

1 **Trial sequential analysis for assessing imprecision in GRADE evaluations –** 2 **protocol for a methodological study**

Joachim Birch Milan¹ (ORCID: 0000-0001-7093-5432) / Johanne Periera Ribeiro^{2,3} (ORCID: 0000-0001-6019-022X), Christian Gunge Riberholt⁴ (ORCID: 0000-0002-6170-1869), Markus Harboe Olsen^{1,5} (ORCID: 0000-0003-0981-0723), and Christian Gluud^{1,6} (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 Center for Evidence-Based Psychiatry, Psychiatric Research Unit, Psychiatry Region Zealand, Region Zealand, Denmark

³ Department of Psychology, The Faculty of Health Sciences, University of Southern Denmark, Odense, Denmark

⁴ Department of Neurorehabilitation / Traumatic Brain Injury, Copenhagen University Hospital – Rigshospitalet, Glostrup, Denmark

⁵ Department of Neuroanaesthesiology, The Neuroscience Centre, Copenhagen University Hospital – Rigshospitalet, Copenhagen, Denmark

⁶ Department of Regional Health Research, The Faculty of Health Sciences, University of Southern Denmark, Odense, Denmark

3

4

5 **Abstract**

6 **Background:** Assessing statistical imprecision of summary estimates is an essential element in
7 evaluating the strength of evidence. The GRADE framework recommends assessing imprecision by
8 confidence intervals (CI) in relation to thresholds of interest and in selected cases to assess the
9 relationship between the acquired information size and the calculated optimal information size (OIS).

10 Trial sequential analysis (TSA) calculates TSA-adjusted confidence intervals after it has calculated a
11 required information size. In a recent methodological study of 544 systematic reviews and meta-
12 analysis reports of clinical trials with TSA, we gathered data regarding the methods used for grading
13 imprecision, specifically regarding the impact of TSA (the METSA project).

14 The questions regarding GRADE imprecision were initially superficially defined and were
15 substantiated only during the project and in the preparations for this protocol. With this add-on
16 study, we investigate the methods of grading imprecision in the GRADE framework by authors of
17 systematic reviews and meta-analysis reports of clinical trials with TSA.

18 **Methods:** The outlined methodological study will be pre-planned but designed with knowledge
19 about existing but not yet reviewed data on the study questions. The METSA project was not initially
20 designed for the questions raised in the current study protocol, which warrants a critical review of
21 the collected data. We aim to improve precision and accuracy of the collected data regarding
22 imprecision methodology by a redesign of selected data fields, adding new data fields to the data
23 extraction form, and a subsequent revision of the existing data extraction accordingly. For each
24 individual study, we will extract or review data regarding the specified methodology including
25 methods for calculating CI and OIS, and thresholds of interest (definitions of important benefit
26 and/or harm). For each topic, we will assess completeness in transparency of described methods and
27 protocolisation, including coherence with the protocol (if relevant).

28 **Results:** We will report frequencies of observed methods, lack of transparency, and protocolisation.
29 From data gathered in the METSA project, we will report the proportion of imprecision assessments
30 that may have differed in their conclusions if the results of the TSA had been used. Informed by our
31 findings, we will outline new suggestions on how to grade imprecision using TSA.

32 **Conclusion:** This protocol outlines a methodological study of methods and reporting characteristics
33 imprecision assessment within the GRADE framework in recent systematic reviews and meta-analysis
34 reports of clinical trials utilising TSA.

35 **Keywords:** Trial sequential analysis; imprecision; type I error; type 2 error; GRADE; systematic
36 review; meta-analysis; research on research, reproducibility

37 **Introduction**

38 Conclusions from systematic reviews with meta-analysis
39 depend on the confidence of the summary estimates (H
40 Schünemann et al., 2013). The Cochrane Handbook
41 encourages the use of The Grading of Recommendations,
42 Assessment, Development, and Evaluations framework
43 (GRADE) (HJ Schünemann, Higgins, et al., 2022) to
44 evaluate the certainty of evidence in systematic reviews,
45 clinical guidelines, population studies, and more (H
46 Schünemann et al., 2013; HJ Schünemann, Higgins, et al.,
47 2022; HJ Schünemann, Vist, et al., 2022).

