

as as

2-MANUSCRIPT.docx

Ankara Yıldırım Beyazıt Üniversitesi

Document Details

Submission ID

trn:oid::3117:499323616

Submission Date

Sep 16, 2025, 10:42 AM GMT+3

Download Date

Sep 16, 2025, 10:44 AM GMT+3

File Name

2-MANUSCRIPT.docx

File Size

33.1 KB

11 Pages

2,253 Words

14,612 Characters

7% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Match Groups

-  **14 Not Cited or Quoted 7%**
Matches with neither in-text citation nor quotation marks
-  **0 Missing Quotations 0%**
Matches that are still very similar to source material
-  **0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
-  **0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 6%  Internet sources
- 4%  Publications
- 0%  Submitted works (Student Papers)

Integrity Flags

0 Integrity Flags for Review

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.

Match Groups

- **14 Not Cited or Quoted 7%**
Matches with neither in-text citation nor quotation marks
- **0 Missing Quotations 0%**
Matches that are still very similar to source material
- **0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
- **0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 6% ■ Internet sources
- 4% ■ Publications
- 0% ■ Submitted works (Student Papers)

Top Sources

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.

1	Publication	Boaz Karmazyn, Christopher L. Newman, Megan B Marine, Mathew R Wanner et a...	1%
2	Internet	www.frontiersin.org	1%
3	Internet	jksronline.org	<1%
4	Internet	myresearchspace.uws.ac.uk	<1%
5	Internet	www.ncbi.nlm.nih.gov	<1%
6	Internet	explodingtopics.com	<1%
7	Publication	Berna Akkus Yildirim, Baver Tutun, Gorkem Durak, Emre Batuhan Yildirim, Emre ...	<1%
8	Publication	Cagatay Bolgen, Birsun Unal Daphan. "Is CT Still the Gold Standard in Semicircula...	<1%
9	Publication	Xinhang Li, Zhenxu Li, Shengfa Pan, Li Zhang, Yanbin Zhao, Xin Chen, Yu Sun, Feif...	<1%
10	Internet	eurjther.com	<1%

11 Internet

lucris.lub.lu.se <1%

12 Internet

www.biomedcentral.com <1%

13 Publication

Kenichi Tanaka, Yosuke Okada, Saeko Umezu, Ryoma Hashimoto et al. "Effects of ... <1%

Introduction

Artificial intelligence (AI) is becoming an important part of modern medicine, and large language models (LLMs) such as ChatGPT are now widely used by both clinicians and patients (1,2). These systems can provide information, generate treatment suggestions, and support education. In anesthesia, where clinical decisions must be made quickly and accurately, the potential value of such tools is high. However, the reliability and scientific accuracy of LLMs in perioperative care are not yet fully understood. (3-5).

Anesthesia is a suitable field to test AI because it involves many acute and life-threatening conditions such as anaphylaxis, malignant hyperthermia, airway obstruction, and cardiovascular collapse. In these situations, correct recognition and rapid intervention are essential. Mistakes can result in serious harm. Previous studies have shown that AI can assist in protocol standardization, monitoring, and perioperative risk assessment, but the problem of incorrect or incomplete answers remains a major limitation (6,7)

Recent evaluations in different medical fields demonstrated that ChatGPT can often produce guideline-based responses, but it may fail in complex or ambiguous cases (8-12). In anesthesiology, such failures may lead to incorrect drug recommendations or delayed recognition of emergencies. Therefore, it is necessary to evaluate ChatGPT in real-life-like clinical situations. This type of assessment can help to decide whether it may be used as a decision-support tool or should remain limited to education and training.

One practical way to investigate this is to compare ChatGPT answers with those of experienced anesthesiologists in structured clinical scenarios. Simulation cases are already widely used in anesthesia education and provide a safe and standardized method for testing performance (3). Scenarios such as peri-induction hypotension, intraoperative arrhythmias, respiratory problems,

or postoperative complications can be applied to both experts and ChatGPT for comparison (13,14).

