

GUIDELINE

Open Access



# Reporting guideline for Chatbot Health Advice studies: the CHART statement

Bright Huo<sup>1\*</sup>, Gary Collins<sup>2,3</sup>, David Chartash<sup>4</sup>, Arun Thirunavukarasu<sup>5</sup>, Annette Flanagan<sup>6</sup>, Alfonso Iorio<sup>7</sup>, Giovanni Cacciamani<sup>8,9</sup>, Xi Chen<sup>10,11</sup>, Nan Liu<sup>12</sup>, Piyush Mathur<sup>13</sup>, An-Wen Chan<sup>14</sup>, Christine Laine<sup>15,16</sup>, Daniela Pacella<sup>17</sup>, Michael Berkwitz<sup>18</sup>, Stavros A. Antoniou<sup>19</sup>, Jennifer C. Camaradou<sup>20</sup>, Carolyn Canfield<sup>21</sup>, Michael Mittelman<sup>22</sup>, Timothy Feeney<sup>23,24</sup>, Elizabeth Loder<sup>23,25</sup>, Riaz Agha<sup>26,27</sup>, Ashirbani Saha<sup>28</sup>, Julio Mayol<sup>29</sup>, Anthony Sunjaya<sup>30</sup>, Hugh Harvey<sup>31</sup>, Jeremy Y. Ng<sup>32</sup>, Tyler McKechnie<sup>1</sup>, Yung Lee<sup>1,33</sup>, Nipun Verma<sup>34</sup>, Gregor Stiglic<sup>35</sup>, Melissa McCradden<sup>36</sup>, Karim Ramji<sup>37</sup>, Vanessa Boudreau<sup>1</sup>, Monica Ortenzi<sup>38</sup>, Joerg Meerpohl<sup>39,40</sup>, Per Olav Vandvik<sup>40,41</sup>, Thomas Agoritsas<sup>7,41,42</sup>, Diana Samuel<sup>43</sup>, Helen Frankish<sup>44</sup>, Michael Anderson<sup>45,46</sup>, Xiaomei Yao<sup>28</sup>, Stacy Loeb<sup>47</sup>, Cynthia Lokker<sup>7</sup>, Xiaoxuan Liu<sup>48</sup>, Eliseo Guallar<sup>49</sup>, Gordon Guyatt<sup>7,41</sup> and The CHART Collaborative

## Abstract

**Background** The Chatbot Assessment Reporting Tool (CHART) is a reporting guideline developed to provide reporting recommendations for studies evaluating the performance of generative artificial intelligence (AI)-driven chatbots when summarizing clinical evidence and providing health advice, referred to as Chatbot Health Advice (CHA) studies.

**Methods** CHART was developed in several phases after performing a comprehensive systematic review to identify variation in the conduct, reporting, and methodology in CHA studies. Findings from the review were used to develop a draft checklist that was revised through an international, multidisciplinary modified asynchronous Delphi consensus process of 531 stakeholders, three synchronous panel consensus meetings of 48 stakeholders, and subsequent pilot testing of the checklist.

**Results** CHART includes 12 items and 39 subitems to promote transparent and comprehensive reporting of CHA studies. These include Title (subitem 1a), Abstract/Summary (subitem 1b), Background (subitems 2ab), Model Identifiers (subitems 3ab), Model Details (subitems 4abc), Prompt Engineering (subitems 5ab), Query Strategy (subitems 6abcd), Performance Evaluation (subitems 7ab), Sample Size (subitem 8), Data Analysis (subitem 9a), Results (subitems 10abc), Discussion (subitems 11abc), Disclosures (subitem 12a), Funding (subitem 12b), Ethics (subitem 12c), Protocol (subitem 12d), and Data Availability (subitem 12e).

**Conclusion** The CHART checklist and corresponding methodological diagram were designed to support key stakeholders including clinicians, researchers, editors, peer reviewers, and readers in reporting, understanding, and interpreting the findings of CHA studies.

\*Correspondence:

Bright Huo  
brighthuo@dal.ca

Full list of author information is available at the end of the article



© The Author(s) 2025. **Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

## Key messages

- CHART was developed by performing a systematic review, Delphi consensus of 531 international stakeholders, and several consensus meetings among an expert panel comprised 48 members.
- The CHART statement outlines 12 key reporting items for Chatbot Health Advice studies in the form of a checklist and methodology diagram.
- All stakeholders including clinicians, researchers, and journal editors should encourage the transparent reporting of Chatbot Health Advice studies.

**Keywords** LLMs, Generative AI, Reporting standards

## Background

Artificial intelligence (AI) has made great strides toward clinical applications in healthcare, with deep learning algorithms performing comparably to current gold standards in several areas in patient care [1, 2]. With the introduction of large language models (LLMs) into mainstream use, there has been an explosive rise in the number of studies evaluating the performance of generative artificial intelligence (AI)-driven chatbots in summarizing evidence and providing health advice [3], termed Chatbot Health Advice (CHA) studies. Investigators typically develop prompts to query generative AI models through a chat-based interface for the purpose of summarizing clinical evidence or obtaining health advice including but not limited to health promotion, prevention, screening, diagnosis, treatment, and/or general health information. For example, physicians may query generative AI-driven chatbots to identify whether their patient should receive colorectal cancer screening [4]. Similarly, a patient may ask questions about their upcoming surgery for gastroesophageal reflux disease [5]. The intense interest in using generative AI-driven chatbots for health advice has generated numerous CHA studies in a short timeframe [6]. Investigators may include clinicians, scientists, or patients, bringing different technical expertise and personal perspectives to study methodology including prompt engineering and model response evaluation.

These studies represent a growing genre of medical AI research [7]. At least 137 CHA studies were published less than a year after the release of ChatGPT in November 2022, but the completeness of reporting among these studies has been highly variable [6]. For instance, few articles elaborate on the development of their prompts, while fewer than 40% of articles report key elements of their query strategy including the date of their search, the number of chat sessions used, or the number of prompts [6]. Raw prompts and model output are infrequently reported, and most articles present an insufficient amount of information to identify the model and

chatbot under evaluation [6]. This problem is important because inadequate reporting impairs the ability of readers to interpret the validity and reliability of study findings [8]. Flaws in the design, data collection, or conduct of a study may lead to erroneous conclusions or raise the risk of patient harm, particularly if generative AI-driven models are used for health purposes [9]. Complete and standardized reporting facilitates critical appraisal and may help identify applications with genuine potential to improve health care, building trust in the use of generative AI models in clinical practice among clinicians, patients, and the general public [9].

