

25 September 2023

Reliability, usability, and coverage of AMSTAR 2 assessing 544 systematic reviews and meta-analysis reports – protocol for a descriptive analytic study

Joachim Birch Milan¹ (ORCID: 0000-0001-7093-5432), Christian Gunge Riberholt² (ORCID: 0000-0002-6170-1869), Markus Harboe Olsen^{1,3} (ORCID: 0000-0003-0981-0723), Zheng Yii Lee^{4,5} (ORCID: 0000-0003-4505-7476), Johanne Pereira Ribeiro^{6,7} (ORCID: 0000-0001-6019-022X), Charles Chin Han Lew⁸ (ORCID: 0000-0001-6410-3859), and Christian Gluud^{1,9} (ORCID: 0000-0002-8861-0799)

Corresponding author: joachim.birch.milan@regionh.dk

¹ Copenhagen Trial Unit, Centre for Clinical Intervention Research, The Capital Region, Copenhagen University Hospital – Rigshospitalet, Copenhagen, Denmark

² Department of Neurorehabilitation / Traumatic Brain Injury, Copenhagen University Hospital – Rigshospitalet, Glostrup, Denmark

³ Department of Neuroanaesthesiology, The Neuroscience Centre, Copenhagen University Hospital – Rigshospitalet, Copenhagen, Denmark

⁴ Department of Anaesthesiology, Faculty of Medicine, University of Malaya, Kuala Lumpur, Malaysia.

⁵ Department of Cardiac Anesthesiology & Intensive Care Medicine, Charité Berlin, Germany

⁶ Center for Evidence-Based Psychiatry, Psychiatric Research Unit, Psychiatry Region Zealand, Region Zealand, Denmark

⁷ Department of Psychology, The Faculty of Health Sciences, University of Southern Denmark, Odense, Denmark

⁸ Department of Dietetics & Nutrition, Ng Teng Fong General Hospital, Singapore

⁹ Department of Regional Health Research, The Faculty of Health Sciences, University of Southern Denmark, Odense, Denmark

25 September 2023

6 **Abstract**

7 **Background:** The AMSTAR 2 tool (second version of Assessing Methodological Quality in Systematic
8 Reviews) is useful for critical appraisal of systematic reviews of clinical trials. In a methodological
9 study of systematic reviews and meta-analysis reports of randomised clinical trials which used the
10 sequential meta-analysis trial sequential analysis (the METSA study), we used the AMSTAR 2 to assess
11 the overall quality of each included study. With the study outlined in this protocol, we aim to explore
12 the interrater reliability of the AMSTAR 2, qualitatively describe our experiences using the tool, and
13 discuss the tool's coverage of critical domains.

14 **Methods:** In the METSA study, we investigated statistical methodology and transparency in 544
15 systematic reviews and meta-analysis reports of randomised clinical trials which used trial sequential
16 analysis (TSA). All systematic reviews (with a protocol) were assessed with AMSTAR 2 by two
17 independent authors (n=270). Meta-analysis reports – defined as not having a protocol – were
18 automatically rated as 'critically low confidence' and did not undergo further AMSTAR 2 assessment.
19 Disagreement on the AMSTAR 2 rating was resolved through discussion between the authors.
20 Principal issues were discussed at weekly meetings. Thoughts on the usability and coverage of
21 AMSTAR 2 was shared at these meetings and noted throughout and will be collected post-hoc for
22 the current study. Here, we will analyse the level of agreement on the initial ratings by raw agreement
23 rates and Cohen's kappa and test for trends concerning the effect of the consensus process (rating
24 up or down confidence) as well as the overall effect of assessor experience. We will compare the
25 AMSTAR 2 rating with the assessments of TSA transparency performed during the METSA study.

26 **Conclusion:** This methodological study will provide insights in some of the characteristics of
27 AMSTAR 2, including interrater reliability and usability in the context of assessing 270 systematic
28 reviews of clinical trials. We will provide group consensus-based suggestions regarding usability and
29 coverage.