10 We aimed to evaluate the accuracy and clinical validity of ChatGPT's responses to anesthesia-related clinical scenarios by directly comparing them with expert anesthesiologists' assessments.

3 **Materials and Methods**

This study did not require approval from an institutional review board because it did not involve human subjects, patient data, or interventions in clinical care. All scenarios were hypothetical case simulations created for educational and research purposes. The evaluation process was limited to expert opinion and artificial intelligence responses, without any patient participation.

Study Design

11 We conducted a prospective comparative study in which clinical scenarios commonly encountered in anesthesiology were presented to both ChatGPT (OpenAI, San Francisco, USA) and anesthesiologists. The aim was to evaluate the scientific accuracy and clinical reliability of ChatGPT responses by direct comparison with expert assessments. This design allowed for a structured and reproducible approach to measure agreement levels and to identify discrepancies.

Scenarios

A total of sixteen standardized scenarios were developed, covering different categories of perioperative practice. These included peri-induction hypotension, anaphylaxis, malignant hyperthermia, arrhythmias, airway complications such as laryngeal edema, respiratory events including pneumothorax and hypoventilation, and postoperative complications such as delirium or pulmonary embolism. Each scenario was constructed based on existing literature, practice guidelines, and simulation training materials. Clinical details were presented in a stepwise manner, including patient demographics, anesthetic plan, intraoperative events, vital signs, and progression of symptoms. The reason for selecting a total of sixteen scenarios was based on the consensus of both expert groups, representing the most common and clinically significant situations encountered in anesthesiology practice. This ensured that the study was structured to cover widespread and critical perioperative events reflective of real-life cases.

Selection of Evaluators

Two anesthesiologists with at least 10 years of independent clinical experience in tertiary care centers were selected as expert evaluators. Their clinical background included both general and subspecialty anesthesia practice. Evaluators were blinded to each other's responses to avoid bias. ChatGPT was provided with the same scenarios under identical conditions, without additional contextual prompts beyond the case description.

Evaluation Criteria and Guidelines

The evaluators (AAK,XXX) judged ChatGPT's responses using a structured framework. Key criteria included: (a) accuracy of diagnosis, (b) appropriateness of treatment recommendation, (c) compliance with international guidelines (American Society of Anesthesiologists [ASA] practice parameters, European Society of Anaesthesiology and Intensive Care [ESAIC] guidelines, and World Health Organization recommendations), and (d) clarity and applicability to real clinical practice. Each item was rated on a Likert scale ranging from "inaccurate" to "fully accurate." (15).

Data Collection

All responses were collected in written form. ChatGPT answers were generated in English without manual editing. Expert responses were recorded independently in structured forms. Data were anonymized, coded, and stored in electronic format. Disagreements between evaluators were resolved by consensus discussions.

Evaluated Parameters

The main parameters assessed included: (a) diagnostic correctness, (b) identification of whether treatment was required, and (c) recognition of first-line treatment options. Additional parameters included response completeness, adherence to guidelines, and the presence of

potentially unsafe recommendations. These metrics allowed both qualitative and quantitative assessment of ChatGPT performance.

5 **Statistical Analysis**

8 Descriptive statistics were applied to summarize evaluator ratings and ChatGPT responses. Categorical variables, such as diagnostic correctness and treatment appropriateness, were presented as percentages. Agreement between ChatGPT and expert anesthesiologists was assessed using Cohen's kappa coefficient. Weighted kappa was additionally applied for Likert-scale ratings, where partial agreement between adjacent categories (e.g., "partially correct" vs. 12 "fully correct") was taken into account. The degree of agreement was interpreted according to the classification by Landis and Koch, with κ values of 0.00–0.20 considered slight, 0.21–0.40 1 fair, 0.41–0.60 moderate, 0.61–0.80 substantial, and 0.81–1.00 almost perfect agreement. 13 Continuous variables, such as response length, were expressed as mean \pm standard deviation. A p-value < 0.05 was considered statistically significant. All analyses were conducted using SPSS 2 version 26.0 (IBM Corp., Armonk, NY, USA).

