

1

2

3

4 **A scoping review protocol examining the application of large language models in**

5 **healthcare education and public health learning spaces**

6 Henry Ndukwe¹, Emmanuel Otukpa²

7

8 1. School of Pharmacy and Medical Sciences, Griffith University, Australia.

9 2. Health and Wellbeing theme, African Population Health Research Center (APHRC),

10 Nairobi Kenya.

11

12

13

14

15 ***Corresponding Author**

16 E-mail: eotukpa@aphrc.org

17

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

18 **Abstract**

19 **Objective:** Through this scoping review, we aim to explore and synthesize existing
20 knowledge and evidence on the learning approaches for incorporating LLMs into healthcare
21 education and public health learning spaces. Specifically, we will attempt to investigate
22 methods for auditing prompts for accuracy; tailoring prompts to improve task-specific
23 accuracy and utility; and exploring how end-user feedback is used to refine and optimize
24 LLM prompts over time. This review will provide a comprehensive understanding of how
25 LLMs are being tailored and improved in these fields, contributing to the development of
26 evidence-based strategies for their implementation. It will also identify areas for future
27 research and innovation.

28 **Introduction:** The increasing integration of large language models (LLMs) into healthcare
29 and public health practice and research highlights their potential to revolutionize service
30 delivery, decision-making, and patient care. Despite these advancements, understanding how
31 LLMs can be effectively tailored, audited, and refined for healthcare-specific tasks remains a
32 critical area of inquiry. Key issues include, the accuracy of generated information, and their
33 relevance to the medical and public health fields.

34 **Inclusion criteria:** Inclusion criteria will focus on studies addressing LLM applications in
35 healthcare and public health, prompt engineering techniques, prompt auditing methods, and
36 processes geared towards integrating user feedback. Articles that do not focus on healthcare
37 or public health contexts and lack relevance to LLM learning approaches will be excluded.

38 **Methods:** The review is guided by the JBI methodology for scoping reviews complemented
39 by updates from Levac et al. Databases including PubMed, Scopus, IEEE Xplore, and Web of
40 Science will be searched for peer-reviewed articles, conference proceedings, and grey

41 literature published in English and French from 2015 to 2025. Data extraction will include
42 information on study characteristics, LLM models, prompt engineering strategies, auditing
43 methodologies, and user feedback mechanisms. We will synthesize to identify trends, gaps,
44 and best practices in leveraging LLMs to generate baseline data for auditing prompts that
45 optimize AI learning and education needs in the healthcare and public health sector.

46 **Keywords:** LLMs, Artificial Intelligence, healthcare, public health, prompt engineering,
47 auditing, user feedback, artificial intelligence.

48 **Introduction**

49 **Background and context**

50 Large language models (LLMs) represent a novel innovation in computing and the part of
51 what constitutes reactive artificial intelligence, offering the potential to revolutionize
52 healthcare and public health through its natural language processing and language generation
53 [1]. These models, exemplified by Generative Pre-trained Transformers (GPT) and its
54 counterparts, can process and generate human-like text, enabling applications rooted in
55 clinical decision-making, patient education, administrative efficiency, and medical and public
56 health research [2,3]. As healthcare systems increasingly explore digital solutions to enhance
57 care delivery, the role of LLMs is becoming more prominent. However, maximizing their
58 potential while ensuring reliability and safety requires a nuanced understanding of their
59 integration into healthcare-specific contexts [4].

60 The integration of LLMs into healthcare and public health education presents both
61 opportunities and challenges. On the one hand, LLMs offer unprecedented capabilities for
62 personalized learning, real-time clinical decision support, and enhanced access to medical
63 knowledge and public health data [5]. They can assist in creating customized educational

64 content, simulating clinical scenarios, and providing immediate feedback to learners using
65 artificial intelligence (AI) models. And on the other hand, LLMs may help address the
66 growing burden of healthcare-information overload and inaccurate online searching by
67 synthesizing complex public and medical health prompts to illicit accurate and digestible
68 output formats [4,5].

69 Despite these potential benefits, several challenges arise when considering the incorporation
70 of LLMs into healthcare and public health education. The traditional learning framework has
71 a shifting focus from scarcity towards the ubiquity of information via LLMs; fostering trust
72 that would ensure user safety on a continuous basis. Future strategies that optimize efficiency
73 with the leverage of AI-driven tools would require adaptable mechanism of feedback and
74 continuous improvement. Ethical concerns surrounding data privacy, bias in AI systems, and
75 the risk of over-reliance on technology must be carefully addressed [4,6]. There are also
76 questions around the quality and accuracy of information generated by LLMs, particularly in
77 rapidly evolving (as well as niche) medical fields [6]. Furthermore, educators and
78 policymakers face the task of developing appropriate frameworks for integrating LLMs into
79 curricula for learning and teaching or clinical practice modules [6,7].

