



UNIVERSITY OF LEEDS

This is a repository copy of *Risk-of-bias assessment using RoB2 was useful but challenging and resource-intensive: observations from a systematic review.*

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/200718/>

Version: Accepted Version

Article:

Crocker, T orcid.org/0000-0001-7450-3143, Lam, N orcid.org/0000-0001-8591-444X, Jordão, M orcid.org/0000-0003-2108-2677 et al. (6 more authors) (2023) Risk-of-bias assessment using RoB2 was useful but challenging and resource-intensive: observations from a systematic review. *Journal of Clinical Epidemiology*, 161. pp. 39-45. ISSN 0895-4356

<https://doi.org/10.1016/j.jclinepi.2023.06.015>

© 2023, Elsevier. This is an author produced version of an article published in *Journal of Clinical Epidemiology*. Uploaded in accordance with the publisher's self-archiving policy. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 Risk-of-bias assessment using RoB2 was useful but challenging and
2 resource-intensive: observations from a systematic review

3 [Author names and affiliations](#)

4 Thomas Frederick Crocker^a, Natalie Lam^{a1}, Magda Jordão^a, Caroline Brundle^a, Matthew
5 Prescott^a, Anne Forster^a, Joie Ensor^{b2}, John Gladman^c, Andrew Clegg^a

6 ^aAcademic Unit for Ageing and Stroke Research (University of Leeds), Bradford Institute for Health
7 Research, Bradford Teaching Hospitals NHS Foundation Trust, Bradford, UK

8 ^bCentre for Prognosis Research, Keele School of Medicine, Keele University, Keele,
9 Staffordshire, UK

10 ^cCentre for Rehabilitation & Ageing Research, Academic Unit of Injury, Inflammation and Recovery
11 Sciences, University of Nottingham and Health Care of Older People, Nottingham University
12 Hospitals NHS Trust, Nottingham, UK

13 Corresponding author: Thomas Frederick Crocker

14 Bradford Institute for Health Research, Bradford Royal Infirmary, Duckworth Lane, Bradford,
15 BD9 6RJ. UK

16 e: tom.crocker@bthft.nhs.uk

17 t: +44 1274 383406

18 Present addresses:

19 1. Mental Health and Addiction Research Group, Department of Health Sciences, Faculty of
20 Science, ARRC Building, University of York, Heslington, York, UK

21 2. Institute of Applied Health Research, College of Medical and Dental Sciences, University of
22 Birmingham, Birmingham, UK

23 [Declaration of interests:](#)

24 Andrew Clegg declares funding through the NIHR HTA programme, NIHR Programme Grants for
25 Applied Research, NIHR HS&DR programme, NIHR Applied Research Collaboration Yorkshire &
26 Humber, and Health Data Research UK; Anne Forster declares NIHR Senior Investigator Award 2017-
27 present, NIHR Programme Grant 10% of salary, NIHR HS&DR grant 8% of salary, HTA grant 5% of
28 salary, National Institute for Health (USA) payment for panel membership 2021, 2022, participation
29 in Programme Steering Committees for NIHR 202339 Improving the lives of stroke survivors with
30 data, and NIHR202020 Research Title Personalised Exercise-Rehabilitation FOR people with Multiple
31 long-term conditions (multimorbidity)-The PERFORM trial, University of Leeds Governor
32 representative on the Governors Board of Bradford Teaching Hospitals NHS Foundation Trust,
33 member of HSDR Researcher-Led panel, member of NIHR Doctoral Fellowship Panel member of
34 Policy Research Unit assessment panel. Other authors declare no potential conflicts of interest.

1 Abstract

2 Objective

3 To report our experience using Version 2 of the Cochrane risk-of-bias tool for randomised trials
4 (RoB2).

5 Study design and setting

6 Two reviewers independently applied RoB2 to results of interest in a large systematic review of
7 complex interventions and reached consensus. We recorded time taken, and noted and discussed
8 our difficulties using the tool, and the resolutions we adopted. We explored time taken with
9 regression analysis and summarised our experience of implementing the tool.