30 **Keywords:** systematic review, AMSTAR 2, meta-analysis, evidence-based medicine

25 September 2023

31 **Introduction**

32 Systematic reviews (SR) of randomised clinical trials are generally considered the highest level of
33 evidence in clinical science (Garattini et al., 2016). The validity of SRs hinges on the methodological
34 robustness of the SR. Methodological issues in SRs prevail and therefore, a thorough, valid, and
35 systematic approach to critical appraisal of SRs is essential for evidence-based medicine (Garattini et
36 al., 2016). AMSTAR (A MeaSurement Tool to Assess systematic Reviews) represents such a systematic
37 approach and has become a popular tool in addressing issues of individual systematic reviews (De
38 Santis et al., 2023; Shea et al., 2007).

39 The revised AMSTAR tool (AMSTAR 2) was published in 2017 to increase the number of critical
40 domains covered and to be more user friendly with easier response categories and better guidance
41 (Shea et al., 2017). However, AMSTAR 2 continues to be inappropriately applied, suggesting a need
42 for even clearer guidance on AMSTAR assessment and reporting (Pieper et al., 2018).

43 In a methodological study of 544 systematic reviews (with a verifiably pre-planned protocol) and
44 meta-analysis reports (without a verifiably pre-planned protocol) of clinical trials which applied trial
45 sequential analysis (the METSA study), we investigated statistical methodology and transparent
46 reporting of trial sequential analysis (Riberholt et al., 2022). Trial sequential analysis (TSA) is a meta-
47 analysis method based on Lan-DeMets alpha spending boundaries that controls the risk of false
48 positives due to repeated testing (a concept best known from interim analyses in single trials)
49 (Wetterslev et al., 2017). We used the AMSTAR 2 to assess the overall quality of each included SR. In
50 the outlined study, we will share our experience of using the AMSTAR 2 for critical appraisal of SR in
51 the METSA study, including assessment of reliability, usability, and coverage, and hope to contribute
52 to the further development of the AMSTAR system.

25 September 2023

53 **Methods**

54 This protocol outlines a post-hoc descriptive analytic study of the AMSTAR 2 reliability, usability, and
55 coverage. The aim of the outlined study was not defined in the METSA project protocol, and the
56 methods applied are defined post-hoc of the METSA project (Riberholt et al., 2022).

57 **Data material**

58 The AMSTAR 2 assessment was performed as part of the METSA study, which is a methodological
59 study of 544 systematic reviews and meta-analysis reports of clinical trials using trial sequential
60 analysis (Riberholt et al., 2022).

61 In brief, we searched MEDLINE and the Cochrane Database for Systematic Reviews for SR and meta-
62 analysis reports of clinical trials which utilised trial sequential analysis published between January
63 2018 and January 2022. For each included study, we extracted characterising data (country of
64 publication, population, intervention, comparator, and outcomes, number of trials included, etc.) and
65 assessed the study using AMSTAR 2. For each study, we extracted data regarding TSA on one
66 dichotomous outcome analysis (n=439) and one continuous outcome analysis (n=185), if applicable
67 (total n = 624). All tasks regarding literature search, data extraction, and AMSTAR 2 assessment were
68 performed in duplicate by study authors using predefined criteria in a standardized data extraction
69 form.

70 **Method of AMSTAR 2 assessment**

71 AMSTAR 2 was incorporated in our standardised data extraction form in REDCap (Research Electronic
72 Data Capture) (Harris et al., 2009) for this project.

73 Each included study was assessed independently by two authors from the assessor group. The
74 authors assigned themselves for study assessment on an ad-hoc basis. After completed data
75 extraction for each included study, the two authors sought consensus on the final rating through
76 discussion. Persistent disagreements or principal issues were discussed and resolved at weekly
77 research meetings. AMSTAR 2 assessment was always performed before any other data extraction,
78 to minimise the impact of the latter on the former, although consensus was sought only after
79 completed data extraction.