80 **Challenges and gaps in knowledge**

81 Despite the promise of having a standard output scoring system for generative AI prompts,
82 challenges remain in deploying LLMs effectively in healthcare learning space [8]. Key
83 concerns include the accuracy of generated outputs, biases embedded in models, and their
84 contextual relevance to specific healthcare needs [8,9]. Addressing these challenges
85 necessitates the development of precise prompt engineering techniques to tailor LLM
86 responses for specialized tasks. Furthermore, auditing LLM outputs for accuracy, fairness,
87 and effectiveness is critical to building trust and ensuring equitable healthcare delivery

88 [10,11]. However, the methodologies and best practices for conducting such audits in
89 healthcare and public health contexts are not well-documented [11,12].

90 **Role of user feedback**

91 User feedback, particularly from healthcare and allied professionals, plays a vital role in
92 refining LLM performance[14]. Feedback mechanisms provide insights into how LLMs can
93 better align with clinical workflows and patient care priorities. While some studies have
94 explored integrating user feedback in LLM applications, there remains a lack of clarity on
95 how this process is operationalized and its impact on improving model outputs over time
96 [10,14].

97 **Rationale for the Scoping review**

98 Given the rapid evolution of LLMs and their potential adoption in healthcare, a
99 comprehensive synthesis of existing knowledge is essential to understand ways for proper
100 prompt engineering for generative AI. The ubiquitous nature of generative AI output poses a
101 persistent risk of misinformation premised around LLMs' learning algorithms, which are
102 built to optimize outputs based on user input (prompts) [15,16]. It is crucial to understand
103 possible oversight strategies to mitigate AI-information overload in learning environments.
104 To that end, this review will explore how LLMs have been deployed and audited in
105 healthcare education and public health learning spaces. By investigating the methods used to
106 audit LLM prompts; strategies to enhance task-specific accuracy and utility, and the
107 integration of user feedback can be proposed to refine LLM outputs within healthcare,
108 clinical and public health settings.

109 **Objectives**

110 In this scoping review, we aim to explore the current state of LLM integration in healthcare
111 and public health education, identify best practices, and highlight areas requiring further
112 research and development. We seek to provide actionable insights for researchers, healthcare
113 professionals, and policymakers to optimize the use of LLMs in healthcare and public health
114 education. The identification of best practices and knowledge gaps, it will contribute to
115 advancing the safe and effective implementation of LLM technologies in these critical fields

116 **Review questions**

- 117 1. What methods are used to audit prompts for accuracy in healthcare and public health
118 LLMs?
- 119 2. How are prompts tailored to specific tasks to improve the accuracy and utility of LLM
120 outputs?
- 121 3. How have end-user feedback been integrated into the process of refining LLM prompts
122 over time?

123 **Materials and Methods**

124 We intend to conduct this scoping review in accordance with the JBI methodology for
125 scoping reviews [17]. The methodology entails a systematic approach to searching, screening,
126 and reporting that include the following stages: (1) identification of the research question (s);
127 (2) identification of relevant databases and studies; (3) selection of studies; (4) data
128 extraction; (5) interpretation, summarization and dissemination of the results.

129 **Inclusion criteria**

130 We intend to focus on studies whose focus is on the application or development of large
131 language models (LLMs) usage and application/integration in healthcare or public health
132 contexts. Examples of relevant contexts include clinical decision support, patient education,
133 administrative processes, and public health interventions and education modules. In terms of
134 the scope of LLM integration, the studies must explore learning approaches for integrating
135 LLMs, including prompt engineering, auditing methodologies, and user feedback
136 mechanisms.

137 Specific publication inclusion criteria will focus on peer-reviewed articles, conference
138 proceedings, and grey literature (e.g., reports, white papers). Additionally, all study designs
139 are eligible, including experimental, observational, qualitative, and mixed-methods studies.
140 Studies focusing on frameworks, methodologies, and case studies will also be considered. We
141 will also consider any study published in English or French language. And with a timeframe
142 for inclusion from 2015 to 2025, reflecting a period of rapid advancements in LLM
143 technology and its applications.

144 In terms of relevance, studies must address at least one of the following core areas:

- 145 • Methods for auditing LLM prompt for accuracy, fairness, and effectiveness.
- 146 • Techniques for tailoring LLM prompts to healthcare-specific tasks.
- 147 • Approaches for integrating user feedback from healthcare professionals to refine and
148 improve LLM outputs.

