

1 Protocol for the development of an artificial intelligence 2 extension to the Consolidated Health Economic Evaluation 3 Reporting Standards (CHEERS) 2022

4 1. Manuscript details

5 1.1. Authors

6 Claire Hawksworth¹, Jamie Elvidge¹, Saskia Knies², Antal Zemlenyi³, Zsuzsanna Petykó³, Pekka
7 Siirtola⁴, Gunjan Chandra⁴, Divya Srivastava⁵, Alastair Denniston⁶, Anastasia Chalkidou¹, Julien
8 Delaye⁷, Petros Nousios⁸, Manuel Gomes⁹, Tuba Saygin Avsar¹, Junfeng Wang¹⁰, Stavros Petrou¹¹,
9 Dalia Dawoud^{1,12}

10 1. National Institute for Health and Care Excellence (NICE), UK; 2. Zorginstituut Nederland (National
11 Health Care Institute), The Netherlands; 3. Syreon Research Institute, Hungary; 4. University of
12 Oulu, Finland; 5. The London School of Economics and Political Science, UK; 6. University of
13 Birmingham, UK; 7. EURORDIS – Rare Diseases Europe, Belgium; 8. Tandvårds-och
14 läkemedelsförhållningsverket (The Dental and Pharmaceutical Benefits Agency), Sweden; 9. University
15 College London, UK; 10. Utrecht University, Netherlands; 11. University of Oxford, UK; 12. Cairo
16 University, Egypt.

17 1.2. Abstract

18 **Introduction:** AI interventions for health care are on the rise. Decisions about coverage and
19 reimbursement are often informed by Health Technology Assessment (HTA) bodies, who rely on
20 Health Economic Evaluations (HEEs) to estimate the value for money (cost effectiveness) of
21 interventions. Transparent reporting of HEEs ensures they can be used for decision making.
22 Reporting guidance exists to support this, such as the Consolidated Health Economic Reporting
23 Standards (CHEERS) checklist. We aim to identify consensus about specific items should be
24 reported by HEEs that evaluate AI interventions and, if such items are identified, to develop them
25 into an extension to CHEERS: “CHEERS-AI”.

26 **Methods and analysis:** The project will have 4 phases:

- 27 • Phase 1 is a literature review to help identify potential AI-related reporting items.
- 28 • Phase 2 commences a Delphi process, with a series of surveys to elicit the importance of the
29 potential AI-related reporting items.
- 30 • Phase 3 is a consensus-generation meeting to agree on the final extension items.

31 • Phase 4 is dissemination of the project's outputs.

32 **Ethics and dissemination:** This study has received ethical approval from Newcastle University
33 Ethics Committee (reference: 28568/2022). The findings will be available in as an open access
34 article and disseminated through blogs, newsletters, and presentations.

35 **1.3. Funding statement**

36 This study is supported by the Next Generation Health Technology Assessment (HTx) project. The
37 HTx project has received funding from the European Union's Horizon 2020 research and innovation
38 programme under grant agreement N° 825162. This dissemination reflects only the views of the
39 authors and the Commission is not responsible for any use that may be made of the information it
40 contains.

41 2. Introduction

42 In recent times there has been a rapid increase in the development of technologies with an artificial
43 intelligence (AI) component for health care interventions. This is evidenced in the number of
44 approvals given by regulatory bodies. Between 1997 and 2021, the Food and Drug Administration in
45 the United States approved 350 AI technologies with 91% of them approved since 2015 (1). In
46 2021, the European Medicines Agency (EMA) led a report on behalf of the International Coalition of
47 Medicines Regulatory Authorities documenting a horizon scanning exercise in AI and highlighting
48 regulatory challenges (2). This was in response to these new technologies increasingly challenging
49 regulatory frameworks and a need for recommendations on how to adapt them.

