186/s41239-023-00408-3
38. Liu M, Ren Y, Nyagoga LM, Stonier F, Wu Z, Yu L. Future of education in the era of generative artificial intelligence: consensus among Chinese scholars on applications of ChatGPT in schools. *Future Educat Res.* **2023**;1:72–101. doi:10.1002/fer3.10
39. McCarthy AM, Maor D, McConney A, Cavanaugh C. Digital transformation in education: critical components for leaders of system change. *Social Scien Humanit Open.* **2023**;8:100479. doi:10.1016/j.ssaho.2023.100479
40. Spector JM, Ma S. Inquiry and critical thinking skills for the next generation: from artificial intelligence back to human intelligence. *Smart Learn Envir.* **2019**;6:8. doi:10.1186/s40561-019-0088-z
41. Essel HB, Vlachopoulos D, Essuman AB, Amankwa JO. ChatGPT effects on cognitive skills of undergraduate students: receiving instant responses from AI-based conversational large language models (LLMs). *Comput Educat.* **2024**;6:100198. doi:10.1016/j.caeai.2023.100198
42. Mijwil M, Abotaleb M, Guma ALI, Dhoska K. Assigning Medical Professionals: chatGPT's Contributions to Medical Education and Health Prediction. *Mesopotamian J Artificial Intelli Health.* **2024**;2024:76–83. doi:10.58496/MJAIH/2024/011
43. Sallam M. Bibliometric top ten healthcare-related ChatGPT publications in the first ChatGPT anniversary. *Narra J.* **2024**;4:e917. doi:10.52225/narra.v4i2.917
44. Yilmaz Muluk S, Olcucu N. The Role of Artificial Intelligence in the Primary Prevention of Common Musculoskeletal Diseases. *Cureus.* **2024**;16:e65372. doi:10.7759/cureus.65372
45. Yilmaz Muluk S, Olcucu N. Comparative Analysis of Artificial Intelligence Platforms: chatGPT-3.5 and GoogleBard in Identifying Red Flags of Low Back Pain. *Cureus.* **2024**;16:e63580. doi:10.7759/cureus.63580
46. Reddy S. Generative AI in healthcare: an implementation science informed translational path on application, integration and governance. *Implement Sci.* **2024**;19:27. doi:10.1186/s13012-024-01357-9