In response to the growing need for reporting standards for evaluating CHA studies for clinical purposes [10], we developed the Chatbot Assessment Reporting Tool (CHART). This reporting standard is an international, multidisciplinary initiative registered with the Enhancing the QUALity and Transparency Of health Research (EQUATOR) Network [11] and was announced in December 2023 [3]. This article describes the methodology used to identify, evaluate, and gain consensus on the checklist items and diagram that comprise CHART. We aimed to develop robust guidance to promote high methodological rigor and transparent reporting of CHA studies evaluating the performance of generative AI-driven chatbots when summarizing clinical evidence and providing health advice. The terminology used in this reporting guideline is listed in Table 1.

## Methods

We formed a steering group responsible for overseeing the development of CHART. We developed CHART in alignment with the EQUATOR Network's framework according to the highest methodological standards for reporting guideline development [8] and published the protocol in May 2024 [7].

To inform the development of CHART, we conducted a comprehensive systematic review to identify information reported in CHA studies. The review protocol was prospectively registered on the Open Sciences Framework:

**Table 1** Glossary

Term	Definition
Artificial intelligence (AI)	The science of developing computer systems that can perform complex tasks approximating human cognitive performance
Base model	A pre-existing generative AI model
Chat session	An interface in a computing device through which communication takes place between a chatbot and its user through text-based prompts
Chatbot Health Advice (CHA) study	Any research study evaluating the performance of chatbots when summarizing health evidence and/or providing clinical advice
Fine-tuned model	A base model that has been manipulated through various methods of algorithmic tuning to alter its performance with specificity; methods include but are not limited to reinforcement learning or distillation
Generative AI-driven chatbot	A program that permits users to interact with an AI model (such as an LLM) that is designed to respond to user prompts
Ground truth	The reference standard, or criteria, on which the model is evaluated to define successful performance
Large language model (LLM)	A type of AI model comprising large neural networks trained over large amounts of text usually to produce an output of continuations of text from corresponding prompts known as next word prediction. LLMs are a subset of generative AI models
Multimodal LLM	LLMs with the capacity to integrate input from various data types including text speech and/or visual sources
Natural language processing (NLP)	A branch of information science that seeks to enable computers to interpret and manipulate human text
Parameter	A variable that is tuned iteratively or automatically to optimize the intended outcome of the algorithm. Parameters may be at the model level to optimize tuning (hyperparameters) or “weights” within the model linking layer to layer (parameters)
Post-implementation/deployment	Refers to alteration of the generative AI model following its release
Pre-implementation/deployment	Refers to alteration of the generative AI model prior to its release
Prompt	The input provided by users when interfacing with a generative AI-driven chatbot, leading to input interaction with the AI model
Prompt engineering	An iterative testing phase where various pieces of text are inputted into a chatbot to achieve an output informing the development of study prompts
Query	The act of communicating with a generative AI-driven chatbot by inputting a prompt into the chatbot which might be a question, comment, or phrase to elicit specific desired outputs from the generative AI model
Response	The output of the generative AI-driven chatbot
Tuned model	A base model that has been altered to provide focused responses by means other than fine-tuning, including but not limited to retrieval augmented generation, which seeks to alter performance rather than the model
Zero shot	A machine learning paradigm in which the task (such as classification) is performed without explicit training, fine-tuning, or other optimization

<https://osf.io/cxsk3>. The systematic review was devised according to methodological guidance from the Joanna Briggs Institute [12]. A systematic literature search was performed with the support of a health sciences librarian using Medline via Ovid, Embase via Elsevier, and Web of Science on October 27th, 2023. Full search syntax from all database searches are provided in the supplementary section of our systematic review [6]. We screened 7752 articles to identify 137 eligible articles of interest. Considerable variation in methodology and reporting was observed, and we identified 120 candidate checklist items for CHART (Appendix 1). Full details on this process can be found in our protocol [7]. To evaluate these candidate checklist items for inclusion in the CHART checklist, we invited an advisory committee to perform a modified Delphi consensus process and formed an expert panel to conduct synchronous consensus meetings. Full details on this recruitment process can be found in the protocol [7]. We considered “experts” as individuals who have made

important contributions academically to their discipline, with an emphasis on individuals that have participated in reporting guideline development previously.

#### Modified delphi consensus survey

The steering group invited 1043 members globally to form an advisory committee to participate in a Delphi survey, comprising clinicians, epidemiologists, research methodologists, generative AI researchers, journal editors, chatbot researchers, ethicists, regulatory experts, policy experts, and patient partners. We identified potential committee members using a multi-pronged approach through co-authors published in the top medical journals, public and internal calls through affiliate journals, as well as through snowballing via all members of our expert panel. To identify the top 10 journals across all specialties, we used the journal ranking feature in Scimago. Full details are listed in our protocol [6]. Via convenience sampling, we included four editors

from the top journals identified. We invited members by email and provided project details as well as our correspondence article and study protocol [3, 7]. Members voluntarily registered to participate in our Delphi consensus survey by providing basic demographic information, as well as details of their prior research experience and content expertise. We presented candidate checklist items to the advisory committee using the online Delphi consensus platform Welphi, *Decision Eyes* ([www.welphi.com](http://www.welphi.com)). Members rated candidate checklist items as one of the following: “include, maybe include, uncertain, maybe exclude, or exclude.” They also suggested additional checklist items. After the first round of voting, advisory committee members engaged in a second round of voting via a modified Delphi consensus survey. Members were able to view the results from the first round and review comments supporting voting considerations. During the second Delphi round, members voted on the same checklist items as well as any additional checklist items from the first round. Advisory committee members were also able to suggest additional checklist items during the second round, generating a total of 28 additional candidate checklist items across both Delphi rounds. A total of 531/1043 (50.9%) members participated in both Delphi consensus rounds, rating a total of 140 candidate checklist items for review by the expert panel (Appendix 1).

