

1 **The accuracy of large language models in labelling neurosurgical ‘case-control studies’ and risk of bias**
2 **assessment: protocol for a study of interrater agreement with human reviewers.**

3 **Authors:** Joanne Igoli^{1,2,#}, Temidayo Osunronbi^{1,3,#}, Olatomiwa Olukoya^{1,4}, Jeremiah Oluwatomi Itodo Daniel¹,
4 Hillary Alemenzohu¹, Alieu Kanu¹, Alex Mwangi Kihunyu¹, Ebuka Okeleke¹, Henry Oyoyo¹, Oluwatobi
5 Shekoni¹, Damilola Jesuyajolu¹, Andrew F Alalade⁵

6 1. Neurosurgery section, Surgery Interest Group of Africa, Lagos, Nigeria.

7 2. Deanery of Clinical Sciences, The University of Edinburgh, Edinburgh, United Kingdom.

8 3. Department of Neurosurgery, Salford Royal NHS Foundation Trust, Manchester, United Kingdom.

9 4. The National Hospital for Neurology and Neurosurgery, London, United Kingdom.

10 5. Department of Neurosurgery, Royal Preston Hospital, Lancashire Teaching Hospitals NHS Foundation Trust,
11 Preston, United Kingdom.

12 # Joint-first authors: contributed equally.

13 **Please address all correspondence to:**

14 Dr Joanne Igoli

15 Deanery of Clinical Sciences, The University of Edinburgh, Edinburgh, United Kingdom. EH16 4SB

16 Email: s1408071@sms.ed.ac.uk

17 **Short title:** The accuracy of large language models in assessing neurosurgical ‘case-control studies’.

18 **Abstract**

19 **Introduction:** Accurate identification of study designs and risk of bias (RoB) assessment is crucial for evidence
20 synthesis in research. However, mislabelling of case-control studies (CCS) is prevalent, leading to a downgraded
21 quality of evidence. Large Language Models (LLMs), a form of artificial intelligence, have shown impressive
22 performance in various medical tasks. Still, their utility and application in categorising study designs and assessing
23 RoB needs to be further explored. This study will evaluate the performance of four publicly available LLMs
24 (ChatGPT-3.5, ChatGPT-4, Claude 3 Sonnet, Claude 3 Opus) in accurately identifying CCS designs from the
25 neurosurgical literature. Secondly, we will assess the human-LLM interrater agreement for RoB assessment of
26 **NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

27 Methods: We identified thirty-four top-ranking neurosurgical-focused journals and searched them on
28 PubMed/MEDLINE for manuscripts reported as CCS in the title/abstract. Human reviewers will independently
29 assess study designs and RoB using the Newcastle-Ottawa Scale. The methods sections/full-text articles will be
30 provided to LLMs to determine study designs and assess RoB. Cohen's kappa will be used to evaluate human-
31 human, human-LLM and LLM-LLM interrater agreement. Logistic regression will be used to assess study
32 characteristics affecting performance. A p -value < 0.05 at a 95% confidence interval will be considered
33 statistically significant.

34 Conclusion If the human-LLM agreement is high, LLMs could become valuable teaching and quality assurance
35 tools for critical appraisal in neurosurgery and other medical fields. This study will contribute to validating LLMs
36 for specialised scientific tasks in evidence synthesis. This could lead to reduced review costs, faster completion,
37 standardisation, and minimal errors in evidence synthesis.

38

39 **Keywords**: Case control study; Neurosurgery; Artificial intelligence; Large Language Model; ChatGPT

40

41 **Introduction**

42 Observational studies, including cross-sectional, cohort, and case-control studies, are ideal for neurosurgery
43 research when placebo or no-treatment groups are risky or ethically challenging or when randomised controlled
44 trials are impractical due to logistical complexities or inadequacy for addressing clinical questions [1].

45 Cross-sectional studies concurrently evaluate exposure and outcome status at a single time point without
46 longitudinal follow-up [2]. Cohort studies divide participants based on exposures or treatments and follow them
47 over a period, either prospectively or retrospectively, to compare outcomes between the groups [2]. Case-control
48 studies (CCS) compare individuals with (case) and without (control) a particular outcome, retrospectively
49 examining differences in exposure risk factors [2].