50 AI is a broad term to encompass iterative, 'learning' algorithms that use data and high computing
51 power to make interpretations, predictions or decisions (2). Some AI technologies are fixed, and
52 others are adaptive. Various subsets of AI, such as machine learning (ML), are being used
53 throughout the drug discovery process for target validation, identification of biomarkers, and
54 analysis of clinical trial data (3). As well as assisting with the drug development process, AI is also
55 featuring in the end product, and it is these health technologies that are the focus of this paper.
56 Examples of AI health interventions include systems for screening and triage, diagnosis, prognosis,
57 decision support, and treatment recommendation (4,5).

58 To ensure their appropriate use in healthcare pathways, we need to understand what benefits new
59 AI technologies bring, and at what cost. There are established methods to do this for
60 pharmacological and diagnostic interventions, but AI algorithms may be distinct from more
61 traditional interventions in numerous challenging ways. Firstly, they have the potential to learn over
62 time, meaning the relationship between intervention and outcome may not be fixed. This has
63 implications when considering future benefits, such as choosing a suitable method or assumption
64 for long-term treatment outcomes. We often see an assumption that the treatment effect of a
65 medicine wanes over time, but how might healthcare decision makers appropriately value on an AI
66 intervention that might get *more* effective over time? Secondly, the user is most often a health care
67 professional rather than a patient, and the degree to which the clinician employs the results of the AI
68 intervention may vary, particularly when its purpose is a decision-support tool. Thirdly, trial data
69 normally underpin a health technology assessment (HTA) and reimbursement decisions. However,
70 to date, AI technologies have not typically been subjected to interventional trials, meaning various
71 data sources or assumptions will be required to inform a value assessment. Although randomised
72 controlled trials (RCTs) are increasingly being conducted to evaluate the clinical efficacy of
73 interventions with an AI component, there are concerns relating to their design and reporting. To try
74 to address these concerns, AI extensions to reporting checklists have been developed; for example,
75 for protocols (SPIRIT-AI) (5) and trials (CONSORT-AI) (4).

76 In addition, the ways in which AI-based intervention are developed arguably create an extra,
77 inherent layer of uncertainty. Their function and attainment depend on the data sets used to train
78 and validate their underlying algorithms. This development, or learning, step precedes any study of
79 efficacy relative to the standard of care, which tends to be the primary source of potential
80 uncertainty for more traditional interventions.

81 Health economic evaluations (HEEs) assessing the cost effectiveness of health interventions are
82 often used by HTA bodies to make their reimbursement recommendations. HTA bodies will
83 increasingly be expected to assess the value of health technologies that use AI. For example, the
84 National Institute of Health and Care Excellence (NICE) in the UK recently updated their Evidence
85 Standards Framework to reflect and include adaptive AI and data-driven technologies (6,7). The
86 usefulness of a published HEE to decision makers depends on how well it is conducted and
87 reported. Reporting guidelines can improve their transparency and completeness. A prominent HEE
88 reporting checklist is the Consolidated Health Economic Evaluation Reporting Standards (CHEERS)
89 statement. It was originally published in 2013 to help authors accurately report details of the HEE,
90 including the health intervention, what was being compared and in what context, how the evaluation
91 was undertaken, and what the findings were (8). This checklist outlined minimum reporting
92 standards and the increased transparency allows decision-makers such as HTA bodies and payers
93 to judge the quality and appropriateness of the HEE for their decision problem, facilitating trust in
94 the results. The CHEERS statement was updated in 2022 (9) and now comprises a 28-item
95 checklist including methodological approach, data identification, model inputs, assumptions,
96 uncertainty analysis, and conflicts of interest.

97 CHEERS 2022 does not include any reporting items that are specific to potential AI components of
98 an intervention, but the authors of CHEERS 2022 explicitly “encourage those who see opportunities
99 to expand CHEERS 2022 items or create additional reporting guidance that provides clarification in
100 specific areas to work with members of the CHEERS Task Force to develop CHEERS extensions in
101 these areas”. As noted above, extensions for AI health interventions have already been developed
102 for other checklists, demonstrating a system wide need and motivation for improving best practice
103 around data collection and transparency. Including AI-specific items in the reporting of HEEs is a
104 logical step to contribute to this standard setting for AI interventions. It will help to ensure that all
105 relevant information required for decision-making is available to decision-makers.