### Expert panel consensus

The steering group assembled an international, multi-disciplinary panel comprising a balanced representation of 48 relevant stakeholders including clinicians, statisticians, research methodologists, reporting guideline developers, generative AI researchers, journal editors, chatbot researchers, ethicists, regulatory experts, policy experts, and four patient partners. The distribution of stakeholders among the panel is presented in the supplementary material. The steering group used a prespecified threshold of 80% agreement for inclusion to show majority consensus based on prior work [7, 13]. We identified items with at least 80% consensus with the selection of either “include” and “maybe include” together, or “exclude” and “maybe exclude” and posed to the panel whether to include or exclude suggested items. Items not meeting 80% consensus were posed to the panel for further discussion. We also presented raw scores including absolute and relative and frequencies to the expert panel to support their interpretation and decision-making. We held synchronous discussions over three separate panel consensus meetings on Zoom spanning 12 twelve collective hours on June 30th, August 5th, and September 2nd, 2024. Items on which the expert panel disagreed with the advisory committee, as well as items voted as “unsure” by the advisory committee, were discussed among panel

members until consensus was reached. Panel members were able to suggest changes to the phrasing of checklist items, as well as suggest additional checklist items. After extensive discussion, the expert panel reached consensus on 12 checklist items (Appendix 2) and 9 abstract checklist items (Appendix 3). A fillable methodological diagram can be found in Appendix 4. A list of panel members can be found in Appendix 5. No items or subitems required voting, as contentious items were discussed thoroughly until consensus was achieved.

### Pilot testing

Following the panel consensus meetings, draft checklist items were presented to authors of separate, prior CHA studies via an iterative process for pilot testing. Groups of five authors used the draft CHART checklist to evaluate 10 published CHA studies and provide feedback in each round until saturation was reached with respect to no new comments or areas for improvement. Pilot testers were provided with feedback from each round of testing to inform their evaluations. Authors were physicians or CHA study researchers and were not affiliated with the articles under evaluation. We instructed pilot testers to flag any item or subitem that they perceived as unclear, or inappropriate for further assessment by the steering group and re-evaluation by the panel if needed. However, we received positive feedback regarding the length, content, and user experience with the checklist. No items or subitems were flagged as inappropriate. Minor changes were made to the checklist including the phrasing of items, the order of items, and the formatting of the fillable document to optimize user experience with the checklist. No additional items or subitems were suggested. Saturation was reached after two rounds of pilot testing. Full details regarding our methodology can be found in our research protocol [7].

### Deviations from the protocol

Based on feedback from the expert multidisciplinary panel, we broadened the scope beyond LLMs to include any applications using generative AI due to the dynamically evolving nature of AI research in medicine. Moreover, two expert subgroups were assembled after the panel reviewed the candidate checklist items after the first consensus meeting. First, an expert generative AI subgroup met to evaluate and revise the terminology and checklist items used in this reporting guideline. Second, an expert data analysis subgroup reviewed checklist items related to statistical analysis. The results of both subgroups were presented to the expert panel and were reviewed for approval and discussed at subsequent panel consensus meetings. Finally, due to the complex nature of the conduct and reporting of CHA studies, we developed the checklist

items and accompanying diagram for CHART over three separate synchronous, 4-h panel consensus meetings rather than two, as initially planned in our protocol [7]. Further guidance and points of emphasis are detailed in the CHART Explanation and Elaboration article [14].

## Results

The CHART methodological diagram can be seen in Fig. 1. The CHART checklist consists of 12 items comprising 39 subitems for the complete and transparent reporting of CHA studies. Items relate to Title & Abstract (item 1), Introduction (item 2), Methods (items 3–9), Results (item 10), Discussion (item 11), and Open Science (item 12). Table 2 lists the CHART checklist items. Table 3 lists the CHART abstract checklist items.

The Delphi advisory committee and the expert panel both emphasized the importance of several checklist items. Specific examples are highlighted here, but the thorough reporting of all items listed in Table 2 is recommended. Delphi and panel members both voiced that authors must adequately identify the generative AI model and chatbot which they evaluated (items 3 and 4). This includes model identifiers, whether it is an open-source or proprietary model, and whether the model was novel or a base model (Table 2). Our expert stakeholders further stressed that authors must report the details involved during prompt engineering as well as the query strategy applied by investigators (item 5 and 6). This includes the process used to develop prompts, the members of the study team involved, and the dates and locations of queries (Table 2). Our panelists also underscored the necessity of explicitly defining a reference standard and describing the performance evaluation process (item 7). Stakeholders emphasized the importance of providing a sample size, which includes the number of independent responses from one or more generative AI-driven chatbot(s). Panelists also identified that the sample size of training data points may also be relevant if authors evaluate a novel or tuned model. Additionally, panelists stressed the importance of reporting the training data used, the ethical approval process undertaken, measures to safeguard the privacy of patient data, the permission or licensing obtained for the use of training data, and whether the training data can be accessed (item 12) (Table 3).

## Discussion

CHART was developed in accordance with the highest methodological standards through a comprehensive systematic review of CHA studies, a modified asynchronous Delphi process conducted by an international, multidisciplinary advisory committee, and three synchronous international, multidisciplinary expert panel consensus meetings [7]. Detailed rationale for each subitem are

described in our Explanation and Elaboration article [14]. The CHART checklist outlines essential items for the reporting of CHA studies which typically evaluate the performance of generative AI-driven chatbots when summarizing clinical evidence or providing health advice. At the time of writing, substantial advancements are being made in other forms of generative AI such as large multimodal models (LMMs), to which our reporting checklist—developed in the context of studies evaluating LLM performance—may not fully apply [15]. Thus, due to the rapidly evolving nature of these studies, a dynamic process must be in place for the monitoring and updating of this reporting guideline [16].