149

150 **Search strategy**

151 Our search strategy will aim to locate both published and unpublished studies. A three-step
152 search strategy will be utilized in this review. First, we will conduct an initial limited search
153 of MEDLINE (PubMed) and CINAHL (EBSCO) to identify articles on the topic. The text
154 words contained in the titles and abstracts of relevant articles, and the index terms used to
155 describe the articles will be used to develop a full search strategy for reporting the name of
156 the relevant databases/information sources (*see Appendix 1*). The search strategy, including
157 all identified keywords and index terms, will be adapted for each included database and/or
158 information source. The reference list of all included sources of evidence will be screened for
159 additional studies. We will search relevant peer-reviewed, English and French-language
160 articles published between January 1, 2015, and January 31, 2025, without methodological
161 restrictions, in several electronic databases, as well as sources with broad specificity (Web of
162 Science and Google Scholar). Searches will extend to grey literature sources, including
163 institutional reports, white papers, and preprints on platforms like arXiv.

164 **Source of evidence selection (databases)**

165 Searches will be conducted in electronic databases for bibliographic sources including
166 Medline, Embase, and Scopus. We will include varied electronic data sources including Web
167 of Science, Google Scholar and IEEE Xplore.

168 **Search terms**

169 We will employ adjacency search and combination of keywords and Medical Subject
170 Headings (MeSH) terms will be used, including: "*large language models*", "*prompt*
171 "*engineering*", "*healthcare*", "*public health*", "*auditing methods*", "*user feedback*", "*artificial*
172 "*intelligence*". The Boolean operators (AND, OR, NOT) will combine the above terms to

173 refine quantity and quality of search hits. A detailed search strategy is provided in Appendix
174 1.

175 **Study selection**

176 We will use the online tool Covidence®, which allows for simultaneous title, abstract and full
177 text article reviews; to screen through articles to export included titles to Excel ® for
178 analysis. Two researchers will independently assess articles for inclusion by screening the
179 titles, abstracts, and full texts of studies returned through the search process. Where there are
180 disagreements between the two independent reviewers on the eligibility of a paper for
181 inclusion, a third reviewer will adjudicate using same inclusion criteria to resolve the conflict.

182 **Data extraction**

183 We will use a standardized data charting form (*Appendix 2*) to extract relevant data from
184 included studies. The following details will be extracted:

- 185 1. **Study Characteristics:** Title, authors, year, and country of publication.
- 186 2. **LLM Details:** Type of LLM, application context, and specific tasks addressed.
- 187 3. **Methodologies:** Techniques for prompt engineering, auditing methods, and feedback
188 mechanisms.
- 189 4. **Outcomes:** Measures of effectiveness, fairness, and task-specific utility.
- 190 5. **User Feedback:** Processes for incorporating feedback from healthcare professionals
191 and the impact on model refinement.

192 **Data synthesis and presentation**

193 We will analyze the data using descriptive statistics and thematic analysis, with results
194 organized in tables and charts and presented into themes that reflect the review objectives.
195 Our thematic analysis will identify patterns and trends in the application of LLMs.

196 Quantitative data, where applicable, will be summarized descriptively. We will also conduct
197 a narrative synthesis to integrate findings across studies, focusing on methods, outcomes, and
198 identified gaps.

199 **Ethics and dissemination**

200 Ethical approval is not required because primary data collection is not involved in this study
201 but rather analyzing both published and grey literature. However, the findings of this study
202 will be disseminated through peer-reviewed publications and conferences as well as in
203 relevant stakeholder fora. In case of any amendments to the protocol following its
204 publication, we will provide the date of each amendment, describe the change(s), and report
205 the rationale for the change(s) in future publications arising from this protocol.

206 **Strengths and limitations**

207 The strength of this review lies within the systematic approach to synthesizing the diverse
208 evidence on the integration of large language models (LLMs) in healthcare, clinical and
209 public health, with a focus on prompt engineering, auditing, and user feedback mechanisms,
210 which, is a relatively niche concept. By utilizing a broad range of sources, including peer-
211 reviewed studies and grey literature, the review will provide a comprehensive understanding
212 of current practices, trends, and gaps in the field. Its focus on healthcare-specific applications
213 ensures relevance to real-world policy relevant challenges, while the inclusion of feedback
214 mechanisms highlights its alignment with user-centered design principles. Conversely,
215 potential limitations include the restriction to English and French language publications,
216 which may exclude relevant studies in other languages, and reliance on available literature
217 that may underrepresent unpublished or proprietary methods used by private collectives.
218 Additionally, the rapidly evolving nature of LLM technologies means that findings may
219 quickly become outdated, necessitating continuous updates to maintain relevance.