Results

A total of sixteen standardized anesthesia scenarios were evaluated by ChatGPT and two expert anesthesiologists. ChatGPT correctly identified the diagnosis in 14 out of 16 cases (88%) and appropriately determined whether treatment was required in 15 out of 16 cases (93%). The accuracy of recommending the correct first-line treatment was slightly lower, with full agreement observed in 13 cases (81%). The overall concordance rate across all domains was 87%. These data are summarized in Table 1.

9 Inter-rater reliability between the two human experts was excellent, with a Cohen's kappa value of 0.82, indicating almost perfect agreement. When ChatGPT's answers were compared with each expert, substantial agreement was observed ($\kappa = 0.74$ with Expert 1 and $\kappa = 0.71$ with Expert 2). These findings suggest that ChatGPT achieved a performance level comparable to experienced anesthesiologists in most scenarios. These findings are summarized in Table 2.

Scenario-specific analysis revealed several notable points. In acute conditions such as anaphylaxis, malignant hyperthermia, and pulmonary embolism, ChatGPT provided fully accurate diagnostic and therapeutic recommendations consistent with international guidelines. In vasovagal syncope, the model correctly recognized the diagnosis and initial management but did not specify atropine dosing, which was considered a partial gap. For postoperative delirium, ChatGPT emphasized antipsychotic use earlier than experts, who prioritized opioid dose reduction, hydration, and environmental modification. In aspiration pneumonia, ChatGPT correctly identified the condition but prematurely suggested antibiotic initiation, while experts highlighted airway management and oxygenation as the immediate priority. These results are summarized in Table 3.

Overall, ChatGPT demonstrated high diagnostic reliability and substantial concordance with expert evaluations, particularly in life-threatening scenarios requiring urgent intervention.

However, minor discrepancies in therapeutic sequencing and omission of specific drug doses highlighted its current limitations for use as an independent decision-making tool in clinical anesthesiology.

Discussion

The present study evaluated the performance of ChatGPT in simulated anesthesia scenarios by directly comparing its answers with those of two experienced anesthesiologists. The main findings demonstrated that ChatGPT provided correct diagnoses in 88% of cases, identified the need for treatment in 93% of cases, and recommended the correct first-line treatment in 81% of cases. The overall concordance rate was 87%, and inter-rater reliability showed almost perfect agreement between the two experts ($\kappa = 0.82$) and substantial agreement between ChatGPT and both experts ($\kappa = 0.74$ and $\kappa = 0.71$, respectively). These results indicate that ChatGPT can generate responses that are broadly consistent with expert reasoning in critical perioperative conditions but still exhibits important limitations in therapeutic precision.

In reviewing the literature, our results align with prior reports that large language models can often generate guideline-based outputs, particularly in acute scenarios with well-established management protocols. For example, in their study, Gilson et al. reported that ChatGPT achieved 60% accuracy on the United States Medical Licensing Examination, indicating that it could replicate a large portion of clinically relevant knowledge (17). Similarly, in their analysis, Kung et al. found that ChatGPT provided coherent and clinically appropriate responses in internal medicine board questions, though gaps were noted in pharmacology and therapeutic decision-making (18). In anesthesiology specifically, Wan et al. demonstrated that ChatGPT could offer reasonable advice in airway management algorithms but tended to omit critical details such as drug doses and alternative approaches (19). Our study expands on these observations by systematically testing ChatGPT in diverse intraoperative scenarios, including hemodynamic instability, arrhythmias, airway emergencies, and postoperative complications, thereby providing a more comprehensive view of its clinical reliability.

The numerical results of this study further highlight both the promise and the shortcomings of ChatGPT. The fact that ChatGPT achieved an 88% diagnostic accuracy suggests that its large

training corpus enables recognition of common anesthetic patterns such as anaphylaxis, malignant hyperthermia, and pulmonary embolism. In these scenarios, where international guidelines provide standardized treatment algorithms, the model was able to reproduce the expected responses with high fidelity. This was evident in its recommendation of epinephrine and fluid resuscitation in anaphylaxis, dantrolene in malignant hyperthermia, and anticoagulation in pulmonary embolism, which matched expert assessments without major deviation. However, in 19% of scenarios, the first-line treatment was either incomplete or suggested prematurely. For example, ChatGPT recommended antipsychotics early in postoperative delirium, while experts emphasized non-pharmacological interventions first. Similarly, in aspiration pneumonia, the model correctly identified the diagnosis but suggested antibiotics earlier than airway management, deviating from guideline priorities.