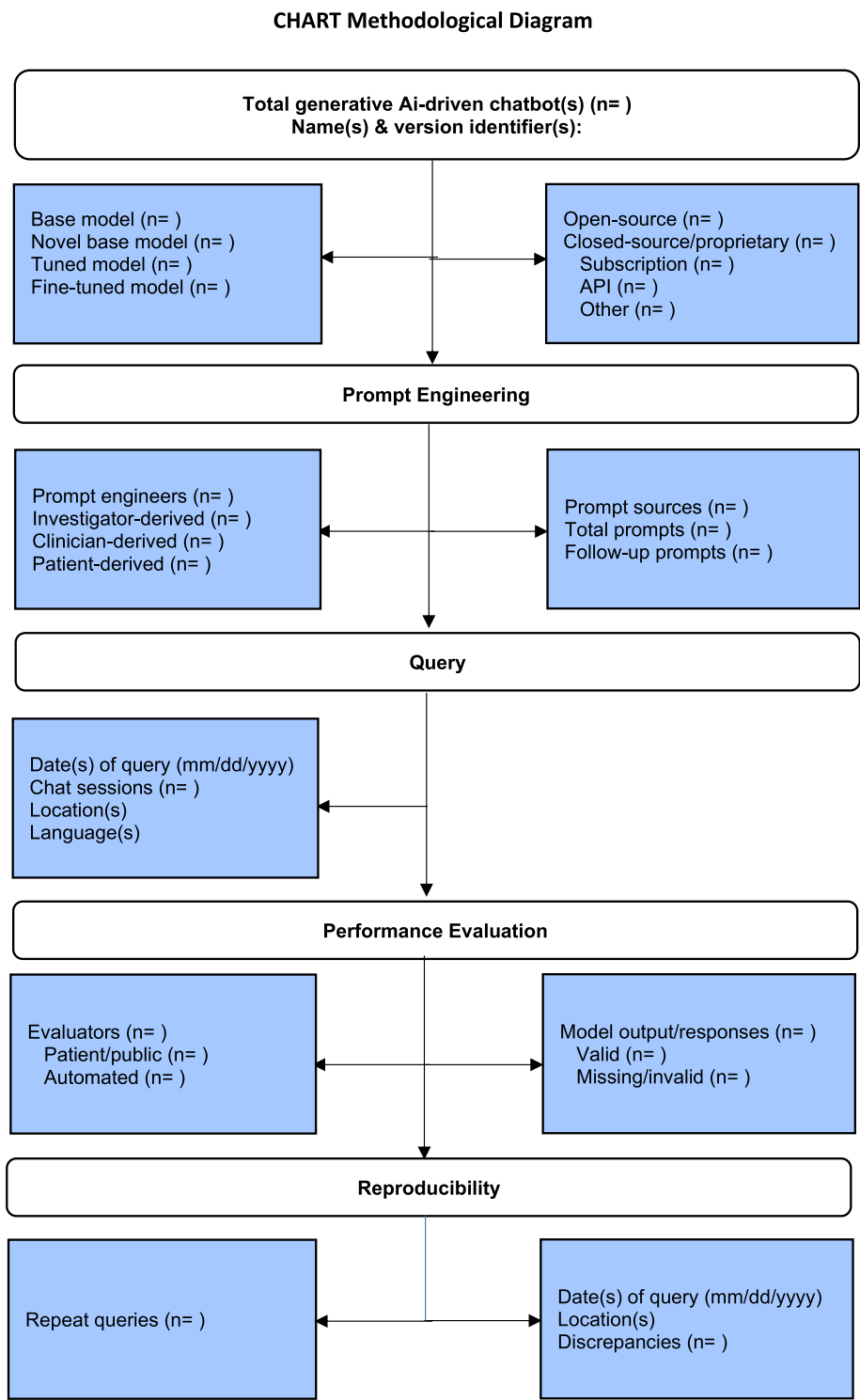
## Applicability and scope

The CHART checklist applies to CHA studies where generative AI-driven chatbots are queried and their responses are reported and evaluated. The CHART checklist does not apply to CHA studies applying randomization techniques (randomized controlled trials), nor studies that follow patients over time (prospective cohort studies). Future CHART extensions of relevant checklists for various study designs are planned, but in the interim authors are encouraged to apply both the CHART checklist and relevant reporting guidelines according to the appropriate study design such as CONSORT or STROBE [17, 18]. Authors using applications in the field of artificial intelligence more broadly (but not generative AI) are encouraged to use more generic reporting guidelines [13, 19, 20]. Authors using generative AI models for medical writing are encouraged to apply the CANGARU reporting guidelines, which are in development [21]. CHART applies to the current landscape of CHA studies and will evolve as a living reporting guideline.

## How to use CHART

We suggest that authors use the CHART checklist early in the writing of CHA studies to ensure all items in the checklist have been reported somewhere in their manuscript. Many of the recommendations in the CHART checklist have a natural order and sequence in a CHA study, but some may not. We do not prescribe a specific format or dictate where each individual reporting recommendation should appear in a CHA study, because this order might also depend on journal formatting policies. A downloadable and editable checklist can be found in the supplementary material. Authors are recommended to complete the checklist indicating the page number where each subitem has been reported. The completed checklist can then be submitted alongside the CHA study manuscript. A detailed Explanation and Elaboration paper accompanies the CHART checklist and explains why the reporting of each item is recommended [14].





**Fig. 1** The CHART Methodological Diagram

**Copyright protections and fair use doctrine**  
The accuracy of LLMs is significantly influenced by the nature of the data on which they were trained [10, 22].

This principle is the first of four according to the fair use doctrine, which are addressed throughout the CHART checklist as they relate to CHA studies. The first

**Table 2** CHART Checklist

HEADING	#	CHART CHECKLIST ITEM	Page #*
Title & Abstract			
Title	1a	State that the study is assessing one or more generative AI-driven chatbots for clinical evidence or health advice.	
Abstract/Summary	1b	Apply a structured format, if applicable.	
Introduction			
Background	2a	State the scientific background, rationale, and healthcare context for evaluating the generative AI-driven chatbot(s), referencing relevant literature when applicable.	
	2b	State the aims and research questions including the target audience, intervention, comparator(s), and outcome(s).	
Methods			
Model Identifiers	3a	State the name and version identifier(s) of the generative AI model(s) and chatbot(s) under evaluation, as well as their date of release or last update.	
	3b	State whether the generative AI model(s) and chatbot(s) are open-source or closed-source/proprietary.	
Model Details	4a	State whether the generative AI model was a base model or a novel base model, tuned model, or fine-tuned model.	
	4b	If a base model is used, cite its development in sufficient detail to identify the model.	
	4c	If a novel base model, tuned model, or fine-tuned model is used, describe the pre- and/or post-implementation/deployment data and parameters.	
Prompt Engineering	5a	Describe the evolution of study prompt development.	
	5ai	Describe the sources of prompts.	
	5aii	State the number and characteristics of the individual(s) involved in prompt engineering.	
	5aiii	Provide details of any patient and public involvement during prompt engineering.	
	5b	Provide study prompts.	
Query Strategy	6a	State route of access to generative AI model.	
	6b	State the date(s) and location(s) of queries for the generative AI-driven chatbot(s) including the day, month, and year as well as city and country.	
	6c	Describe whether prompts were input into separate chat session(s).	
	6d	Provide all generative AI-driven chatbot output/responses	
Performance Evaluation	7a	Define the ground truth or reference standard used to define successful generative AI-driven chatbot performance.	
	7b	Describe the process undertaken for generative AI-driven chatbot performance evaluation.	
	7bi	State the number and characteristics of team members involved in performance evaluation.	
	7bii	Provide details of any patients and public involvement during the evaluation process.	
	7biii	State whether evaluators were blinded to the identity of the generative AI-driven chatbot(s) under assessment.	
Sample Size	8	Report how the sample size was determined.	
Data Analysis	9a	Describe statistical analysis methods, including any evaluation of reproducibility of generative AI-driven chatbot responses.	
	9ai	Report the measures used for performance evaluation.	
Results			
	10a	Report the performance evaluation undertaken including the alignment between generative AI-driven chatbot output and ground truth or reference standard using quantitative or mixed methods approaches as applicable.	
	10b	For responses deviating from the ground truth or reference standard, state the nature of the difference(s).	
	10c	Report the evaluation for potentially harmful, biased, or misleading responses.	
Discussion			
	11a	Interpret study findings in the context of relevant evidence.	
	11b	Describe the strengths and limitations of the study.	
	11c	Describe the potential implications for practice, education, policy, regulation, and research.	
Open Science			
Disclosures	12a	Report any relevant conflicts of interest for all authors.	
Funding	12b	Report sources of funding and their role in the conduct and reporting of the study.	