106 **3. Methods**

107 Our research approach was guided by the EQUATOR (Enhancing the QUALity and Transparency Of
108 health Research) Network's recommended steps for developing a health research reporting
109 guideline (10) and methods used to develop other related extensions (CHEERS 2022, CONSORT-
110 AI and SPIRIT-AI). The guideline extension is registered on the EQUATOR Network website (11).
111 The structure and writing of this protocol were guided by the recently published protocol for the
112 SPIRIT-SURROGATE and CONSORT-SURROGATE extensions (12).

113 A project management group led by NICE is organising and conducting the project with oversight
114 from a Steering Group. The Steering Group is a multi-disciplinary and international group with
115 representation from University of Oulu, Finland; Zorginstituut Nederland (National Health Care
116 Institute) and Utrecht University, the Netherlands; Syreon Research Institute, Hungary; Tandvårds-
117 och läkemedelsförmånsverket (The Dental and Pharmaceutical Benefits Agency), Sweden; and
118 The London School of Economics and Political Science, the University of Birmingham, University
119 College London and University of Oxford, UK. The Steering Group includes a representative from
120 the CHEERS Task Force to provide expert input. The Steering Group was formed in December
121 2022.

122 This study has been supported by Next Generation Health Technology Assessment ([HTx](#)), which is
123 a Horizon 2020 project supported by the European Union, lasting for 5 years from January 2019. Its
124 main aim is to create a framework for the next generation of HTA to support patient-centred,
125 societally oriented, real-time decision making on access to and reimbursement for health
126 technologies throughout Europe.

127 **5.1. Phase 1: Systematic literature review**

128 We conducted a systematic literature review in the summer of 2022 to assess the methodological
129 and reporting quality of HEEs of AI-based technologies. This updated a previously published review
130 by Voets et al (13). Our search was performed on 17th June 2022 and found 21 HEE studies
131 published in the preceding 15 months. This review was used to identify potential AI-extension items.
132 Members of the Steering Group were also able to contribute potential AI-extension items, based on
133 their knowledge and experience of AI, HEE, and reporting guideline development. This led to a
134 'long-list' of potential items for an AI extension to CHEERS 2022. The review and Steering Group
135 also helped to identify subject matter experts who could participate in the Delphi study. This group
136 are referred to as the Expert Panel (EP).

137 **5.2. Phase 2: Consensus-generation surveys (Delphi process)**

138 This phase will involve participants rating long-list candidate items generated in phase 1 and
139 suggesting additional items not included in the long-list. There will also be the opportunity to revise
140 wording for the items. Proposed timelines are to open the first survey round in May 2023 to coincide
141 with the ISPOR 2023 conference in Boston, US. The consensus process will dictate the number of
142 necessary survey rounds, but it is anticipated that the whole project will complete in during 2023.

143 **5.2.1. Survey design and setting**

144 This methodology follows that used for the development of CHEERS 2022 (9). CHEERS 2022
145 employed a modified Delphi process. Delphi is a widely recognised and used method for
146 consensus-building and revolves around the following key steps: identification of factors,
147 anonymous surveys among subject matter experts to elicit importance, integration and controlled
148 feedback and presentation of aggregated data at consensus meetings (14). The survey will be
149 developed in [Snap Surveys](#) software.

150 We will conduct a minimum of 2 survey rounds and will consider additional rounds if necessary. This
151 approach has been taken for other guideline extensions (4,5,9,12).

152 **5.2.2. Sample size, recruitment, and inclusion criteria**

153 We will recruit an EP representing the following key stakeholder groups: health economists, AI
154 methodologists and academics, industry, policy makers, HTA experts, ethicists, patient
155 representatives, journal editors, healthcare professionals, payers, and research funders. This is
156 consistent with the EQUATOR Network's guidance (10) and groups who participated in related
157 extension. We anticipate inviting over 100 EP members.