50 Unlike other observational studies, CCS is best suited for investigating rare outcomes or those with long latency
51 periods, leading to its increasing use in neurosurgery [3]. However, they have limitations such as recall bias and
52 the inability to determine incidence and absolute risk or establish temporality [1–3]. Previous research indicates
53 a significant prevalence of misclassified 'CCS' in neurosurgery literature, ranging from 41% to 63% [1–3].
54 Mislabelling of CCS is not unique to neurosurgery, with mislabelling rates reaching as high as 30% to 97% in

55 other fields [4–6]. Cohort studies are most frequently mislabelled as CCS, leading to a downgrading of evidence
56 quality since cohort studies represent the highest level of evidence among observational studies [1–3].

57 Moreover, mislabelled CCS often report odds ratios instead of relative risks, leading to distorted effect size
58 measurements, particularly in systematic reviews and meta-analyses [3]. Hence, accurate labelling of study
59 designs is crucial for stakeholders, including readers, authors, and editors. In addition, assessing the risk of bias
60 (RoB) is critical to systematic reviews. This process involves reviewing and understanding each eligible study,
61 which relies on a solid grasp of study methods and RoB assessment tools. However, RoB assessment is labour-
62 intensive and prone to human error, which may introduce biases in the conclusions of an evidence synthesis [7].

63 The recent upsurge in excitement about artificial intelligence (AI) has increased its impact on every aspect of
64 healthcare [8]. Large Language Models (LLMs), a subset of AI, are trained on extensive amounts of text data to
65 understand, generate, and process human-like language for various natural language processing tasks [8, 9]. Many
66 healthcare professionals have begun to use LLMs such as ChatGPT and Claude as advanced search tools for
67 complex medical information. These models exhibit emergent properties resembling human-level intelligence and
68 have demonstrated impressive performance on various medical speciality exams, including neurosurgery [10, 11],
69 and have even succeeded in challenging tests like the United States Medical Licensing Examination [9].
70 Additionally, some machine learning systems, such as the RobotReviewer, have shown high accuracy in
71 evaluating the risk of bias in clinical trials [12]. However, the potential of LLMs, an advanced AI tool, in
72 categorising study designs and assessing RoB in neurosurgery research still needs to be explored. Leveraging
73 LLMs in these tasks may lead to reduced review costs, faster completion times, and decreased errors in the
74 assessment process.

75 This study aims to evaluate the performance of four publicly available LLMs (ChatGPT-3.5 [OpenAI/Microsoft],
76 ChatGPT-4 [OpenAI/Microsoft], Claude 3 Sonnet [Anthropic], and Claude 3 Opus [Anthropic]) in accurately
77 identifying the design of 'case-control studies' in the neurosurgical literature. It also seeks to identify predictive
78 study characteristics that affect LLM performance. Additionally, we will evaluate human-LLM agreement in
79 overall and domain-level risk of bias (RoB) assessment using the Newcastle-Ottawa Scale for CCS [13].

80

81 **Materials and Methods**

82 Search strategy

83 There are no official lists of all neurosurgical journals in the literature, considering the ongoing introduction of
84 new journals. We conducted an online Google search using the phrase ‘top neurosurgery journals’ to compile a
85 speciality list of journals for neurosurgeons. This search yielded a Google Scholar [14] and Welch Medical Library
86 [15] list of neurosurgical journals, which we complemented with additional journals from a previously published
87 article [3]. We excluded the journals that have stopped publications and those with a nursing theme. Our search
88 strategy featured 34 PubMed-indexed journals (**Appendix 1**).

89 A PubMed/MEDLINE search was performed for all the articles in these thirty-four indexed journals from database
90 inception to 8 June 2024, using the search terms ‘case-control’, ‘case control’, ‘case controlled’, or ‘case-
91 controlled’ in the title or abstract.

92 The human reviewers

93 The assessment team will include a consultant/attending neurosurgeon (AFA), two neurosurgical
94 trainees/residents with training in critical appraisal/ postgraduate certificate in health research and statistics/
95 Masters of Neurosurgery by Research (TO, OO), and nine medical students/ clinicians who will be trained on
96 critical appraisal prior to commencing this study.

97 Eligibility criteria

98 Only original research articles reported as ‘case-control’ in the titles or abstracts will be included. Reviews,
99 commentaries, letters, genetic studies, animal studies, and cost-effectiveness studies will be excluded. Similarly,
100 articles will be excluded if they lack the term ‘case-control’/ ‘case control’/ ‘case controlled’/ ‘case-controlled’ in
101 their abstract/title or if this term was used in reference to another study. Studies with ambiguous study design
102 labels in their abstract/ title and/or those that use multiple study designs will be excluded (for example: ‘cross-
103 sectional case-control study’, ‘case-control cohort study’, ‘systematic review/ meta-analysis and case-control
104 study’). In addition, articles that are neurology-focused instead of neurosurgery-focused will be excluded.