**Table 2** (continued)

HEADING	#	CHART CHECKLIST ITEM	Page #*
Ethics	12c	Describe the process undertaken for ethical approval.	
	12ci	Describe the measures taken to safeguard data privacy of patient health information, as applicable.	
	12cii	State whether permission/licensing was obtained for the use of original, copyrighted data.	
Protocol	12d	Provide a study protocol.	
Data availability	12e	State where study data, code repository, and model parameters can be accessed.	

**Table 3** The CHART Abstract Checklist

HEADING	CHART Checklist #	ITEM	Page #
Background	2a	State the scientific background, rationale, and healthcare context for evaluating the generative AI-driven chatbot(s), referencing relevant literature when applicable.	
	2b	State the aims and research questions including the target audience, intervention, comparator(s), and outcome(s).	
Methods			
Model Identifiers	3a	State the name and version identifier(s) of the generative AI model(s) and chatbot(s) under evaluation, as well as their date of release or last update.	
	3b	State whether generative AI model(s) and chatbot(s) are open-source versus closed-source/proprietary.	
Model Details	4a	State whether the generative AI model was a base model or a novel base model, tuned model, or fine-tuned model.	
Prompt Engineering	5a	Describe the evolution of study prompt development.	
	5ai	Describe the sources of prompts.	
	5aii	State the number and characteristics of the individual(s) involved in prompt engineering.	
	5aiii	Provide details of any patient and public involvement during prompt engineering.	
Query Strategy	6a	State route of access to generative AI model.	
	6b	State the date(s) and location(s) of queries for the generative AI-driven chatbot(s) including the day, month, and year as well as city and country.	
Performance Evaluation	7a	Define the ground truth or reference standard used to define successful generative AI-driven chatbot performance.	
	7b	Describe the process undertaken for the performance evaluation of the generative AI-driven chatbot(s).	
Sample Size	8	Report how the sample size was determined.	
Data Analysis	9a	Describe statistical analysis methods, including any evaluation of reproducibility of generative AI-driven chatbot responses.	
Results			
	10a	Report the alignment between generative AI-driven chatbot output and ground truth or reference standard using quantitative or mixed methods approaches as applicable.	

principle refers to the purpose and character of use of the model [10]. The second principle is the nature of the original training data [10, 23]. While many LLMs will be trained on non-medical data, it is essential that factual, evidence-based information must be prioritized in the healthcare setting [10]. The third principle pertains to the amount and substantiality of original material used to train the generative AI model [10], and clarity regarding the origin of training data and permission or license to use content or data protected by copyright is

recommended. Finally, the fourth principle relates to the impact on original work, where generative AI models may be trained with copyrighted data [10]. We address these principles in the CHART checklist by encouraging authors to state the purpose of the study, and whether they are evaluating a pre-existing base model rather than one that is a novel base model, a tuned model, or a fine-tuned model (items 3 and 4). The CHART checklist promotes open science practices and calls for authors to share their code and training datasets to optimize



transparency and mitigate uncertainty over data provenance (item 12e). The CHART checklist further uses an evidence-based approach by encouraging authors to state the source of their prompts, their definition of successful model/bot performance, and the process behind performance evaluation (item 5, item 7). The CHART checklist recommends that authors state whether permission or license was obtained by investigators for use of the original work (item 12cii). Readers may also identify the presence of copyrighted data as authors share their coding and training data (item 12e).

### Bias and patient safety

In the setting of model development, the output of generative AI models such as LLMs are further impacted by the presence of bias in their training datasets [10]. This introduces the risk of LLMs producing misleading or harmful information when applied for the purposes of patient care. These biases may pertain to many factors including, but not limited to race or ethnicity, sex or gender, language, and culture [24, 25]. This risk further highlights the importance of the Open Science checklist item (item 12) in CHART because the risk of bias from data used to develop LLM-driven chatbots may be identified and/or mitigated by open coding and training data sharing [25]. Furthermore, data used to train generative AI models may pose a threat to data security and patient privacy. The use of identifiable patient data during model training is of particular concern, as sensitive information may be inadvertently disclosed in the absence of appropriate data security measures [10, 26]. The risk for data breaches must be met accordingly with robust cybersecurity measures [10]. This concept underscores the importance of the CHART checklist item related to steps taken to ensure safeguarding of patient health information (item 12ci). The push for clinically integrating generative AI models necessitates human oversight of the ethical and safe inclusion of patients and their health information to provide guidance for the safe conduct of CHA studies [27, 28]. Although we recognize the importance of making advancements by including patients in CHA studies to develop more patient-centered studies (item 5biii, item 7bii), we encourage authors to report whether ethics approval was obtained in these instances for the responsible conduct of their study (item 12c).