158 Our approach to recruit EP members involves a multi-faceted approach. The Steering Group will
159 identify participants. This purposive sampling will utilise a snowball sampling method where invited
160 participants will be allowed to invite additional participants, meaning the total number of survey
161 recipients should far exceed the those identified by the Steering Group. The survey will elicit the
162 profession of the recipient to ensure that all respondents are part of one or more of the key
163 stakeholder groups. We will also approach relevant professional groups such as the ISPOR
164 Machine Learning Task Force. We will utilise authors identified in the phase 1 as another source of
165 potential participants and will coordinate completions of the survey with the ISPOR conference,
166 taking place in May 2023 in Boston, US. All EP members will be sent an introductory email and
167 participant information sheet.

168 We will collect descriptive demographic data at the start of the survey, including stakeholder group,
169 country of work and years of relevant experience to indicate understanding of AI in healthcare and
170 HEE. Inclusion criteria are the key stakeholder groups previously specified. There are no exclusion
171 criteria, but this targeted recruitment should result in identification of suitable EP members.

172 There is guidance on the minimum number of survey responses to allow statistical rigour, with 30
173 commonly cited (15). By identifying and inviting over 100 experts to participate, we will allow
174 sufficient headroom for non-response and attrition between survey rounds.

175 **5.2.3. Data collection, analysis, and consensus definition**

176 The survey will be developed in consultation with the Steering Group, including a pilot prior to the
177 launch to ensure usability. All participants will be sent a link to the survey which will start with study
178 information and a tick box for consent. The first survey will be open for a 3-week window,
179 commencing May 2023. The second survey will be sent approximately 4 weeks after closure of the
180 first. Response rates will be monitored during the survey window and email reminders sent to
181 participants to increase response rates. Records will be kept of the number approached, and non-
182 responses.

183 The EP will be asked to vote on the relevance of candidate items when reporting a HEE of an AI-
184 based intervention. They will be asked to use a 9-point Likert rating scale, consistent with CHEERS
185 2022 and other reporting extensions. A 'don't know' option will also be available for each item, in
186 case a participant feels unable to provide a rating. Potential items will be grouped according to
187 standard sections of HEEs (e.g., title, abstract, methods, discussion), and each item will have an
188 accompanying definition and rationale for inclusion. We will employ the consensus definition used
189 for CHEERS 2022 (see Figure 1). After survey round 1, any items that were scored lower than 7 by
190 at least 70% of respondents will be excluded; that is, we will conclude that consensus has been
191 reached that those items are not relevant reporting standards for HEEs of AI interventions. Those
192 items will be 'rejected' and will not proceed to survey round 2 in their original form. Participants will
193 have the opportunity to comment on the wording and propose new items. If suitable revised wording
194 of original items has been proposed, then revised items may be included in survey round 2.