80 Included studies that did not have a documentable pre-defined protocol were all considered meta-
81 analysis reports (MAR) and were rated as of 'critically low confidence' (274/544). The individual
82 AMSTAR 2 items were not assessed further, as further assessment would not impact the overall
83 rating. Therefore, these studies will not be included for the current analysis. Studies with a
84 documentable pre-defined protocol were considered systematic reviews (SR) and were all assessed
85 using each of the 16 items in AMSTAR 2.

86 We did not modify the AMSTAR 2 tool, however, in our data extraction form, we added an automated
87 calculator for each item concluding 'Yes', 'Partially yes', or 'No', corresponding to the original
88 AMSTAR 2 tool. To each item, we further added a multiple-choice field, e.g., 'Did the "PICO" question
89 reveal any moderate or critical weaknesses?', with the answer options 'Yes, critical weakness', 'Yes,
90 moderate weakness', and 'No'. We further added an optional comment field under each item, where
91 weaknesses noted for each item could be described. The presence of critical flaws and non-critical

25 September 2023

92 weaknesses were listed in an auto-generated table at the bottom of the form for easy overview. The
93 number of critical flaws and the occurrence of 'Yes', 'Partially yes' and 'No' were not calculated or in
94 other ways analysed, and the rating relied on an overall assessment, as is recommended by the
95 AMSTAR 2 guidance document.

96 **The assessors**

97 The AMSTAR 2 assessor group consisted of the 13 data extractors from the METSA study. At the time
98 of data extraction and consensus, four assessors were medical students, one was a medical doctor,
99 two were clinical dietitians with PhDs, two were Masters of Public Health in PhD-programmes, two
100 were physiotherapists with PhDs, one was physiotherapist in a PhD-programme, one was a
101 psychologist in a PhD-programme, and one was a medical doctor in a PhD-programme. Some
102 assessors were familiar with the AMSTAR 2 tool, however, no one in the assessor group had formal
103 experience with applying the tool, except CGR, CG, and JPR. CG and CGR instructed the assessors in
104 use of AMSTAR 2 prior to initiation of each assessor's participation in the data extraction process.

105 **Statistical analysis plan**

106 Data from the METSA database in REDCap will be exported and analysed in the latest available stable
107 version of R (R Core Team, 2022). Meta-analysis reports (without a pre-published protocol) are
108 excluded from the analysis.

109 Interrater reliability will be analysed by calculating raw agreement rates for each AMSTAR level (after
110 consensus) and Cohen's kappa and weighted kappa coefficients. We will further calculate raw
111 agreement rates for each variable in the AMSTAR assessment (each checkable answer option, the
112 calculated AMSTAR 'conclusion' for each item; 'Yes', 'Partially yes', or 'No', and the appended
113 question of whether weaknesses were identified for each item).

114 To assess whether the individual assessors had distinct rating tendencies, e.g. more positive, or
115 negative, we will compare each reviewer's initial rating with the corresponding rating of each study.
116 An initial rating that was identical to the corresponding rating will receive a score of 0, while an
117 assessment that was more positive, e.g. 'Moderate' against 'Low' will receive a score of +1, 'High'
118 against 'Critically low' will receive +3, 'Low' against 'Critically low' will receive -1, etc. We will then
119 calculate the mean for each assessor. This method may be biased as assessors may have tended to
120 co-assess with, e.g. other positive assessors. To partially account for this, we will provide an overview
121 of co-assessments in a network graph.

122 A mixed effects ordinal regression will be used to ascertain if the overall rating tendency is influenced
123 by experience gained over time, by using the rating as an ordinal outcome, assessor as random
124 effects, and accumulated number of systematic reviews assessed as a fixed effect.

125 We will analyse the change in rating after consensus was performed, to see if the consensus process
126 generally made the AMSTAR ratings more positive or more negative. For each individual study, we
127 will calculate the sum of difference (e.g. two initial ratings at 'High' and 'Low' which after consensus
128 are changed to 'Moderate' have a change of -1 and +1, respectively, with a sum of 0. If the rating
129 had been changed into 'Low', the changes would be -2 and 0, respectively, summing to -2) and

25 September 2023

130 provide a table for the frequency of each possible sum of difference (-5 to +5) subgrouped by
131 difference between initial ratings (1 to 3).