47. Newton PM, Xiromeriti M. ChatGPT performance on MCQ exams in higher education. A pragmatic scoping review. *EdArXiv*. 2023. doi:10.35542/osf.io/sytu3
48. Sallam M, Barakat M, Sallam M. A Preliminary Checklist (METRICS) to Standardize the Design and Reporting of Studies on Generative Artificial Intelligence-Based Models in Health Care Education and Practice: development Study Involving a Literature Review. *Interact J Med Res*. 2024;15:e54704. doi:10.2196/54704
49. Kung TH, Cheatham M, Medenilla A, 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
50. Newton P, Xiromeriti M. ChatGPT performance on multiple choice question examinations in higher education. A pragmatic scoping review. *Assess Eval Higher Educ*. 2024;1–18. doi:10.1080/02602938.2023.2299059
51. Liu X, Wu J, Shao A, et al. Uncovering Language Disparity of ChatGPT on Retinal Vascular Disease Classification: cross-Sectional Study. *J Med Internet Res*. 2024;26:e51926. doi:10.2196/51926
52. Rosol M, Gąsior JS, Laba J, Korzeniewski K, Młyniczak M. Evaluation of the performance of GPT-3.5 and GPT-4 on the Polish Medical Final Examination. *Sci Rep*. 2023;13:20512. doi:10.1038/s41598-023-46995-z
53. Siebielec J, Ordak M, Oskroba A, Dworakowska A, Bujalska-Zadrozny M. Assessment Study of ChatGPT-3.5's Performance on the Final Polish Medical Examination: accuracy in Answering 980 Questions. *Healthcare*. 2024;12. doi:10.3390/healthcare12161637
54. Guillen-Grima F, Guillen-Aguinaga S, Guillen-Aguinaga L, et al. Evaluating the Efficacy of ChatGPT in Navigating the Spanish Medical Residency Entrance Examination (MIR): promising Horizons for AI in Clinical Medicine. *Clin Pract*. 2023;13:1460–1487. doi:10.3390/clinpract13060130
55. Sallam M, Al-Mahzoum K, Alshuaib O, et al. Language discrepancies in the performance of generative artificial intelligence models: an examination of infectious disease queries in English and Arabic. *BMC Infect Dis*. 2024;24:799. doi:10.1186/s12879-024-09725-y
56. Sallam M, Barakat M, Sallam M. Pilot Testing of a Tool to Standardize the Assessment of the Quality of Health Information Generated by Artificial Intelligence-Based Models. *Cureus*. 2023;15:e49373. doi:10.7759/cureus.49373
57. Douglas M, Wilson J, Ennis S. Multiple-choice question tests: a convenient, flexible and effective learning tool? A case study. *Innov Educat Teach Int*. 2012;49:111–121. doi:10.1080/14703297.2012.677596
58. Bloom BS, Krathwohl DR. *Taxonomy of Educational Objectives: The Classification of Educational Goals*. Longmans, Green; 1956:403.
59. Seaman M. Bloom's Taxonomy: its Evolution, Revision, and Use in the Field of Education. *CurricuTeach Dial*. 2011;13:29–131A.
60. Liu Q, Wald N, Daskon C, Harland T. Multiple-choice questions (MCQs) for higher-order cognition: perspectives of university teachers. *Innov Educat Teach Int*. 2023;1–13. doi:10.1080/14703297.2023.2222715
61. Karanja E, Malone LC. Improving project management curriculum by aligning course learning outcomes with Bloom's taxonomy framework. *J Interna Educat Busin*. 2021;14:197–218. doi:10.1108/JIEB-05-2020-0038
62. Herrmann-Werner A, Festl-Wietek T, Holderried F, et al. Assessing ChatGPT's Mastery of Bloom's Taxonomy using psychosomatic medicine exam questions. *medRxiv*. 2023. doi:10.1101/2023.08.18.23294159
63. Herrmann-Werner A, Festl-Wietek T, Holderried F, et al. Assessing ChatGPT's Mastery of Bloom's Taxonomy Using Psychosomatic Medicine Exam Questions: mixed-Methods Study. *J Med Internet Res*. 2024;26:e52113. doi:10.2196/52113
64. Sallam M, Al-Salahat K. Below average ChatGPT performance in medical microbiology exam compared to university students. *Frontiers in Education*. 2023;8:1333415. doi:10.3389/educ.2023.1333415
65. OpenAI. GPT-3.5. Available from: <https://openai.com/>. Accessed November 27, 2023.
66. Microsoft O Bing is your AI-powered copilot for the web. Available from: <https://www.bing.com/search?q=Bing+AI&showconv=1&FORM=hpcodx>. Accessed November 27, 2023.
67. Google. Bard. Available from: <https://bard.google.com/chat>. Accessed November 27, 2023.
68. Burtis CA, Ashwood ER, Bruns DE, Tietz NW. *Tietz Textbook of Clinical Chemistry and Molecular Diagnostics*. 5th ed. Saunders: St. Louis, Mo; 2013.
69. Bishop ML, Fody EP, Schoeff LE. *Clinical Chemistry: Principles, Techniques, and Correlations*. Eighth edition. ed. Philadelphia: Wolters Kluwer; 2018:736.
70. McPherson RA, Pincus MR. *Henry's Clinical Diagnosis and Management by Laboratory Methods*. 24. ed. Philadelphia: Elsevier; 2021:pagescm.
71. Sallam M, Al-Salahat K, Al-Ajlouni E. ChatGPT Performance in Diagnostic Clinical Microbiology Laboratory-Oriented Case Scenarios. *Cureus*. 2023;15:e50629. doi:10.7759/cureus.50629
72. Lo CK. What Is the Impact of ChatGPT on Education? A Rapid Review of the Literature. *Education Sciences*. 2023;13:410. doi:10.3390/educsci13040410
73. Sallam M, Salim NA, Al-Tammemi AB, et al. ChatGPT Output Regarding Compulsory Vaccination and COVID-19 Vaccine Conspiracy: a Descriptive Study at the Outset of a Paradigm Shift in Online Search for Information. *Cureus*. 2023;15:e35029. doi:10.7759/cureus.35029
74. Jung LB, Gudera JA, Wiegand TLT, Allmendinger S, Dimitriadis K, Koerte IK. ChatGPT Passes German State Examination in Medicine With Picture Questions Omitted. *Dtsch Arztebl Int*. 2023;120:373–374. doi:10.3238/arztebl.m2023.0113
75. Yanagita Y, Yokokawa D, Uchida S, Tawara J, Ikusaka M. Accuracy of ChatGPT on Medical Questions in the National Medical Licensing Examination in Japan: evaluation Study. *JMIR Form Res*. 2023;7:e48023. doi:10.2196/48023
76. 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
77. Gobira M, Nakayama LF, Moreira R, Andrade E, Regatieri CVS, Belfort R Jr. Performance of ChatGPT-4 in answering questions from the Brazilian National Examination for Medical Degree Revalidation. *Rev Assoc Med Bras*. 2023;69:e20230848. doi:10.1590/1806-9282.20230848
78. Hofmann Hayden L, Guerra Gage A, Le Jonathan L, et al. The Rapid Development of Artificial Intelligence: GPT-4's Performance on Orthopedic Surgery Board Questions. *Orthopedics*. 2023;2023:1–5. doi:10.3928/01477447-20230922-05
79. Mbakwe AB, Lourentzou I, Celi LA, Mechanic OJ, Dagan A. ChatGPT passing USMLE shines a spotlight on the flaws of medical education. *PLOS Digital Health*. 2023;2:e0000205. doi:10.1371/journal.pdig.0000205
80. Zhai X, Nyaaba M, Ma W. Can AI Outperform Humans on Cognitive-demanding Tasks in Science? *SSRN*. 2023. doi:10.2139/ssrn.4451722
81. Ghosh A, Bir A. Evaluating ChatGPT's Ability to Solve Higher-Order Questions on the Competency-Based Medical Education Curriculum in Medical Biochemistry. *Cureus*. 2023;15:e37023. doi:10.7759/cureus.37023