220 **Acknowledgements**

221 The authors would like to acknowledge Griffith University, and the African Population and

222 Health Research Center(APHRC) for their institutional support towards this innovative

223 collaboration.

224

225 **Funding**

226 No funding source to disclose

227 **Declarations**

228 None to declare

229

230 **Author contributions**

231 HN and EO initiated the conception of the review and EO led the drafting of the protocol
232 while HN provided critical review.

233 **Conflicts of interest**

234 None to disclose

235

237 **References**

- 238 [1] Marr B. Understanding the 4 Types of Artificial intelligence. Bernard Marr 2021.
239 <https://bernardmarr.com/understanding-the-4-types-of-artificial-intelligence/> (accessed
240 January 14, 2025).
- 241 [2] Qiu J, Lam K, Li G, Acharya A, Wong TY, Darzi A, et al. LLM-based agentic systems in
242 medicine and healthcare. *Nat Mach Intell* 2024;6:1418–20.
243 <https://doi.org/10.1038/s42256-024-00944-1>.
- 244 [3] Thirunavukarasu AJ, Ting DSJ, Elangovan K, Gutierrez L, Tan TF, Ting DSW. Large
245 language models in medicine. *Nat Med* 2023;29:1930–40.
246 <https://doi.org/10.1038/s41591-023-02448-8>.
- 247 [4] Ng KKY, Matsuba I, Zhang PC. RAG in Health Care: A Novel Framework for
248 Improving Communication and Decision-Making by Addressing LLM Limitations.
249 *NEJM AI* 2025;2:AIra2400380. <https://doi.org/10.1056/AIra2400380>.
- 250 [5] Mirzaei T, Amini L, Esmailzadeh P. Clinician voices on ethics of LLM integration in
251 healthcare: a thematic analysis of ethical concerns and implications. *BMC Med Inform*
252 *Decis Mak* 2024;24:250. <https://doi.org/10.1186/s12911-024-02656-3>.
- 253 [6] Abd-alrazaq A, AlSaad R, Alhuwail D, Ahmed A, Healy PM, Latifi S, et al. Large
254 Language Models in Medical Education: Opportunities, Challenges, and Future
255 Directions. *JMIR Med Educ* 2023;9:e48291. <https://doi.org/10.2196/48291>.
- 256 [7] Safranek CW, Sidamon-Eristoff AE, Gilson A, Chartash D. The Role of Large Language
257 Models in Medical Education: Applications and Implications. *JMIR Med Educ*
258 2023;9:e50945. <https://doi.org/10.2196/50945>.
- 259 [8] Jiang X, Yan L, Vavekanand R, Hu M. Large Language Models in Healthcare Current
260 Development and Future Directions. *Gener. AI Res., Remote, Hong Kong SAR China:*
261 2023. <https://doi.org/10.5281/zenodo.12705655>.

- 262 [9] Yang R, Tan TF, Lu W, Thirunavukarasu AJ, Ting DSW, Liu N. Large language models
263 in health care: Development, applications, and challenges. *Health Care Sci* 2023;2:255–
264 63. <https://doi.org/10.1002/hcs2.61>.
- 265 [10] Joshi I, Shahid S, Venneti S, Vasu M, Zheng Y, Li Y, et al. CoPrompter: User-Centric
266 Evaluation of LLM Instruction Alignment for Improved Prompt Engineering 2024.
267 <https://doi.org/10.48550/arXiv.2411.06099>.
- 268 [11] Marvin G, Hellen N, Jjingo D, Nakatumba-Nabende J. Prompt Engineering in Large
269 Language Models. In: Jacob IJ, Piramuthu S, Falkowski-Gilski P, editors. *Data Intell.*
270 *Cogn. Inform.*, Singapore: Springer Nature; 2024, p. 387–402.
271 https://doi.org/10.1007/978-981-99-7962-2_30.
- 272 [12] Ingeniería de instrucciones: una metodología para optimizar interacciones con
273 Modelos de Lenguaje de IA en el campo de ingeniería | DYNA n.d.
274 <https://revistas.unal.edu.co/index.php/dyna/article/view/111700> (accessed December 27,
275 2024).
- 276 [13] Arvidsson S, Axell J. Prompt engineering guidelines for LLMs in Requirements
277 Engineering 2023.
- 278 [14] Meskó B. Prompt Engineering as an Important Emerging Skill for Medical
279 Professionals: Tutorial. *J Med Internet Res* 2023;25:e50638.
280 <https://doi.org/10.2196/50638>.
- 281 [15] Torkington. These are the 3 biggest emerging risks the world is facing. *World Econ*
282 *Forum* 2024. <https://www.weforum.org/stories/2024/01/ai-disinformation-global-risks/>
283 (accessed January 21, 2025).
- 284 [16] Lawton G. Generative AI Ethics: 8 Biggest Concerns and Risks. *Search Enterp AI*
285 2024. [https://www.techtarget.com/searchenterpriseai/tip/Generative-AI-ethics-8-biggest-](https://www.techtarget.com/searchenterpriseai/tip/Generative-AI-ethics-8-biggest-concerns)
286 [concerns](https://www.techtarget.com/searchenterpriseai/tip/Generative-AI-ethics-8-biggest-concerns) (accessed January 21, 2025).