132 To test whether AMSTAR 2 is potentially insensitive to transparent reporting of statistical methods,
133 we will test the correlation between the AMSTAR rating and the TSA transparency ratings of each
134 study by ordinal regression. If possible, we will test the correlation between AMSTAR rating and the
135 secondarily collected GRADE imprecision transparency assessments of each study (protocol:
136 10.5281/zenodo.8318950).

137 For all frequentist analyses, we will not perform null hypothesis significance testing, but will calculate
138 95% confidence intervals where relevant.

139 **Qualitative evaluation of AMSTAR 2 usability and coverage**

140 All assessors were on multiple occasions encouraged to take note of issues or challenges regarding
141 usability or coverage of the AMSTAR 2 assessments and report these at the weekly meetings or in a
142 shared project document. We will read all comments in the comment fields that we added to the
143 AMSTAR 2 segment of our data extraction form to identify comments indicating assessment issues
144 or challenges. All assessors will be requested to provide feedback on coverage and usability of the
145 AMSTAR 2 tool, either in written or oral communication.

25 September 2023

146 **Reporting of results**

147 We will report and discuss the results of all conducted statistical analyses as defined in the statistical
148 analysis plan.

149 We will report raw agreement rates for each variable in the AMSTAR assessment and the overall
150 rating, and Cohen's kappa for the overall rating. We will report the frequencies of disagreement levels
151 (0 levels = no disagreement, 1 level = minor disagreement, 2 levels = major disagreement, 3 levels
152 = extreme disagreement).

153 For each reviewer, we will provide the mean difference between the reviewers initial rating and the
154 corresponding rating. We will provide a visual overview of how the individual assessors teamed up
155 in a network graph.

156 We will provide a table for the frequency of each possible sum of differences between initial ratings
157 and consensus rating (-5 to +5) subgrouped by difference between initial ratings (1-3).

158 For the ordinal regression models, we will report odds ratios (with 95% CI) and measures of
159 goodness-of-fit. We will provide a plot of each assessor's ratings in a chronological order.

160 We will further provide a qualitative description of the feedback provided during the assessment
161 process as well as feedback received after the initiation of the outlined study. We will also report if
162 any comments made directly in the AMSTAR assessment indicate an issue with usability or coverage.

25 September 2023

163 **Discussion**

164 The AMSTAR 2 tool is generally considered a useful, valid, and reliable tool for critical appraisal of
165 systematic reviews of randomised clinical trials. However, previous reports by AMSTAR 2 users
166 suggest a need for improved usability and guidance (Pieper et al., 2018). Additionally, AMSTAR 2 has
167 been critiqued for being superficial in the description of included domains, e.g. conflicts of interest
168 (Lundh et al., 2020), lacking clear reasoning behind the definitions of critical domains (Li et al., 2022)
169 and additionally lacking guidance on some domains (De Santis et al., 2023). In our outlined study,
170 we will provide a detailed discussion of our assessor groups opinions on the usability and coverage,
171 e.g. a discussion of the AMSTAR 2 instrument in relation to assessment of trial sequential analysis
172 (the focus of the METSA project) and the GRADE guidelines (Schünemann et al., 2013).

173 Systematic review methodology is a field in constant development and so, continuous updating and
174 improvement of the AMSTAR tool is warranted. The outlined study aims to contribute with insights
175 into the further development of the AMSTAR tool.

176 **Interpretation of results**

177 We will seek to identify outliers in the overview of agreement rates for each variable, as these may
178 indicate particularly challenging items or domains. From the agreement rates and kappa scores for
179 the overall rating, we will discuss the reliability of the AMSTAR 2 tool.