48 In systematic reviews of randomised clinical trials the
49 certainty of the evidence used to generate summary
50 estimates for a given outcome is assessed with GRADE by five domains: (1) risk of bias, (2) publication
51 bias, (3) imprecision, (4) indirectness, and (5) inconsistency (Balshem et al., 2011). In the GRADE
52 framework, the certainty of the evidence is initially assumed high and subsequently downgraded to
53 either moderate, low, or very low certainty, according to the five domains. Each domain can
54 downgrade the certainty one, two, or – as recently suggested for imprecision – three levels (G H
55 Guyatt, Oxman, Kunz, Brozek, et al., 2011; G H Guyatt, Oxman, Kunz, Woodcock, et al., 2011b, 2011a;
56 G H Guyatt, Oxman, Montori, et al., 2011; Gordon H Guyatt et al., 2011; H Schünemann et al., 2013;
57 Zeng et al., 2022).

58 Imprecision is a different term for statistical uncertainty and in this context primarily arises from
59 underpowered meta-analyses, which are comparable to interim analyses of randomised clinical trials
60 (Bender et al., 2008; Brok et al., 2009; Kjaergard et al., 2001). Assessment of imprecision is a key
61 domain in the GRADE framework as recommendations about interventions cannot be made based
62 on imprecise summary estimates.

63 In the METSA project (Riberholt et al., 2022), we assessed 544 systematic reviews and meta-analysis
64 reports (see Box 1) regarding the use of trial sequential analysis (TSA). The focus of the METSA project
65 was assessing transparency and completeness in reporting of TSA and the related conclusions, e.g.
66 assessing imprecision in the GRADE framework.

67 In the outlined methodological study, we focus on the methodology and reporting characteristics of
68 imprecision assessments in the GRADE framework in systematic reviews and meta-analysis reports
69 with TSA with the purpose of contributing to the further development and promotion of the GRADE
70 framework.

71

72

Box 1

We define a *systematic review* (SR) as a protocolised approach to evidence synthesis, using verifiably pre-defined eligibility criteria and synthesis methodology.

We define a *meta-analysis report* (MAR) as methodologically comparable to a systematic review but lacking a predefined publicly available protocol.

73 **Methods**

74 **Study design**

75 This protocol outlines a pre-planned, secondary report on the current practice of grading imprecision
76 within the GRADE framework in systematic reviews and meta-analysis reports of clinical trials, which
77 utilised TSA. The aim of the outlined study was defined in the METSA project protocol (Riberholt et
78 al., 2022), but the methods applied are defined post-hoc.

79

80 **Data material**

81 We will adapt the existing METSA project database. The METSA project database, which is available
82 at zenodo.org (DOI: 10.5281/zenodo.8318331), contains data on 544 systematic reviews and meta-
83 analysis reports, 300 (55%) of which applied GRADE. A complete description of the METSA project
84 methodology is available at (Riberholt et al., 2022) [REF also? to main paper See my comment above].
85 In brief, we searched MEDLINE and the Cochrane database for systematic reviews and meta-analysis
86 reports of randomised clinical trials which utilised TSA published between January 2018 and January
87 2022. For each included study, we extracted baseline data and assessed the study using AMSTAR 2
88 (Shea Beverley et al., 2017). We extracted data regarding TSA on one dichotomous outcome analysis
89 (n=439) and one continuous outcome analysis (n=185), if applicable (total TSAs n = 624). All tasks
90 regarding literature search, data extraction, and AMSTAR assessment were performed in duplicate
91 by study authors using predefined criteria in a standardised data extraction form.