10 Results

11 We assessed risk of bias in 860 results of interest in 113 studies. Staff resource averaged 358
12 minutes per study (SD 183). Number of results ($\beta=22$) and reports ($\beta=14$) per study and experience
13 of the team ($\beta=-6$) significantly affected assessment time. To implement the tool consistently we
14 developed cut-points for missingness and considerations of balance regarding missingness, assumed
15 some concerns with intervention deviations unless otherwise prevented or investigated, some
16 concerns with measurements from unblinded self-reporting participants, and judged low risk of
17 selection for certain dichotomous outcomes despite the absence of an analysis plan.

18 Conclusion

19 The RoB2 tool and guidance are useful but resource-intensive and challenging to implement. Critical
20 appraisal tools and reporting guidelines should detail risk-of-bias implementation. Improved
21 guidance focussing on implementation could assist reviewers.

22

1 **Keywords**

2 RoB2; risk of bias; process duration; research methods; certainty assessment; systematic reviews

3 **Running title**

4 Risk-of-bias assessment using RoB2 was useful but challenging and resource-intensive

5 **Word count**

6 3085

7 **What is New?**

8 **Key findings**

9 It took 5h 58m of staff time to assess risk of bias per study on average using version 2 of the
10 Cochrane tool for assessing risk-of-bias in randomised trials (RoB2) in a systematic review of
11 community-based complex interventions.

12 Variation in time per study was largely explained by models including the number of results of
13 interest, the experience of the reviewers with use of the RoB2 tool, and, for individual assessments,
14 number of reports.

15 Despite the extensive guidance we had difficulty implementing aspects of most domains.

16 **What this adds to what is known**

17 This is the first research to identify the impact of number of results of interest and number of
18 reports on the time taken to assess RoB2, and adds to a limited body of published evidence about
19 how reviewers have implemented the RoB2 tool.

20 **What is the implication, what should change now**

21 Improved guidance is needed to further assist review authors in implementing the RoB2 tool and to
22 encourage reporting details of their approach to implementation.

1 Introduction

2 Systematic reviews that synthesise evidence of the effectiveness of interventions are a cornerstone
3 of clinical guidelines and evidence-based medicine. Evaluating risk of bias is an essential element of
4 systematic reviews and one step towards establishing the degree of confidence or certainty in the
5 synthesis [1, 2]. Methods for evaluating risk of bias have evolved from study quality checklists to an
6 increasing focus on factors that affect the internal validity of the results of studies. In 2008 the
7 Cochrane Collaboration published a tool to evaluate risk of bias in randomised controlled trials
8 (RCTs) which became the standard for such assessments [3-5].

9 Version 2 of the Cochrane risk-of-bias tool for randomised trials (RoB2) was published in 2019 [6, 7].
10 This revised version, developed through an expert consensus process, introduced substantial
11 changes. For example, each relevant result of a study was now assessed instead of the study as a
12 whole and an overall judgement was to be reached. There was some restructuring of the domains of
13 bias in response to theoretical developments. Additionally, a series of signalling questions (SQs) and
14 accompanying algorithm were provided for each domain to try to improve reviewers' agreement [8].

15 RoB2 has been widely cited in the three years since its publication (over 8000 citations according to
16 Google Scholar) suggesting widespread uptake. A recent meta-epidemiological study identified 196
17 completed systematic reviews that applied RoB2 [9]. Two studies have estimated the time taken to
18 apply RoB2, finding it demanding with problematic reliability, and recommending development of
19 operational criteria specific to the review to improve implementation [10, 11].

20 We used RoB2 to assess risk of bias in a large, robust systematic review with network meta-analysis.
21 This article reports our experiences, including details to help reviewers planning to use RoB2, how
22 we operationalised aspects of it, and recommendations for those maintaining and developing the
23 tool.

24 2 Methods

25 2.1 Overarching systematic review

26 We undertook a systematic review and network meta-analysis of community-based complex
27 interventions to sustain independence in older people. The review followed standard procedures,
28 was prospectively registered on PROSPERO (CRD42019162195) and the protocol published [12].

29 Briefly, the methods of relevance to our use of RoB2 were as follows.

1 Following a comprehensive search, we included randomised controlled trials (RCTs) or cluster-RCTs
2 that measured outcomes at least 24 weeks after baseline. Participants were older people living at
3 home (≥ 65 years). Eligible interventions were community-based complex interventions targeted at
4 the individual, focused on sustaining independence. Eligible comparators were usual care, “placebo”
5 or attention control, or a different complex intervention which met our criteria.