180 From the analyses of effects of experience and distinct rating tendencies per reviewer, we will discuss
181 the potential impact of assessor selection and the relative importance of training and experience
182 before using AMSTAR 2 for critical tasks, such as guideline formations. The findings of these analyses
183 will potentially be biased, particularly if we see that the assessors tended to co-assess with other
184 specific assessors.

185 With the table of frequencies of each possible sum of difference (-5 to +5) subgrouped by difference
186 between initial ratings (1-3), we will be able to identify common patterns in the effect of consensus
187 process, e.g. if large disagreements (e.g. 'High' and 'Critically low') commonly result in a 'compromise'
188 (e.g. 'Low', which would sum to -1) or choosing either initial rating (e.g., 'High', which would sum to
189 +3 or 'Critically low', which would sum to -3). We may also observe that the consensus rating is
190 sometimes lower than either of the initial ratings (e.g. initial ratings of 'Moderate' and 'Low' with
191 consensus rating 'Critically low', summing to -3).

192 If we find that there is no correlation between the AMSTAR rating and the TSA transparency rating
193 in each study, this could indicate that the AMSTAR tool puts too little emphasis on detailed
194 transparent reporting.

195 **Limitations**

196 During the project, the assessor group developed a discourse on AMSTAR assessment through the
197 consensus processes and the research meetings, which is expected to increase interrater agreement
198 rates. Therefore, our findings may not be applicable to agreement between naïve assessors or
199 assessors having a different group discourse.

25 September 2023

200 In the case of AMSTAR assessment, interrater agreement is not directly tied to the tool's validity. The
201 domains covered by AMSTAR 2 are complex and each individual researcher can validly hold differing
202 opinions on whether a methodological choice is a critical flaw or not.

203 **Conclusion**

204 In this protocol, we describe a planned methodological study that will quantitatively and
205 qualitatively assess reliability, usability, and coverage of the AMSTAR 2 assessment tool in the
206 context of assessing 544 systematic reviews and meta-analysis reports of clinical trials. The study
207 results will provide a basis for possibly making suggestions to recommended amendments of the
208 AMSTAR 2, contributing to the further development.

25 September 2023

209 **Additional information**

210 **Project status**

211 None of the data regarding AMSTAR in METSA has been viewed or analysed, except for the results
212 provided in this protocol (proportion of studies at each AMSTAR level).

213 **Ethical considerations**

214 The outlined study is performed on public, non-sensitive data.

215 **Author contributions**

216 JBM, CGR, MHO and CG are responsible for study conception and design. JBM drafted the protocol
217 manuscript.

218 All authors critically revised and approved the final version. The corresponding author attests that all
219 listed authors meet authorship criteria and that no others meeting the criteria have been omitted.

220 **Sources of funding and conflicts of interest**

221 Neither the outlined study nor the METSA project received external financial support. The authors
222 have nothing to declare.

223 **Data and source code availability**

224 The METSA project database is available at zenodo.org (DOI: 10.5281/zenodo.8318331). The source
225 code used for the outlined study will be made available at zenodo.org.