105 The titles/abstracts will be screened independently by pairs of authors using the Rayyan software, with a third
106 author (TO) resolving any discrepancies.

107 Data extraction

108 Data extraction from the eligible full texts will be performed by a pair of authors (TO and other authors), with a
109 third author (OO) resolving any discrepancies. The following data will be extracted based on previous related
110 publications [1, 3]:

- 111 - Journal name (The journal results will be presented anonymously in the resulting publication).
- 112 - Year of publication (<2008, 2008 - 2019, >2019). The STROBE statement was published in 2007, and the
113 last publication on the mislabelling of case-control studies in neurosurgery was published in 2019 [2, 3]. This
114 forms the rationale for the year categories.
- 115 - Topic (spine, trauma, vascular, functional/epilepsy, neuro-oncology, paediatrics, skull base, pituitary,
116 hydrocephalus and other).
- 117 - Country of origin (based on where the study took place; the first author's country will be used when the study
118 location is not specified). Countries will be grouped by the number of case-control studies published (Group
119 A: countries with >10 case-control studies; Group B: countries with 5 to 10 case-control studies; Group C:
120 countries with <5 case-control studies). The countries will also be grouped by continents (Africa, Antarctica,
121 Asia, Australia, Europe, North America, and South America).
- 122 - Presence or acknowledgement of a case-control expert in the study (such as a statistician, epidemiologist, or
123 one with a master's degree or equivalent in public health)
- 124 - Study design characteristics:
- 125 ○ Aim of study (a): Outcome assessment
 - 126 ○ Aim of study (b): Risk factors assessment
 - 127 ○ Used logistic regression analysis.
 - 128 ○ Reported odds ratio (OR)
 - 129 ○ Used survival analysis/Kaplan-Meier curves.
- 130 - Terminology of the study:
- 131 ○ The word 'cohort' was used in the methods, results, or discussion sections.
 - 132 ○ The word 'outcome' was used in the results section.
 - 133 ○ The word 'prospective' or 'prospectively' was used in the methods section.
 - 134 ○ The word 'retrospective' or 'retrospectively' was used in the methods section.

135 Assessment of study design and risk of bias by human reviewers

136 The assessment of the study design of the eligible full text articles will be performed by a pair of authors (TO and
137 other authors), with a third author (OO) resolving any discrepancies. The human assessors will classify the studies
138 as 'true case-control studies' or 'non-case-control studies'. A study will be deemed a true case-control study if it
139 comprises three fundamental elements [1]: 1) compares a group of patients with a disease or who have experienced
140 an event with a control group lacking the disease or event; 2) a retrospective evaluation from the time point of a

141 known outcome is made; and 3) focuses on identifying risk factors/associations/causality of the disease or event.
142 The ‘non-case-control studies’ design will be specified as prospective cohort studies, retrospective cohort studies,
143 cross-sectional studies, case series, case reports, randomised clinical trials, and other.

144 The Newcastle-Ottawa Scale (NOS) will be used to evaluate the risk of bias (RoB) in the true case-control studies
145 [13]. The true case-control studies will be divided into five groups, and the RoB assessment will be performed by
146 a pair of authors, with a third author (TO or OO) adjudicating any discrepancies. Studies with NOS scores of 0-3,
147 4-5, 6-7, and 8-9 will be considered unsatisfactory, satisfactory, good, and very good quality, respectively.

148 Assessment of study design and risk of bias assessment by LLMs

149 For each eligible article obtained from the abstract/title screening, the methods section will be copied and imputed
150 separately into each LLM (ChatGPT-3.5 [OpenAI/Microsoft], ChatGPT-4 [OpenAI/Microsoft], Claude 3 Sonnet
151 [Anthropic], and Claude 3 Opus [Anthropic]) and the LLMs will be prompted with this question: ‘Some authors
152 may or may not correctly label their study design. Using the hierarchy of evidence, with a rationale, what is the
153 actual specific study design in the text below?’ To facilitate the assessment of the LLM-LLM intrarater agreement,
154 we will obtain LLM assessments in duplicate, i.e., two different authors (TO and OO) will separately use the
155 LLMs independently for the assessment of study design.