6 Two researchers independently screened records (title and abstract), assessed eligibility, and
7 extracted data.

8 Our outcomes of interest comprised six dichotomous outcomes and seven continuous outcomes.
9 Each comparison between two trial arms (e.g. experimental and control interventions) for an
10 outcome of interest at a particular timepoint for which an effect estimate was reported or could be
11 calculated is referred to henceforth as a *result of interest*. Results of interest were extracted for
12 three timeframes.

13 2.2 Risk-of-bias assessment using RoB2

14 Two reviewers independently assessed risk of bias (RoB) in each result of interest from each
15 included study, using version 2 of the Cochrane risk-of-bias tool for randomised trials (RoB2) [6-8].
16 We were interested in the effect of assignment to the intervention (‘intention-to-treat’ effect).
17 Disagreements were resolved by consensus between the reviewers or through discussion with the
18 Programme Management Group (PMG) which included expert clinicians, trialists and statisticians.

19 For individually randomised studies, we assessed risk of bias in five domains: (1) the randomization
20 process; (2) deviations from intended interventions; (3) missing outcome data; (4) measurement of
21 the outcome; and (5) selection of the reported result. Domain 1 was assessed at the study level and
22 the other domains were assessed at the result level. Cluster-RCTs, were assessed similarly, except
23 with two domains of allocation bias in place of domain 1: (1a) the randomization process, and (1b)
24 identification or recruitment of participants [13].

25 For each domain, we made a judgement of high risk of bias, low risk of bias, or some concerns. We
26 used the SQs and algorithm and considered whether to override the algorithm result, recording our
27 reasons and supporting evidence. We reached an overall judgement at least as severe as the most
28 severe domain risk.

29 Reviewers who conducted RoB2 assessments read the guidance and watched the Cochrane RoB2:
30 Learning Live webinars [8, 13, 14]. Two reviewers had previous experience of using the original
31 Cochrane risk-of-bias tool (TC, NL) and three were new to assessing risk of bias (CB, MJ, MP).

1 2.3 Evaluation of RoB2 usage

2 We recorded details of the time taken per study to conduct assessments (per reviewer) and reach
3 consensus. We discussed and noted our difficulties with using the tool, and the resolutions we
4 adopted.

5 We produced graphs, summary statistics, and conducted multivariable linear regression to explore
6 whether time taken to assess risk of bias for a study (per individual, per consensus meeting, overall)
7 was influenced by:

- 8 • the number of results to be assessed (increasing time);
- 9 • experience using RoB2 (decreasing time);
- 10 • number of reports per study (increasing time).

11 Additionally, we anticipated variability between individual reviewers and so included this as a
12 categorical variable for regressions of individual assessment time.

13 In these analyses we only included studies with time data for two individual assessments. We
14 calculated resource used in person-hours, and full-time equivalent (FTE) as one person working 7.5
15 hours per day, five days per week, with 33 days leave per year (including public holidays): 142.3
16 hours per month.

17 We summarised our approach and reasoning to RoB2 implementation.

18 3 Results

19 3.1 Overarching systematic review

20 Our literature searches produced 40,112 records after deduplication. We included 129 studies
21 consisting of 496 reports, of which 113 reported results of interest (see Appendix A). Among these
22 studies there were 860 results of interest for which we assessed risk of bias,¹ ranging from 1 to 33
23 per study (median 6).

24 3.2 RoB2 Assessments

25 In every result of interest we judged there to be at least some concerns about overall risk of bias
26 (28%), with 72% at high risk of bias. A description of risk of bias by domain is provided in Appendix B.

¹ 34 were unsuitable for inclusion, 826 results are presented in Appendix B.

1 3.3 RoB2 assessment process time

2 Mean time per study to conduct an individual assessment (per reviewer) was 127 minutes (2h 7m;
 3 SD 67) and 54 minutes (SD 43) for the consensus meeting (see Table 1). We had complete timing
 4 data for 99 studies; overall these included 35,472 minutes of worktime (591.2 person-hours or 2.1
 5 months of 2 x FTE work) including each individual assessment and two people in a consensus
 6 meeting (5h 58m per study, 47m per result).