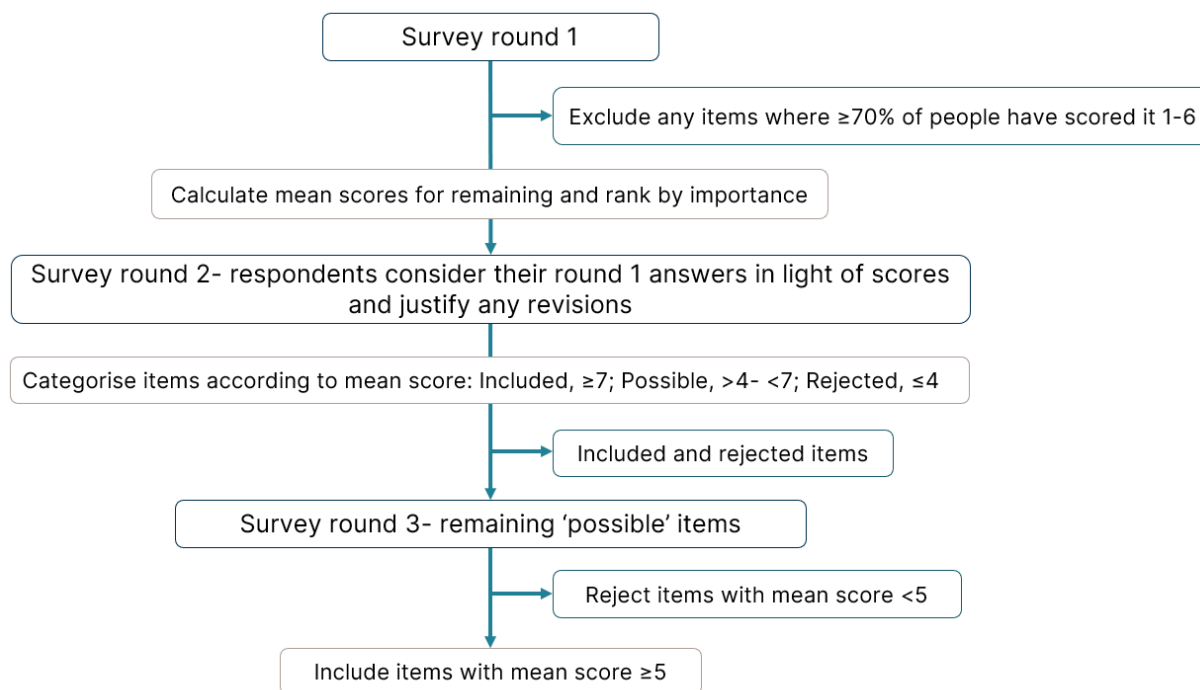


Figure 1. Flowchart showing the process for including and excluding items during the survey rounds.

195

196 Results from survey round 1 will be analysed in MS Excel to identify the mean scores and measure
197 of agreement (proportion of scores 7 or higher). After excluding those meeting the exclusion
198 threshold, the remaining items will be included in survey round 2. They will be presented in order of
199 mean score and the measure of agreement will also be shown. Respondents will have the
200 opportunity to consider their round 1 answers in light of the aggregate results and justify any
201 revisions. Any new items that were suggested by respondents in survey one will also be voted on in
202 the second survey round, along with proposed changes such a new wording or merging of items
203 from free-text responses.

204 After survey round 2, results will be analysed in MS Excel, and items with a mean score of 4 or less
205 will be categorised as 'rejected' (consensus reached). Items with mean scores of 7 or higher will be
206 grouped as 'included' (consensus reached). Items with mean scores above 4 but less than 7 will be
207 grouped as 'possible' (consensus not reached). Any such items will proceed to survey round 3,
208 presented in order of importance (mean score), alongside a measure of agreement (proportion of
209 scores 7 or higher). After this final survey round, items with a mean score of 5 or less will be
210 'rejected' (consensus reached). Items with a mean score above 5 will be 'included' (consensus
211 reached). Therefore, consensus will be achieved for all items after survey round 3, unless

212 participants provide substantial and conflicting free-text responses about the wording, or additional
213 items. Any such items will proceed to a consensus meeting for resolution.

214 **5.3. Phase 3: Consensus-generation meeting (Delphi process)**

215 A consensus meeting will be held virtually, with the aim of concluding the final extension list. We will
216 invite a purposive sample representative of the EP who completed every survey round. The
217 EQUATOR Network has guidance on conducting face to face consensus meetings (10).

218 **5.3.1. Structure and participants**

219 The COVID-19 pandemic has normalised virtual working and we propose to hold the meeting
220 virtually. This is also advantageous in terms of maximising engagement and attendance from a
221 range of geographical locations. The meeting length will be agreed after survey round 3, at which
222 point the number of items that still haven't achieved consensus (and therefore require extensive
223 discussion) will be known. However, the meeting time will be sufficient to allow discussion time for
224 all items. Other extensions have used meetings over two days (4,5,12).

225 At the end of survey round 2, participants will be invited to register their interest in attending the
226 consensus meeting. The Steering Group will purposively select members from this pool considering
227 the need to have an international multidisciplinary group of participants.

228 **5.3.2. Consensus procedure**

229 Attendees at the meeting will ratify all items about which consensus was reached during the first 2
230 survey rounds. Items that proceeded to survey round 3 will be discussed more comprehensively at
231 the meeting. Minor modifications to wording can be included by a simple 50% majority vote during
232 the meeting.