### Monitoring and updates

This reporting guideline will follow and adapt the traditional methodology for a living clinical practice guideline [16]. The update interval for this reporting guideline will apply to individual checklist items, rather than the entire guideline [16]. Core members of the

steering group will perform a systematic search of the literature to continuously survey the literature per living guideline best practices [16] and will meet to discuss any relevant developments in the generative AI field every 6 months for the first 2 years (until 2026). If important changes occur sooner, the group will meet ad hoc as needed. The timing for monitoring and updating the guideline will be reviewed and revised at the time of the next reporting guideline update or by the end of 2026, whichever occurs sooner.

Furthermore, a living expert panel consisting of 14 expert panel members was selected following the third expert panel consensus meeting in accordance with living guideline best practices [16], and comprised panel members committed to making themselves available to meet virtually at very short notice [16]. Living expert panel members represent backgrounds stemming from medicine, epidemiology, data science, health research methodology, reporting guideline methodology, and statistics. If no changes to the reporting guideline are warranted within a given year, the living expert panel will be updated with the activities of the core steering group and will be alerted to any relevant literature or topics within generative AI to monitor and be aware of. This update will occur at a minimum of once per year at a meeting between the core members of the steering group and the living expert panel. Finally, living peer reviewers will be selected following the peer review process for the CHART statement and Elaboration and Explanation articles [16]. They will similarly be provided with an annual update, but will only be contacted if checklist items must be updated. If new candidate checklist items or revisions to existing items are identified by the core members of the steering group, the living expert panel will be convened at its earliest convenience to review the relevant literature. In alignment with living guideline best practices [16], the minimum threshold will be set at 90% agreement among living expert panel members for changing checklist items to mitigate the risk of false positives inherent to frequent updates, while avoiding an excessively high threshold [16]. If applicable, the updated manuscript will be co-published in relevant journals with interest.

### Target users and implications for stakeholders

CHART applies to individuals performing and reviewing CHA studies such as study investigators, peer reviewers, and journal editors for academic purposes, as well as the wider readership of CHA studies including clinicians, statisticians, generative AI researchers, regulatory experts, ethicists, research methodologists, policy makers, hospital managers, funders, patients, and the wider

public. To promote the transparent reporting of CHA studies, we call for clinical journals to adopt CHART: a comprehensive reporting standard developed with high methodological rigor. The main barrier that we anticipate to CHART uptake is the failure to reach the appropriate audience. Therefore, this reporting guideline will be listed on the EQUATOR Network website, and we will disseminate the publication of this reporting guideline widely. CHART will also be presented at peer-reviewed meetings across various medical specialties to optimize the dissemination and reach of the checklist and accompanying diagram. Finally, we will develop a website to house fillable versions of the abstract checklist, the full checklist, and the methodological diagram, which can be found in supplementary Appendices 2–4 of this publication to facilitate the application of CHART by CHA researchers.

Following the publication of previous reporting guidelines, it has been shown that the reporting quality of applicable studies improve [29, 30]. As investigators and journals apply CHART and the completeness of reporting of CHA studies improve, higher quality studies may be produced. Researchers, ethicists, clinicians, and regulators in the clinical generative AI community must then turn toward the validation of generative AI-driven chatbots for the purposes of providing health advice [10]. This may include the prioritization of standardized quality validation metrics, clarifying the role of human involvement in validation studies, validation methodology [31], and the reporting of validation results using the CHART tool. Regulators must further look toward data sensitivity and privacy, ensuring that data security measures are put in place by generative AI developers according to risk category [10]. Funders must invest in the development of high-quality benchmarking and validation studies, as well as highly rigorous CHA studies in the context of the healthcare setting of interest. Funders may also encourage applicants to include a research plan in alignment with the CHART checklist. With studies exhibiting greater transparency and improved methodological rigor, clinicians, patients, and the public will develop progressively increased trust in the clinical integration of generative AI-driven chatbots.

Finally, quality appraisal tools do not exist for CHA studies and remains a future area of study. CHART is a reporting guideline rather than a critical appraisal tool. Still, we hope that attention to CHART's core checklist items will indirectly improve the methodologic rigor of studies in this field [32]. As high-quality evidence builds, the path forward for integrating generative AI into the clinical practice environment will become clearer for both hospital managers and policy makers.

## Conclusion

The transparent reporting of CHA studies is crucial for their interpretation as we move toward the clinical integration of AI technologies. The CHART reporting guideline consists of a 12-item checklist and corresponding methodological diagram to support key stakeholders including clinicians, researchers, editors, peer reviewers, and readers in reporting, understanding, and interpreting the findings of CHA studies.

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12916-025-04274-w>.

- Additional file 1: Appendix 1 Candidate checklist items.
- Additional file 2: Appendix 2 Fillable CHART checklist.
- Additional file 3: Appendix 3 Fillable CHART abstract checklist.
- Additional file 4: Appendix 4 Fillable methodological diagram.
- Additional file 5: Appendix 5 Panel members.

## Acknowledgements

The authors would like to thank the First Cut competition organizers and the Postgraduate Medical Education Committee at McMaster University for financially supporting the development of this project. The authors would also like to thank members of the Advisory Committee for their invaluable time and effort devoted to the development of the Chatbot Assessment Reporting Tool.

## Authors' contributions

BH, DC, AWC, CL, and GG contributed to the conception and design of the work. All authors contributed to data analysis, interpretation, manuscript drafting, revising, and approve of the final version to be published. BH acted as guarantor and is responsible for the overall content (as guarantor), and all authors are accountable for the accuracy of the checklist and methodological diagram.

## Funding

The Chatbot Assessment Reporting Tool (CHART) was funded by the *First Cut* competition and the Postgraduate Medical Education Committee at McMaster University. Neither funding sources were involved in the design, conduct, or reporting of this reporting guideline.

## Data availability

No datasets were generated or analysed during the current study.

## Declarations

### Ethics approval and consent to participate

Ethics approval was submitted to and waived by the Hamilton Integrated Research Ethics Board (HIREB #17025).

### Consent for publication

Not applicable.