7 **Table 1. Descriptive data for time taken to conduct risk of bias assessments and consensus meetings.**

	Individual assessments [*]	Consensus meetings [†]	Resource for overall process [‡]
Studies with complete data (n)	106 [§]	105	99 [¶]
Results of interest per study	7.4 (5.7) [1 to 33]	7.8 (6.0) [1 to 33]	7.6 (5.8) [1 to 33]
Reports per study	3.8 (2.8) [1 to 12]	3.9 (2.8) [1 to 12]	3.9 (2.8) [1 to 12]
Time per study (mins)	127 (67) [32 to 421]	54 (43) [6 to 271]	358 (183) [98 to 976]
Time per result of interest (mins) [#]	17	6.9	47

8 Data are presented as mean (SD) [min to max] unless otherwise stated.

9 * Studies with timing data for two complete individual assessments only. Times are presented per reviewer so total
 10 resource use is double.

11 † Studies with timing data for consensus meetings (following two independent individual assessments). Times are
 12 presented for the meeting in which two reviewers were present so total resource use is double.

13 ‡ Studies with timing data for two complete individual assessments and a consensus meeting. Times are presented for
 14 total resource use counting both reviewers.

15 § Seven studies missing: four with missing timing data for at least one reviewer, and three with more than two reviewers
 16 involved in the assessment.

17 || Eight studies with missing timing data for the consensus meeting.

18 ¶ Fourteen studies not analysable for the overall process. Of the 106 studies with timing data for (only) two complete
 19 individual assessments, seven did not have timing data for the consensus meeting.

20 # Only means are presented as times were recorded per study.

21 3.3.1 Factors influencing process time

22 Graphs and results of regression analyses are presented in Appendix C. The number of results of
 23 interest per study affected time to conduct individual assessments and consensus meetings, adding
 24 22 minutes per result (95% confidence interval (CI) 18 to 26) to overall worktime. Number of reports
 25 affected time to conduct individual assessments but not the consensus meetings, adding 14 minutes
 26 per report (95% CI 6 to 22) to overall worktime. Experience reduced time taken to conduct individual
 27 assessments and consensus meetings. Additionally, there was substantial variation between
 28 individual reviewers.

1 Thirty-two studies were assessed by two experienced reviewers who had each previously assessed
2 at least 25 studies in this review using RoB2. For these studies, overall worktime was 5 hours 15
3 minutes per study or 44 minutes per result on average. Regression analysis estimated this to be 178
4 minutes per study (2h 58m; 95% CI 139 to 218) plus 19 minutes per result (95% CI 15 to 24)
5 (adjusted R^2 .73).

6 3.4 Challenges implementing RoB2

7 Although the guidance for RoB2 is extensive, reviewers found certain aspects of the guidance to be
8 lacking specificity sufficient to operationalise it.

9 3.4.1 Deviations from the intended interventions

10 We were uncertain what evidence should be considered sufficient to indicate there were ‘probably
11 no’ deviations from the intended intervention that arose because of the trial context (SQ 2.3). This
12 signalling question often decided the judgement for the domain because we were unable to answer
13 that both participants and personnel were unaware of their allocation for any of these complex
14 intervention studies (SQ 2.1 and 2.2). The guidance provided an example of trial enrolment and
15 randomisation potentially leading control participants to feel in need and unlucky and thus seek
16 other interventions they would not have otherwise. This seemed like a risk that was always plausible
17 although, perhaps, often unlikely, and one that was very difficult for trialists to investigate.

18 We answered ‘probably no’ for four studies: three where authors had specifically investigated
19 control group behaviours and found no concerns; and one where we considered the stepped-wedge
20 design would make such deviation unlikely. We answered ‘no information’ for most studies,
21 answering ‘probably yes’ for nine studies with evidence of contamination.

22 3.4.2 Missing outcome data

23 We found several difficulties in implementing the guidance regarding assessment of missing
24 outcome data. Firstly, standardising assessment of ‘nearly all’ data being available (SQ 3.1). We
25 operationalised this as at most 5% of participants missing as a proportion of: the number allocated
26 for continuous outcomes; or, recorded cases for binary outcomes (e.g. 1 person missing if 20 people
27 died regardless of sample size).

28 SQ 3.2 asks “Is there evidence that the result was not biased by missing outcome data?” [8]. We
29 judged there was insufficient evidence, even when study authors had conducted multiple imputation
30 with multiple measured variables to correct for bias. Our outcomes of interest were so distal it was
31 unlikely *all* relevant variables were included in such models (we were uncertain whether this could
32 ever be a plausible assumption).