233 The richest discussion will be for any items where consensus has not yet been reached by the end
234 of the survey round 3. These will be any new items that were proposed in free-text responses during
235 survey round 3, and any items that received extensive requests for modification in free-text
236 responses (e.g., merging items), such that a simple 50% majority vote at the meeting would not be
237 appropriate to support inclusion. For these significant modifications, a 70% majority vote during the
238 meeting will be required to include the modified or new items.

239 Items that do not reach consensus (e.g., due to a large proportion of 'don't know' or abstained
240 votes) will be discussed further and voted on again, if appropriate, until consensus is reached or
241 time runs out. The Steering Group will make final decisions soon after the consensus meeting on
242 any outstanding items without consensus.

243 The consensus meeting will be recorded to ensure accurate recall, and minutes will be taken.

244 **5.4. Phase 4: Knowledge translation**

245 This phase includes all activities aiming to publish and publicise the extension. This objective will be
246 integrated and considered throughout all stages of the project.

247 **5.4.1. Pilot testing and revision of final checklist**

248 After the meeting we will conduct a pilot of the finalised extension with invited researchers to check
249 clarity of wording and identify any challenges. These will be people whom we invited but were
250 unable to join our Steering Group or who expressed an interest to join after survey round 1 had
251 started. This exercise will inform writing of the explanation and elaboration documents.

252 **5.4.2. Publications**

253 We aim to publish the CHEERS-AI extension in a high impact open access journal to maximise
254 dissemination. We will also utilise the ISPOR CHEERS Task Force to help publicise the extension.
255 We will seek the endorsement of the extension from journals and editorial groups.

256 **5.4.3. Partner and stakeholder engagement**

257 The project is registered on the EQUATOR website. We aim to have the final extension published
258 on the CHEERS website. The CHEERS statement is endorsed by ISPOR.

259 **5.4.4. Patient and public involvement**

260 We have patient advocacy on the Steering Group led by EURORDIS (Rare Diseases Europe). They
261 will ensure views and perspectives of patients and the public are represented at all stages. Patients
262 and the public are also one of our key stakeholder groups to be represented on the Expert Panel
263 and therefore will be involved in consensus building for the final checklist. EURORDIS and the NICE
264 lead for patient and public involvement will be requested to advise on dissemination activities.

265 **5.4.5. Ethics and dissemination**

266 An ethics application was submitted via the NICE ethical approval process. Ethics approval was
267 received from Newcastle University Ethics Committee who are the awarding body (reference:
268 28568/2022).

269 Expert Panel members will be provided with a participant information sheet and will be asked to
270 provide consent before completing the first survey. We will also obtain electronic written consent
271 before the consensus meeting for participation and recording of the meeting. Participants will have a
272 right to withdraw at any stage of the project. All data will be securely stored although we anticipate

273 that it will not be highly sensitive. We will ask participants if they prefer to opt out of
274 acknowledgement in any publications.