82. Tlili A, Shehata B, Adarkwah MA, et al. What if the devil is my guardian angel: chatGPT as a case study of using chatbots in education. *Smart Learning Environments*. 2023;10:15. doi:10.1186/s40561-023-00237-x
83. Dai W, Lin J, Jin H, et al. Can Large Language Models Provide Feedback to Students? A Case Study on ChatGPT. In Proceedings of the 2023 IEEE International Conference on Advanced Learning Technologies (ICALT), 2023; pp. 323–325.
84. Schleiss J, Laupichler MC, Raupach T, Stober S. AI Course Design Planning Framework: developing Domain-Specific AI Education Courses. *Education Sciences*. 2023;13(954). doi:10.3390/educsci13090954
85. Perkins M. Academic integrity considerations of AI Large Language Models in the post-pandemic era: chatGPT and beyond. *J Univer Teac Learn Pract*. 2023;20. doi:10.53761/1.20.02.07
86. Memarian B, Doleck T. ChatGPT in education: methods, potentials, and limitations. *Comp Human Behav*. 2023;1:100022. doi:10.1016/j.chbah.2023.100022
87. Saylam S, Duman N, Yildirim Y, Satsevich K. Empowering education with AI: addressing ethical concerns. *London J So Scien*. 2023;39–48. doi:10.31039/ljss.2023.6.103
88. Grájeda A, Burgos J, Córdova P, Sanjinés A. Assessing student-perceived impact of using artificial intelligence tools: construction of a synthetic index of application in higher education. *Cogent Educat*. 2024;11:2287917. doi:10.1080/2331186X.2023.2287917
89. Hadi Mogavi R, Deng C, Juho Kim J, et al. ChatGPT in education: a blessing or a curse? A qualitative study exploring early adopters' utilization and perceptions. *Comp Human Behav*. 2024;2:100027. doi:10.1016/j.chbah.2023.100027
90. Alfertshofer M, Hoch CC, Funk PF, et al. Sailing the Seven Seas: a Multinational Comparison of ChatGPT's Performance on Medical Licensing Examinations. *Ann. Biomed. Eng*. 2023. doi:10.1007/s10439-023-03338-3
91. Sallam M, Mousa D. Evaluating ChatGPT performance in Arabic dialects: a comparative study showing defects in responding to Jordanian and Tunisian general health prompts. *Mesopotamian J Artificial Intelli Health*. 2024;2024:1–7. doi:10.58496/MJAIH/2024/001

Advances in Medical Education and Practice

Dovepress

Publish your work in this journal

Advances in Medical Education and Practice is an international, peer-reviewed, open access journal that aims to present and publish research on Medical Education covering medical, dental, nursing and allied health care professional education. The journal covers undergraduate education, postgraduate training and continuing medical education including emerging trends and innovative models linking education, research, and health care services. The manuscript management system is completely online and includes a very quick and fair peer-review system. Visit <http://www.dovepress.com/testimonials.php> to read real quotes from published authors.

Submit your manuscript here: <http://www.dovepress.com/advances-in-medical-education-and-practice-journal>