92 The METSA project was not initially designed for the questions raised in the current study protocol,
93 and upon data revision, we have found that a critical review of the collected data is warranted.

94

95 **Data extraction**

96 All tasks regarding data revision and extraction will be performed in duplicate by two independent
97 authors. Discrepancies will be resolved through a consensus process. A third author will be involved
98 in case of disagreements.

99 We aim to improve precision and accuracy of the collected data regarding imprecision methodology
100 by a redesign of selected data fields, adding new data fields to the data extraction form and a
101 subsequent revision of the existing data extraction accordingly. The amended data extraction form
102 is provided in Supplement. Most changes concern specifications regarding the specified
103 methodology including methods for calculating CI and RIS, and thresholds of interest (definitions of
104 important benefit and/or harm). For each topic, we will assess completeness in transparency of
105 described methods and protocolisation, including coherence with the protocol (if relevant).

106

107 **Analysis and presentation of results**

108 We will report and describe the observed methods for downgrading imprecision, lack of
109 transparency, and protocolisation and report the frequencies of each. The data extraction form has
110 predefined options for certain fields but not for other. We will assess the need for revising the defined
111 options/categories based on the collected data. Specifically, we will report:

- 112 - Frequencies of publications with each identified method for downgrading imprecision, including
113 frequencies of unclear methodology
 - 114 - For relevant methods, the frequency of setting limits for important differences for benefit
115 and harm, respectively
 - 116 - For relevant methods, the frequency of reporting OIS calculation methods (if applicable,
117 we will provide an overview of identified methods)
- 118 - Frequencies of publications in which the specified methodology:
 - 119 - Was not planned in a protocol made public prior to data extraction
 - 120 - Differed from the methodology described in the protocol
- 121 - Frequencies of outcomes downgraded 0, 1, 2, or 3 levels respectively, including frequencies of
122 unclear or non-reproducible gradings of imprecision.

123

124 From data gathered in the METSA project, we will additionally report the proportion of imprecision
125 assessments that may have differed if they had assessed imprecision using the results of their TSA
126 analysis, i.e. had used TSA CI instead of conventional 95% CI (only regarding inclusion of no effect,
127 i.e. not important benefits or harms) or an RIS instead of OIS for downgrading imprecision. We will
128 seek to demonstrate how TSA may contribute to downgrading imprecision with practical examples
129 from the included studies in relation to a discussion of the statistical theory behind TSA. Informed
130 by our findings, we will develop suggestions on how to assess imprecision in the GRADE framework
131 using TSA, including reporting standards. We will discuss these new suggestions in relation to
132 previous suggestions of using TSA for GRADE.

133 Deviations from the outlined analysis and presentation plan will be described in the final study
134 report.

135

136 **Discussion**

137 In this protocol, we outline a methodological study in which we will assess and report the
138 methodology of downgrading imprecision in the GRADE framework in recent systematic reviews
139 and meta-analysis reports of randomised clinical trials which use TSA. The aim of the study is to
140 bring attention to the important topic of imprecision by providing an overview of the methods
141 used in current literature and potentially identifying inadequacies in methodology and reporting. It
142 is our goal to contribute to the further development and promotion of the GRADE framework.

143 This study is not without limitations. Three study authors (JBM, JMPP, CGR) participated in
144 extracting, revising, and reviewing data for the METSA project. The answers to the questions raised
145 in this protocol are partially known by these authors, e.g. there is a significant lack of transparency
146 in the utilised methods, but with uncertainty due insufficient data accuracy and precision. These
147 observations have not been shared elsewhere. The outlined study seeks to answer new questions
148 as well as confirming the assumptions that were based on observations made during the METSA
149 project.

150

151 **Conclusion**

152 This protocol outlines a methodological study of method and reporting of imprecision assessment
153 within the GRADE framework in recent systematic reviews and meta-analysis reports of clinical trials
154 utilising trial sequential analysis.