275 4. References

- 276 1. Miller M. FDA Publishes Approved List of AI/ML-enabled Medical Devices [Internet]. IQVIA
277 Blog. 2021 [cited 2023 May 9]. Available from: [https://www.iqvia.com/locations/united-](https://www.iqvia.com/locations/united-states/blogs/2021/10/fda-publishes-approved-list-of-ai-ml-enabled-medical-devices)
278 [states/blogs/2021/10/fda-publishes-approved-list-of-ai-ml-enabled-medical-devices](https://www.iqvia.com/locations/united-states/blogs/2021/10/fda-publishes-approved-list-of-ai-ml-enabled-medical-devices)
- 279 2. ICMRA. Horizon Scanning Assessment Report - Artificial Intelligence. 2021. 1–37 p.
- 280 3. Vamathevan J, Clark D, Czodrowski P, Dunham I, Ferran E, Lee G, et al. Applications of
281 machine learning in drug discovery and development. Vol. 18. 2019. 463–477 p.
- 282 4. Liu X, Cruz Rivera S, Moher D, Calvert M, Denniston AK, Spirit-ai T, et al. Reporting
283 guidelines for clinical trial reports for interventions involving artificial intelligence: the
284 CONSORT-AI extension. *Nat Med*. 2020;26(September):1364–74.
- 285 5. Rivera SC, Liu X, Chan AW, Denniston AK, Calvert MJ. Guidelines for clinical trial protocols
286 for interventions involving artificial intelligence: The SPIRIT-AI Extension. *BMJ*. 2020;370:1–
287 14.
- 288 6. Unsworth H, Dillon B, Collinson L, Powell H, Salmon M, Oladapo T, et al. The NICE Evidence
289 Standards Framework for digital health and care technologies – Developing and maintaining
290 an innovative evidence framework with global impact. *Digit Heal*. 2021;7:1–20.
- 291 7. Excellence NI for H and C. Evidence standards framework (ESF) for digital health
292 technologies [Internet]. National Institute for Health and Care Excellence. 2022 [cited 2023
293 May 9]. Available from: <https://www.nice.org.uk/corporate/ecd7>
- 294 8. Husereau D, Drummond M, Petrou S, Carswell C, Moher D, Greenberg D, et al. Consolidated
295 Health Economic Evaluation Reporting Standards (CHEERS) statement. *BMJ*.
296 2013;346(March):1–6.
- 297 9. Husereau D, Drummond M, Augustovski F, De Bekker-Grob E, Briggs AH, Carswell C, et al.
298 Consolidated Health Economic Evaluation Reporting Standards 2022 (CHEERS 2022)
299 statement: Updated reporting guidance for health economic evaluations. *BMJ*.
300 2022;376(Cheers):1–7.
- 301 10. Medicine U of OC for S in. EQUATOR Network toolkit for developing a reporting guideline
302 [Internet]. EQUATOR Network. 2018. Available from: [https://www.equator-](https://www.equator-network.org/toolkits/developing-a-reporting-guideline/)
303 [network.org/toolkits/developing-a-reporting-guideline/](https://www.equator-network.org/toolkits/developing-a-reporting-guideline/)
- 304 11. Hawsworth C, Elvidge J, Dawoud D. CHEERS-AI – Consolidated Health Economic

- 305 Evaluation Reporting Standards Artificial Intelligence Extension [Internet]. EQUATOR
306 Network. 2023. Available from: [https://www.equator-network.org/library/reporting-guidelines-](https://www.equator-network.org/library/reporting-guidelines-under-development/reporting-guidelines-under-development-for-other-study-designs/#CHEERS-AI)
307 [under-development/reporting-guidelines-under-development-for-other-study-](https://www.equator-network.org/library/reporting-guidelines-under-development/reporting-guidelines-under-development-for-other-study-designs/#CHEERS-AI)
308 [designs/#CHEERS-AI](https://www.equator-network.org/library/reporting-guidelines-under-development/reporting-guidelines-under-development-for-other-study-designs/#CHEERS-AI)
- 309 12. Manyara AM, Davies P, Stewart D, Weir CJ, Young A, Butcher NJ, et al. Protocol for the
310 development of SPIRIT and CONSORT extensions for randomised controlled trials with
311 surrogate primary endpoints: SPIRIT-SURROGATE and CONSORT-SURROGATE. *BMJ*
312 *Open*. 2022;12(10):1–8.
- 313 13. Voets MM, Veltman J, Slump CH, Siesling S, Koffijberg H. Systematic Review of Health
314 Economic Evaluations Focused on Artificial Intelligence in Healthcare: The Tortoise and the
315 Cheetah. *Value Heal*. 2022;25(3):340–9.
- 316 14. Barrett D, Heale R. What are Delphi studies? *Evid Based Nurs*. 2020;23(3):68–9.
- 317 15. Chuenjitwongsa S. How to conduct a Delphi study. Wales deanery [Internet].
318 2017;27(1173):639–43. Available from:
319 https://meded.walesdeanery.org/sites/default/files/how_to_conduct_a_delphistudy.pdf