### Competing interests

All authors have completed the ICMJE uniform disclosure form at [www.icmje.org/disclosure-of-interest/](http://www.icmje.org/disclosure-of-interest/) and declare: GSC is a National Institute for Health and Care Research (NIHR) Senior Investigator. The views expressed in this article are those of the author(s) and not necessarily those of the NIHR, or the Department of Health and Social Care; AJT has received funding from HealthSense to investigate evidence-based medicine applications of large language models. PM is the co-founder of BrainX LLC; AS has received research funding from the Australian government and is co-founder of BantingMed Pty Ltd; DS is the Acting Deputy Editor for the *Lancet Digital Health*; MM has

received research funding from The Hospital Research Founding Group; TF sits on the executive committee of MDEpiNet; HF is a Senior Executive Editor for The Lancet; CL is the Editor in Chief of Annals of Internal Medicine; AF is Executive Managing Editor and Vice President, Editorial Operations, JAMA and The JAMA Network; TF and EL are journal editors for the BMJ; RA is the Editor in Chief of International Journal of Surgery; GS is an Executive Editor of Artificial Intelligence in Medicine; SL is a paid consultant for Astellas; DP has received research funding from the Italian Ministry of University and Research; MO is a paid consultant for Theator; TA, POV, GG are board member of the MAGIC Evidence Ecosystem Foundation ([www.magicproject.org](http://www.magicproject.org)), a non-for profit organization, which conducts research and evidence appraisal and guideline methodology and implementation, and which provides a authoring and publication software (MAGICapp) for evidence summaries, guidelines and decision aids.

# Author details

<sup>1</sup>Division of General Surgery, Department of Surgery, McMaster University, Hamilton, Canada. <sup>2</sup>UK EQUATOR Centre, University of Oxford, Oxford, UK. <sup>3</sup>Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology & 401 Musculoskeletal Sciences, Botnar Research Centre, University of Oxford, Oxford, UK. <sup>4</sup>Department of Biomedical Informatics and Data Science, Yale University School of Medicine, New Haven, USA. <sup>5</sup>Nuffield Department of Clinical Neurosciences, Medical Sciences Division, University of Oxford, Oxford, UK. <sup>6</sup>JAMA and JAMA Network, American Medical Association, Chicago, USA. <sup>7</sup>Department of Health Research Methods, Evidence, and Impact; Department of Medicine, McMaster University, Hamilton, Canada. <sup>8</sup>USC Institute of Urology and Catherine and Joseph Aresty Department of Urology, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>9</sup>Artificial Intelligence Center at USC Urology, USC Institute of Urology, University of Southern California, Los Angeles, CA, USA. <sup>10</sup>Sports Medicine Center, West China Hospital, Sichuan University, Chengdu, China. <sup>11</sup>Department of Orthopedics and Orthopedic Research Institute, West China Hospital, Sichuan University, Chengdu, China. <sup>12</sup>Duke-NUS Medical School, National University of Singapore, Singapore, Singapore. <sup>13</sup>Cleveland Clinic, Case Western Reserve University, Cleveland, USA. <sup>14</sup>Department of Medicine, Women's College Research Institute, University of Toronto, Toronto, Canada. <sup>15</sup>Annals of Internal Medicine, American College of Physicians, Philadelphia, USA. <sup>16</sup>American College of Physicians, Philadelphia, USA. <sup>17</sup>Department of Public Health, University of Naples Federico II, Naples, Italy. <sup>18</sup>Director, Office of Science Dissemination, Office of Science, Centers for Disease Control and Prevention, Atlanta, GA, USA. <sup>19</sup>Department of General Surgery, Papageorgiou General Hospital, Thessaloniki, Greece. <sup>20</sup>British Psychological Society, University of Plymouth, Plymouth, UK. <sup>21</sup>Innovation Support Unit, Department of Family Practice, University of British Columbia, Vancouver, Canada. <sup>22</sup>Patient SME, Independent Cybersecurity Professional, London, UK. <sup>23</sup>The BMJ, London, UK. <sup>24</sup>Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA. <sup>25</sup>Department of Neurology, Brigham and Women's Hospital, Boston, MA, USA. <sup>26</sup>International Journal of Surgery, London, UK. <sup>27</sup>Eworkflow Ltd, London, UK. <sup>28</sup>Department of Oncology, McMaster University, Hamilton, Canada. <sup>29</sup>Hospital Clinico San Carlos, Instituto de Investigación Sanitaria San Carlos, Facultad de Medicina Universidad Complutense de Madrid, Madrid, Spain. <sup>30</sup>The George Institute for Global Health; Tyree Institute of Health Engineering, UNSW Engineering; School of Population Health, UNSW Medicine and Health, Sydney, Australia. <sup>31</sup>Hardian Health, London, UK. <sup>32</sup>Centre for Journalology, Ottawa Hospital Research Institute, Ottawa, Canada. <sup>33</sup>Digestive Diseases Institute, Cleveland Clinic, Cleveland, OH, USA. <sup>34</sup>Postgraduate Institute of Medical Education and Research, Chandigarh, India. <sup>35</sup>University of Maribor, Maribor, Slovenia. <sup>36</sup>Australian Institute for Machine Learning (AIML), Adelaide, Australia. <sup>37</sup>Phelix AI, Toronto, Canada. <sup>38</sup>Università Politecnica delle Marche, Clinica di Chirurgia Generale e d'Urgenza, Ancona, Italy. <sup>39</sup>Institute for Evidence in Medicine, Medical Center & Faculty of Medicine, University of Freiburg, Freiburg im Breisgau, Germany. <sup>40</sup>Cochrane Germany, Cochrane Germany Foundation, Freiburg, Germany. <sup>41</sup>MAGIC Evidence Ecosystem Foundation, Oslo, Norway. <sup>42</sup>University Hospitals of Geneva, Geneva, Switzerland. <sup>43</sup>The Lancet Digital Health, London, UK. <sup>44</sup>The Lancet, London, UK. <sup>45</sup>NIHR Clinical Lecturer, Health Organisation, Policy, Economics (HOPE), Centre for Primary Care & Health Services Research, The University of Manchester, Manchester, UK. <sup>46</sup>Senior Visiting Fellow, LSE Health, London School of Economics and Political Science, Manchester, UK. <sup>47</sup>New York University Langone Health, New York City, USA.

<sup>48</sup>College of Medicine and Health, University of Birmingham, Birmingham, UK. <sup>49</sup>School of Global Public Health, New York University, New York City, USA.