155

156 **Additional information**

157 **Project status**

158 Data revision regarding the variables relevant for the outlined study has been initiated (n=25) for the
159 purpose of testing the data extraction form.

160 **Ethical considerations**

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

162 **Author contributions**

163 JPR and JBM are responsible for the design of the outlined study and drafted the protocol
164 manuscript.

165 All authors contributed to the design of the outlined study, and critically revised and approved the
166 final version of the protocol manuscript. The corresponding author attests that all listed authors meet
167 authorship criteria and that no others meeting the criteria have been omitted.

168 For contributions to the METSA database, we refer to the METSA publication (awaiting publication).

169 **Sources of funding and conflicts of interest**

170 Neither the outlined study nor the METSA project received external financial support. Java and R
171 implementations of trial sequential analysis was developed by the Copenhagen Trial Unit, for R
172 directed by Anne Lyngholm Sørensen (<https://orcid.org/0000-0002-8265-0394>), PhD student at
173 Section for Biostatistics, Institute of Public Health, Copenhagen University. The authors have nothing
174 to declare.

175 **Data and source code availability**

176 The METSA project database is available at zenodo.org (DOI: 10.5281/zenodo.8318331). The
177 amended METSA project database containing data used for the outlined study will be made available
178 along with any relevant code for analysis at zenodo.org.

179 **Acknowledgements**

180 We are grateful to Mark Asante (<https://orcid.org/0009-0002-8034-4139>) and Buddhheera
181 Kumburegama for contributing to the data revision for this study.

182 **References**

- 183 Balshem, H., Helfand, M., Schünemann, H. J., Oxman, A. D., Kunz, R., Brozek, J., Vist, G. E., Falck-Ytter,
184 Y., Meerpohl, J., Norris, S., & Guyatt, G. H. (2011). GRADE guidelines: 3. Rating the quality of
185 evidence. *Journal of Clinical Epidemiology*, 64(4), 401–406.
186 <https://doi.org/10.1016/j.jclinepi.2010.07.015>
- 187 Bender, R., Bunce, C., Clarke, M., Gates, S., Lange, S., Pace, N. L., & Thorlund, K. (2008). Attention
188 should be given to multiplicity issues in systematic reviews. *Journal of Clinical Epidemiology*,
189 61(9), 857–865. <https://doi.org/10.1016/J.JCLINEPI.2008.03.004>
- 190 Brok, J., Thorlund, K., Wetterslev, J., & Gluud, C. (2009). Apparently conclusive meta-analyses may be
191 inconclusive - Trial sequential analysis adjustment of random error risk due to repetitive testing
192 of accumulating data in apparently conclusive neonatal meta-analyses. *International Journal of*
193 *Epidemiology*, 38(1), 287–298. <https://doi.org/10.1093/ije/dyn188>
- 194 Guyatt, G H, Oxman, A. D., Kunz, R., Brozek, J., Alonso-Coello, P., Rind, D., & Et, a I. (2011). GRADE
195 guidelines: 6. Rating the quality of evidence - imprecision. *Journal of Clinical Epidemiology*,
196 64(12), 1283–93–1283–93. <https://doi.org/10.1016/j.jclinepi.2011.01.012>
- 197 Guyatt, G H, Oxman, A. D., Kunz, R., Woodcock, J., Brozek, J., Helfand, M., & et, a I. (2011a). GRADE
198 guidelines: 7. Rating the quality of evidence - inconsistency. *Journal of Clinical Epidemiology*,
199 64(12), 1294–302–1294–302.
- 200 Guyatt, G H, Oxman, A. D., Kunz, R., Woodcock, J., Brozek, J., Helfand, M., & et, a I. (2011b). GRADE
201 guidelines: 8. Rating the quality of evidence - indirectness. *Journal of Clinical Epidemiology*,
202 64(12), 1303–10–1303–10.
- 203 Guyatt, G H, Oxman, A. D., Montori, V., Vist, G., Kunz, R., Brozek, J., & et, a I. (2011). GRADE guidelines:
204 5. Rating the quality of evidence - publication bias. *Journal of Clinical Epidemiology*, 64(12),
205 1277–82–1277–82.
- 206 Guyatt, Gordon H, Oxman, A. D., Vist, G., Kunz, R., Brozek, J., Alonso-Coello, P., Montori, V., Akl, E. A.,
207 Djulbegovic, B., & Falck-Ytter, Y. (2011). GRADE guidelines: 4. Rating the quality of evidence—
208 study limitations (risk of bias). *Journal of Clinical Epidemiology*, 64(4), 407–415.
209 <https://doi.org/10.1016/j.jclinepi.2010.07.017>
- 210 Kjaergard, L. L., Villumsen, J., & Gluud, C. (2001). Reported Methodologic Quality and Discrepancies
211 between Large and Small Randomized Trials in Meta-Analyses. *Annals of Internal Medicine*,
212 135(11), 982–989. <https://doi.org/10.7326/0003-4819-135-11-200112040-00010>
- 213 Riberholt, C. G., Olsen, M. H., Milan, J. B., & Gluud, C. (2022). Major mistakes and errors in the use of
214 Trial Sequential Analysis in systematic reviews or meta-analyses – protocol for a systematic
215 review. *Systematic Reviews*, 11(1), 114. <https://doi.org/10.1186/s13643-022-01987-4>
- 216 Schünemann, H, Brožek, J., Guyatt, G., & Oxman, A. (2013). *GRADE handbook for grading quality of*
217 *evidence and strength of recommendations. Available from guidelinedevelopment.org/handbook*
218 (Holger Schünemann, J. Brožek, G. Guyatt, & A. Oxman (eds.); Updated Oc). The GRADE Working
219 Group.
- 220 Schünemann, HJ, Higgins, J., Vist, G., Glasziou, P., Akl, E., Skoetz, N., & Guyatt, G. (2022). Chapter 14:

- 221 Completing 'Summary of findings' tables and grading the certainty of the evidence. In J. Higgins,
222 J. Thomas, J. Chandler, M. Cumpston, T. Li, M. Page, & V. Welch (Eds.), *Cochrane Handbook for*
223 *Systematic Reviews of Interventions*. Cochrane.
- 224 Schünemann, HJ, Vist, G., Higgins, J., Santesso, N., Deeks, J., Glasziou, P., Akl, E., & Guyatt, G. (2022).
225 Chapter 15: Interpreting results and drawing conclusions. In J. Higgins, J. Thomas, J. Chandler,
226 M. Cumpston, T. Li, M. Page, & V. Welch (Eds.), *Cochrane Handbook for Systematic Reviews of*
227 *Interventions*. Cochrane.
- 228 Shea Beverley, J., Reeves Barnaby, C., Wells, G. e. o. r. g. e., Thuku, M. i. c. e. r. e., Hamel, C. a. n. d. y.
229 c. e., Moran, J. u. l. i. a. n., Moher, D. a. v. i. d., Tugwell, P. e. t. e. r., Welch, V. i. v. i. a. n., &
230 Kristjansson, E. l. i. z. a. b. e. t. h. (2017). AMSTAR 2: a critical appraisal tool for systematic reviews
231 that include randomised or non-randomised studies of healthcare interventions, or both. *Bmj*,
232 358, j4008–j4008.
- 233 Zeng, L., Brignardello-Petersen, R., Hultcrantz, M., Mustafa, R. A., Murad, M. H., Iorio, A., Traversy, G.,
234 Akl, E. A., Mayer, M., Schünemann, H. J., & Guyatt, G. H. (2022). GRADE Guidance 34: update on
235 rating imprecision using a minimally contextualized approach. *Journal of Clinical Epidemiology*,
236 150(20), 216–224. <https://doi.org/10.1016/j.jclinepi.2022.07.014>
- 237
- 238