These findings echo previous reports on the limitations of LLMs. In their analysis, Rosen et al. observed that ChatGPT often generated plausible but incomplete treatment strategies in emergency medicine vignettes (20). In their study, Pham et al. showed that while ChatGPT could reproduce core ACLS algorithms, it sometimes confused drug sequencing and dosages. Such discrepancies are consistent with our observation that the model is strong in pattern recognition but weaker in therapeutic nuance, particularly when management requires stepwise prioritization rather than simultaneous interventions (21).

The strength of our study lies in its structured, prospective design with predefined scenarios and standardized evaluation criteria. Unlike retrospective content analyses, we presented identical cases to both experts and the AI model, allowing direct comparison. Furthermore, the use of Likert-based scoring and kappa statistics provided quantitative evidence of agreement, with κ values ranging from 0.71 to 0.82 confirming substantial to almost perfect concordance. The prospective nature, randomization of scenario order, and blinded expert assessments reduce bias and enhance the reliability of our conclusions. Additional strengths include the focus on

real-life perioperative emergencies, the systematic comparison across multiple domains (diagnosis, treatment necessity, first-line therapy), and the application of robust statistical methods such as kappa reliability analysis.

4 However, our study has several limitations. The study was conducted in a single center with only two expert evaluators, which may limit generalizability. The sample size of 16 scenarios, while diverse, may not fully capture the breadth of anesthetic practice. Moreover, the evaluations focused on short-term clinical reasoning rather than long-term patient outcomes, as no actual patients were included. Finally, ChatGPT was tested in English only, and its performance might vary across languages and cultural contexts. These limitations suggest caution in extrapolating our findings beyond the controlled simulation environment.

The clinical implications of these results are noteworthy. While ChatGPT achieved high diagnostic accuracy and reasonable concordance with experts, its therapeutic recommendations occasionally lacked detail or sequence accuracy. This supports the view that ChatGPT should not be used as a routine decision-making tool in anesthetic management. Instead, its most appropriate role may be as an adjunct in selected cases, particularly for educational purposes, simulation training, and providing rapid summaries of guideline-based care. Similar to how adjunctive hemostatic agents such as FloSeal® or Surgicel® are not required for every partial nephrectomy but may be considered in selected complex cases, ChatGPT may have a role in complementing—but not replacing—expert judgment in anesthesia practice. In daily clinical use, reliance solely on ChatGPT could pose risks due to its occasional inaccuracies, but when combined with expert oversight, it may enhance efficiency, learning, and decision support.

Future research should expand this work to multicenter designs with larger cohorts of anesthesiologists and a broader array of scenarios. Such studies could stratify performance by scenario complexity, compare multiple LLMs, and investigate whether iterative prompting improves reliability. Longitudinal evaluations could also assess whether repeated exposure to

ChatGPT enhances resident education or simulation training outcomes. Furthermore, incorporating objective outcome measures such as time to recognition of critical events or success in simulated resuscitation would provide deeper insights. Finally, evaluating the model's integration with electronic health records and its potential for real-time perioperative monitoring represents an important future direction.

In conclusion, this study demonstrated that ChatGPT achieved 88% diagnostic accuracy, 93% recognition of treatment need, and 81% concordance in first-line therapy recommendations across sixteen anesthesia scenarios, with overall agreement of 87% and kappa values between 0.71 and 0.82 indicating substantial reliability. While its performance was encouraging in life-threatening conditions such as anaphylaxis and malignant hyperthermia, discrepancies in therapeutic prioritization highlight its limitations as an independent clinical tool. ChatGPT may serve as a valuable adjunct for education and training but should not replace expert judgment in anesthesia practice.