226 **Acknowledgements**

227 None.

25 September 2023

228 **References**

- 229 De Santis, K. K., Pieper, D., Lorenz, R. C., Wegewitz, U., Siemens, W., & Matthias, K. (2023). User
230 experience of applying AMSTAR 2 to appraise systematic reviews of healthcare interventions: a
231 commentary. *BMC Medical Research Methodology*, *23*(1), 63. [https://doi.org/10.1186/s12874-](https://doi.org/10.1186/s12874-023-01879-8)
232 [023-01879-8](https://doi.org/10.1186/s12874-023-01879-8)
- 233 Garattini, S., Jakobsen, J. C., Wetterslev, J., Bertelé, V., Banzi, R., Rath, A., Neugebauer, E. A. M., Laville,
234 M., Masson, Y., Hivert, V., Eikermann, M., Aydin, B., Ngwabyt, S., Martinho, C., Gerardi, C.,
235 Szmigielski, C. A., Demotes-Mainard, J., & Gluud, C. (2016). Evidence-based clinical practice:
236 Overview of threats to the validity of evidence and how to minimise them. *European Journal of*
237 *Internal Medicine*, *32*, 13–21. <https://doi.org/10.1016/j.ejim.2016.03.020>
- 238 Harris, P. A., Taylor, R., Thielke, R., Payne, J., Gonzalez, N., & Conde, J. G. (2009). Research electronic
239 data capture (REDCap)—A metadata-driven methodology and workflow process for providing
240 translational research informatics support. *Journal of Biomedical Informatics*, *42*(2), 377–381.
241 <https://doi.org/10.1016/j.jbi.2008.08.010>
- 242 Li, L., Asemota, I., Liu, B., Gomez-Valencia, J., Lin, L., Arif, A. W., Siddiqi, T. J., & Usman, M. S. (2022).
243 AMSTAR 2 appraisal of systematic reviews and meta-analyses in the field of heart failure from
244 high-impact journals. *Systematic Reviews*, *11*(1), 147. [https://doi.org/10.1186/s13643-022-](https://doi.org/10.1186/s13643-022-02029-9)
245 [02029-9](https://doi.org/10.1186/s13643-022-02029-9)
- 246 Lundh, A., Rasmussen, K., Østengaard, L., Boutron, I., Stewart, L. A., & Hróbjartsson, A. (2020).
247 Systematic review finds that appraisal tools for medical research studies address conflicts of
248 interest superficially. *Journal of Clinical Epidemiology*, *120*, 104–115.
249 <https://doi.org/10.1016/j.jclinepi.2019.12.005>
- 250 Pieper, D., Koensgen, N., Breuing, J., Ge, L., & Wegewitz, U. (2018). How is AMSTAR applied by authors
251 – a call for better reporting. *BMC Medical Research Methodology*, *18*(1), 56.
252 <https://doi.org/10.1186/s12874-018-0520-z>
- 253 R Core Team. (2022). *R: A Language and Environment for Statistical Computing*. [https://www.r-](https://www.r-project.org/)
254 [project.org/](https://www.r-project.org/)
- 255 Riberholt, C. G., Olsen, M. H., Milan, J. B., & Gluud, C. (2022). Major mistakes and errors in the use of
256 Trial Sequential Analysis in systematic reviews or meta-analyses – protocol for a systematic
257 review. *Systematic Reviews*, *11*(1), 114. <https://doi.org/10.1186/s13643-022-01987-4>
- 258 Schünemann, H. J., Brozek, J., Guyatt, G. H., & Oxman, A. D. (2013). *Handbook for grading the quality*
259 *of evidence and the strength of recommendations using the GRADE approach. Updated October*
260 *2013*. <https://gdt.gradepro.org/app/handbook/handbook.html>
- 261 Shea, B. J., Grimshaw, J. M., Wells, G. A., Boers, M., Andersson, N., Hamel, C., Porter, A. C., Tugwell, P.,
262 Moher, D., & Bouter, L. M. (2007). Development of AMSTAR: a measurement tool to assess the
263 methodological quality of systematic reviews. *BMC Medical Research Methodology*, *7*(1), 10.
264 <https://doi.org/10.1186/1471-2288-7-10>
- 265 Shea, B. J., Reeves, B. C., Wells, G., Thuku, M., Hamel, C., Moran, J., Moher, D., Tugwell, P., Welch, V.,
266 Kristjansson, E., & Henry, D. A. (2017). AMSTAR 2: a critical appraisal tool for systematic reviews
267 that include randomised or non-randomised studies of healthcare interventions, or both. *BMJ*,

25 September 2023

- 268 358, j4008. <https://doi.org/10.1136/bmj.j4008>
- 269 Wetterslev, J., Jakobsen, J. C., & Gluud, C. (2017). Trial Sequential Analysis in systematic reviews with
270 meta-analysis. *BMC Medical Research Methodology*, 17(1), 1–18.
271 <https://doi.org/10.1186/s12874-017-0315-7>