156 Subsequently, we will evaluate the LLMs RoB assessment for the author-labelled true CCS. The PDF files of the
157 eligible papers will be imputed separately as attachments into each LLM. The LLMs will be prompted with this
158 question: ‘Given that studies with an overall Newcastle-Ottawa scale (NOS) scores of 0-3, 4-5, 6-7, and 8-9 are
159 considered unsatisfactory, satisfactory, good, and very good quality, respectively, provide a domain-level and
160 overall risk of bias assessment for the following study using the Newcastle-Ottawa scale for case-control studies.’

161 If we are unable to attach the PDF file or the LLM is unable to read the PDF file, we will copy the methods text
162 and the patients/participants characteristics/demographics subsection of the results and impute this into the LLM.
163 To facilitate the assessment of the LLM-LLM intrarater agreement, we will obtain LLM RoB assessments in
164 duplicate — i.e., two authors (TO and OO) will each use the LLMs independently for the RoB assessment.

165 Statistical analysis and reporting:

166 Statistical analyses will be conducted on IBM SPSS Statistics 27 (Windows).

167 LLM-LLM and human-human interrater reliability for the study design and RoB assessments will be assessed
168 using Cohen’s kappa (κ) for categorical data. In the event of LLM-LLM (for example, ChatGPT-3.5 - ChatGPT-

169 3.5) discrepancies, we will reduce the duplicate assessments to a single assessment for each study by randomly
170 choosing one of the assessments for each study.

171 We will calculate the proportion of articles labelled as ‘case-control’ in the title/abstract that are true case-control
172 studies as identified by human reviewers. Furthermore, using the study design determined by human reviewers in
173 this study as a reference, we will calculate the proportion of study design correctly labelled by each LLM.
174 Subsequently, LLM-human inter-rater reliability for the study design and RoB assessments will be assessed using
175 Cohen’s kappa (κ) for categorical data. Kappa values will be interpreted as follows: values ≤ 0 (no agreement),
176 0.01–0.20 (slight agreement), 0.21–0.40 (fair agreement), 0.41– 0.60 (moderate agreement), 0.61–0.80
177 (substantial agreement), and 0.81–1.00 (almost perfect agreement) [16].

178 Simple logistic regression analyses will be conducted to assess the associations between select study
179 characteristics and whether a study was a true case-control (yes/no). These analyses will also be conducted for
180 each LLM to assess the association between the select study characteristics and the accurate labelling of study
181 designs by the LLM (yes/no). A p -value < 0.05 at a 95% confidence interval will be considered statistically
182 significant.

183 **Discussion**

184 To our knowledge, this study will be the first to evaluate interrater agreement between human reviewers and
185 LLMs in labelling study designs and assessing RoB in neurosurgical case-control studies.

186 If the human-LLM interrater agreement is almost perfect, then LLMs could become valuable tools for teaching
187 and quality assurance in critical appraisal and identifying study designs in neurosurgery and other fields. This
188 study is expected to make a significant early contribution to the research exploring the utilisation and validation
189 of general-purpose LLMs trained on vast internet data for specialised scientific tasks. It is anticipated that this
190 study will mark the beginning of a series focused on employing LLMs in evidence synthesis. The investigation
191 into the application of LLMs, particularly for systematic reviews, is poised to bring about significant changes in
192 how evidence synthesis tasks are conducted, who undertakes them, the speed and cost of completion, and the way
193 primary studies are conducted and reported to enhance comprehensibility for artificial intelligence.

194 Limitations

195 This study will not include some non-neurosurgical-specific journals where the neurosurgical community may
196 choose to publish. Thus, the representativeness of the selected articles as a sample of all neurosurgical case-control

197 studies can be questioned. Based on our exclusion criteria, articles lacking explicit mention of “case control” in
198 title or abstract will be excluded. However, the improper use of the term "case-control" might be more prevalent
199 in these articles and missed in our search. Though unlikely, reverse mislabelling could occur, where true case-
200 control studies may not have been labelled as such and thus missed in our search.

201 To evaluate LLMs’ ability in RoB assessment, we will provide only the methods and results section or full articles
202 (where possible) to the LLMs. Human reviewers will have access to the entire text and supplementary materials
203 where available, providing them with more information about each study than LLMs. As a result, the human-
204 LLM interrater agreements we estimate are expected to be conservative estimates of what is achievable.