287 [17] Peters MD, Godfrey C, McInerney P, Munn Z, Tricco AC, Khalil H. Scoping reviews.
 288 In: Aromataris E, Lockwood C, Porritt K, Pilla B, Jordan Z, editors. JBI Man. Evid.
 289 Synth., JBI; 2024. <https://doi.org/10.46658/JBIMES-24-09>.

290 Appendices

291 Appendix 1: Search strategy

292 Database Search Strategy

293 We will use the advance search builder of the PubMed and CINAHL databases, which allows
 294 for the use of wildcards (*) and lengthy search terms. We limited our search to a period
 295 between 1st January 2015 to 31st January 2025 .

296 *Table 1 PubMed Search Strategy*

#1	"Large Language Models" OR LLMs OR GPT-3 OR GPT-4 OR BERT OR RoBERTa OR T5 OR MedPaLM OR BioGPT OR ClinicalBERT OR pubmedBERT OR BioMegatron OR GatorTron OR ALBERT OR CODEX OR Glalactica OR SapBERT OR MIMIC-BERT
#2	"Prompt Engineering" OR "Natural language processing" OR "machine learning" OR "deep learning" OR "language models" OR prompts OR fine-tuning
#3	Medicine OR disease OR epidemiology OR "public health" OR "healthcare systems" OR "healthcare delivery"OR"auditing methods"OR "user feedback",
#4	Diagnosis OR treatment OR "drug discovery" OR "medical imaging" OR health OR Application OR education OR "disease surveillance"
#5	#1 AND #4
#6	#1 AND #2 AND #4
#7	#1 AND #2 AND #3 AND #4

297

298

299

300

301 **Appendix 2: Data extraction instrument**

Key domain	Data point	Description
Study Identification	Title	Indicates the study title
	Authors	Indicates the names of the authors
	Publication year	
	Country /region	Indicates the country of study and geographical region
	Source	This indicates Database, journal, or organization where the study was found
Study characteristics	Study type	Type of study (e.g., experimental, observational, qualitative, case study)
	Objective	Primary aim or research question of the study.
	Context	Description of the healthcare or public health setting
LLM detail	LLM type	Name and version of the large language model used (e.g., GPT, BERT)
	Application	Specific application of the LLM in healthcare or public health (e.g., clinical decision-making).
	Task specific use	Specific tasks for which the LLM was used or tailored.
	Prompt Engineering details	Techniques
	Tailoring/Customizing strategies	Techniques for tailoring prompts to improve accuracy and utility.
	Auditing process	Methods used to assess LLM outputs for accuracy, fairness, and effectiveness.
	Auditing outcomes	Key findings or metrics from the auditing process.
End-user feedback	Feedback processes and mechanism	Processes for gathering feedback from healthcare professionals.
	Feedback implementation	Description of how feedback was used to refine prompts or LLM outputs.

	Feedback impact assessment	Outcomes resulting from feedback integration.
Results and outcomes	Key findings	Major findings relevant to the study objectives
	Effectiveness	Measures of success or effectiveness of LLM application
	Challenges	Any reported challenges or limitations in the study.
Gaps and relevance	Relevance to review objectives	Description of how the study aligns with the review objectives.
	Identified gaps	Gaps or unanswered questions highlighted by the study.
Quality assessment	Randomized Trials/Interventions	Cochrane Risk of Bias assessment tool
	Observational studies	Newcastle-Ottawa scale (NOS) for non-randomized studies
	Qualitative studies	JBI Critical Appraisal Checklist
	Study limitations	Limitations or biases reported by the study authors
	Reviewer observations	Any additional domains or observations identified by the reviewer.

302

303