Received: 11 June 2025 Accepted: 10 July 2025

Published online: 01 August 2025

# References

- Kolbinger FR, Veldhuizen GP, Zhu J, Truhn D, Kather JN. Reporting guidelines in medical artificial intelligence: a systematic review and meta-analysis. *Commun Med*. 2024;4:1.
- Han R, Acosta JN, Shakeri Z, Ioannidis JPA, Topol EJ, Rajpurkar P. Randomised controlled trials evaluating artificial intelligence in clinical practice: a scoping review. *Lancet Digit Health*. 2024;6:e367–73.
- Huo B, Cacciamani GE, Collins GS, McKechnie T, Lee Y, Guyatt G. Reporting standards for the use of large language model-linked chatbots for health advice. *Nat Med*. 2023;29:2988.
- Huo B, McKechnie T, Ortenzi M, Lee Y, Antoniou S, Mayol J, et al. Dr. GPT will see you now: the ability of large language model-linked chatbots to provide colorectal cancer screening recommendations. *Health Technol*. 2024;14:463–9.
- Huo B, Marfo N, Sylla P, Calabrese E, Kumar S, Slater BJ, et al. Clinical artificial intelligence: teaching a large language model to generate recommendations that align with guidelines for the surgical management of GERD. *Surg Endosc*. 2024;38:5668–77.
- Huo B, Boyle A, Marfo N, Tangamornsuksan W, Steen JP, McKechnie T, et al. Large language models for chatbot health advice studies: a systematic review. *JAMA Netw Open*. 2025;8:e2457879.
- The CHART Collaborative. Protocol for the development of the Chatbot Assessment Reporting Tool (CHART) for clinical advice. *BMJ Open*. 2024;14:e081155.
- Moher D, Schulz KF, Simera I, Altman DG. Guidance for developers of health research reporting guidelines. *PLoS Med*. 2010;17:e1000217.
- Collins GS, Moons KGM, Dhiman P, Riley RD, Beam AL, Van Calster B, et al. TRIPOD+AI statement: updated guidance for reporting clinical prediction models that use regression or machine learning methods. *BMJ*. 2024;384:q902.
- Ong JCL, Chang SYH, William W, Butte AJ, Shah NH, Chew LST, et al. Ethical and regulatory challenges of large language models in medicine. *The Lancet Digital Health Elsevier Ltd*. 2024;6:e428–32.
- Altman DG, Simera I, Hoey J, Moher D, Schulz K. EQUATOR: reporting guidelines for health research. *Open Med*. 2008;371:1149–50.
- Munn Z, Peters MDJ, Stern C, Tufanaru C, McArthur A, Aromataris E. Systematic review or scoping review? Guidance for authors when choosing between a systematic or scoping review approach. *BMC Med Res Methodol*. 2018;18:143.
- Liu X, Rivera SC, Moher D, Calvert MJ, Denniston AK. Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: the CONSORT-AI Extension. *BMJ*. 2020;370:m3164.
- The CHART Collaborative. Reporting guidelines for chatbot health advice studies: explanation and elaboration for the Chatbot Assessment Reporting Tool (CHART). *BMJ*. 2025;390:e083305.
- Yin S, Fu C, Zhao S, Li K, Sun X, Xu T, et al. A survey on multimodal large language models. 2023. <http://arxiv.org/abs/2306.13549>. Accessed 24 Jun 2025.
- Akl EA, Meerpohl JJ, Elliott J, Kahale LA, Schünemann HJ, Agoritsas T, et al. Living systematic reviews: 4. Living guideline recommendations. Vol. 91. *J Clin Epidemiol*. 2017;91:47–53.
- Begg C, Cho M, Eastwood S, Horton R, Moher D, Olkin I, et al. Improving the quality of reporting of randomized controlled trials. The CONSORT statement. *JAMA*. 1996;276:637–9.
- von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP. Strengthening the reporting of observational studies in epidemiology (STROBE) statement: guidelines for reporting observational studies. *BMJ*. 2007;335:806–8.
- Rivera SC, Liu X, Chan AW, Denniston AK, Calvert MJ. Guidelines for clinical trial protocols for interventions involving artificial intelligence: the SPIRIT-AI Extension. *The BMJ*. 2020;370:m3210.
- Vasey B, Nagendran M, Campbell B, Clifton DA, Collins GS, Denaxas S, et al. Reporting guideline for the early-stage clinical evaluation of

- decision support systems driven by artificial intelligence: DECIDE-AI. *Nat Med*. 2022;28:924–33.
21. Cacciamani G, Gill I, Collins G. ChatGPT: standard reporting guidelines for responsible use. *Nat*. 2023;618:1–1.
  22. Xie SM, Pham H, Dong X, Du N, Liu H, Lu Y, et al. DoReMi: optimizing data mixtures speeds up language model pretraining. 2023. <https://arxiv.org/abs/2305.10429>. Accessed 24 June 2025.
  23. Ng FYC, Thirunavukarasu AJ, Cheng H, Tan TF, Gutierrez L, Lan Y, et al. Artificial intelligence education: an evidence-based medicine approach for consumers, translators, and developers. *Cell Rep Med*. 2023;4: 101230.
  24. Li H, Moon JT, Purkayastha S, Celi LA, Trivedi H, Gichoya JW. Ethics of large language models in medicine and medical research. *Lancet Dig Health*. 2023;5:e333–5.
  25. Health TLD. Large language models: a new chapter in digital health. *Lancet Dig Health*. 2024;6: e1.
  26. Thirunavukarasu AJ, Ting DSJ, Elangovan K, Gutierrez L, Tan TF, Ting DSW. Large language models in medicine. *Nat Med*. 2023;29:1930–40.
  27. Haltaufderheide J, Ranisch R. The ethics of ChatGPT in medicine and healthcare: a systematic review on large language models (LLMs). *npj Digit Med*. 2024;7:183.
  28. Thirunavukarasu AJ. Large language models will not replace health-care professionals: curbing popular fears and hype. *J R Soc Med*. 2023;116:181–2.
  29. Kane RL, Wang J, Garrard J. Reporting in randomized clinical trials improved after adoption of the CONSORT statement. *J Clin Epidemiol*. 2007;60:241–9.
  30. Turner L, Shamseer L, Altman DG, Schulz KF, Moher D. Does use of the CONSORT statement impact the completeness of reporting of randomised controlled trials published in medical journals? *A Cochrane review Syst Rev*. 2012;1:60.
  31. de Hond A, Leeuwenberg T, Bartels R, van Buchem M, Kant I, GM Moons K, et al. From text to treatment: the crucial role of validation for generative large language models in health care. *Lancet Digital Health*. 2024;6:e441–3.
  32. Logullo P, Maccarthy A, Kirtley S, Collins GS. Reporting guideline checklists are not quality evaluation forms: they are guidance for writing. *Health Sci Rep*. 2020;3: e165.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.