205 **Declarations**

206 **Authors' contributions:** Joanne Igoli, Temidayo Osunronbi, Olatomiwa Olukoya, Damilola Jesuyajolu, and
207 Andrew F Alalade contributed to the study's conception and design. The first draft of the paper was written by
208 Temidayo Osunronbi, and all authors commented on the subsequent versions of the manuscript. All authors read
209 and approved the final manuscript.

210 **References**

- 211 Nesvick CL, Thompson CJ, Boop FA, Klimo P Jr. Case-control studies in neurosurgery. *J Neurosurg.*
212 2014;121(2):285-296. doi:10.3171/2014.5. JNS132329
- 213 Kicielinski KP, Dupépe EB, Gordon AS, Mayo NE, Walters BC. What Isn't a Case-Control Study?.
214 *Neurosurgery.* 2019;84(5):993-999. doi:10.1093/neuros/nyy591
- 215 Esene IN, Mbuagbaw L, Dechambenoit G, Reda W, Kalangu KK. Misclassification of Case-Control Studies in
216 Neurosurgery and Proposed Solutions. *World Neurosurg.* 2018;112:233-242. doi:10.1016/j.wneu.2018.01.171
- 217 Grimes DA. "Case-control" confusion: mislabeled reports in obstetrics and gynecology journals. *Obstet*
218 *Gynecol.* 2009;114(6):1284-1286. doi:10.1097/AOG.0b013e3181c03421
- 219 Mayo NE, Goldberg MS. When is a case-control study a case-control study?. *J Rehabil Med.* 2009;41(4):217-
220 222. doi:10.2340/16501977-0341
- 221 Mihailovic A, Bell CM, Urbach DR. Users' guide to the surgical literature. Case-control studies in surgical
222 journals. *Can J Surg.* 2005;48(2):148-151.
- 223 Könsgen N, Barcot O, Heß S, et al. Inter-review agreement of risk-of-bias judgments varied in Cochrane
224 reviews. *J Clin Epidemiol.* 2020;120:25-32. doi:10.1016/j.jclinepi.2019.12.016
- 225 Sallam M. ChatGPT Utility in Healthcare Education, Research, and Practice: Systematic Review on the
226 Promising Perspectives and Valid Concerns. *Healthcare (Basel).* 2023;11(6):887.
227 doi:10.3390/healthcare11060887
- 228 Kung TH, Cheatham M, Medenilla A, et al. Performance of ChatGPT on USMLE: Potential for AI-assisted
229 medical education using large language models. *PLOS Digit Health.* 2023;2(2):e0000198.
230 doi:10.1371/journal.pdig.0000198

- 231 Hoch CC, Wollenberg B, Lüers JC, et al. ChatGPT's quiz skills in different otolaryngology subspecialties: an
232 analysis of 2576 single-choice and multiple-choice board certification preparation questions. *Eur Arch*
233 *Otorhinolaryngol.* 2023;280(9):4271-4278. doi:10.1007/s00405-023-08051-4
- 234 Ali R, Tang OY, Connolly ID, et al. Performance of ChatGPT, GPT-4, and Google Bard on a Neurosurgery
235 Oral Boards Preparation Question Bank. *Neurosurgery.* 2023;93(5):1090-1098.
236 doi:10.1227/neu.0000000000002551
- 237 Marshall IJ, Kuiper J, Wallace BC. RobotReviewer: evaluation of a system for automatically assessing bias in
238 clinical trials. *J Am Med Inform Assoc.* 2016;23(1):193-201. doi:10.1093/jamia/ocv044
- 239 Wells G, Shea B, O'Connell D, Peterson J. The Newcastle-Ottawa Scale (NOS) for assessing the quality of
240 nonrandomised studies in meta-analyses. Ottawa, ON: Ottawa Hospital Research Institute.
241 https://www.ohri.ca/programs/clinical_epidemiology/oxford.asp [Accessed on 30 April 2023]
- 242 Google Scholar. Top journals in neurosurgery.
243 https://scholar.google.co.uk/citations?view_op=top_venues&hl=en&vq=med_neurosurgery [Accessed on 30
244 April 2024]
- 245 Welch Medical Library. Journals by Subject: Neurosurgery.
246 <https://welch.jhmi.edu/journalsbysubject?s=Neurosurgery> [Accessed on 30 April 2024]
- 247 McHugh ML. Interrater reliability: the kappa statistic. *Biochem Med (Zagreb).* 2012;22(3):276-282.