1 For almost all results where some data was missing it was plausible that missingness *could* have
2 depended on the true value of the outcome (SQ 3.3). This was because outcomes were health
3 outcomes and reasons for missingness typically included mortality, care home admission or
4 withdrawal. SQ 3.4 delineates between a judgement of some concerns and high risk. It asks whether
5 it is *likely* missingness depended on the true value of the missing result; the elaboration indicates we
6 should answer '(probably) yes' if "Reported reasons for missing outcome data provide evidence that
7 missingness in the outcome depends on its true value" [8]. Therefore, when we first applied the tool
8 we invariably answered 'probably yes' to SQ 3.4 based on our reason for SQ 3.3. Almost all results
9 were therefore judged high risk. It seemed inappropriate and at odds with the preceding guidance
10 that we would judge a high risk of bias due to missingness when losses were relatively small and
11 balanced in numbers and reasons; for example, for a continuous outcome where 5.1% of
12 participants had died in both arms.

13 To decide between a judgement of some concerns and high risk we first considered whether more
14 than 45% of participants/cases were missing overall, regardless of how balanced, as an arbitrary
15 upper limit. Secondly, we considered whether differences in total numbers or reasons missing
16 differed by our threshold values of 5%, or if insufficient detail regarding reasons was given. In any of
17 these cases we would make a judgement of high risk. For remaining results where between 5% and
18 45% of participants/cases were missing overall we sought increasing balance in the numbers and
19 reasons of missingness to judge the result some concerns only.

20 We decided against answering 'no information' to SQ 3.4 (which would have resulted in a judgement
21 of high risk) in the common situation where numbers and reasons for losses were presented by
22 group for a trial overall but not for each specific result of interest, although we were uncertain this
23 was how we were supposed to interpret the guidance.

24 3.4.3 Measurement of the outcome

25 For participant-reported outcomes (PROs) ("involving judgement"), we had to judge the likelihood
26 that knowledge of the intervention influenced assessment (high risk) or not (some concerns) (SQ
27 4.5). Based on the guidance we considered the degree of participant judgement for each outcome,
28 whether certain interventions (e.g. alternative medicine) or closely related outcomes (e.g. ADL
29 training for ADL outcomes) should indicate high risk, and whether studies with active comparators
30 should only be rated 'some concerns' despite some of these features. We were concerned that such
31 judgements were predicated on assumptions rather than evidence, risking our own biases
32 influencing the assessment. In the absence of stronger empirical evidence regarding the factors
33 influencing risk of bias for PROs, we decided that knowledge of the intervention could influence

1 assessment but that this was unlikely (some concerns in the absence of other problems). Our
2 reasoning was that unblinded observers were a greater risk than unblinded self-reporting
3 participants, many of the interventions may not be recognised as departures from usual care by
4 participants, and measurements would often be temporally distant from intervention receipt. We
5 were able to judge this domain low risk for some PRO results for which we concluded that the
6 participants were probably unaware of the intervention received, usually in cluster trials.

7 3.4.4 Selection of the reported result

8 Domain five considers selective reporting of results. To reach a judgement of low risk, the algorithm
9 requires the result of interest to have been analysed in accordance with a pre-specified analysis
10 plan. Such plans are rarely available. When a dichotomous outcome that could only be measured in
11 one way (e.g. mortality) was reported as the number of participants, we considered that this was
12 equivalent to a situation where we gathered and reanalysed individual patient data, and therefore
13 judged the risk as low, regardless of a pre-specified plan.

14 4 Discussion

15 This article detailed the substantial resource taken to conduct RoB2 assessments in a systematic
16 review (358 mins/study) and estimated that multiple results and study reports increased time per
17 study, while experience reduced time taken. We have detailed the way we operationalised guidance
18 for missing outcome data, deviations due to trial context, knowledge of the intervention influencing
19 assessment, and selection of the reported result. Overall, we found the signalling questions and
20 algorithms helpful. However, sections of the guidance seemed unnecessarily discursive and
21 theoretical with insufficient practical advice for interpreting the SQs. We also found the wording for
22 some SQs misleading. Additionally, completing the SQs and the supporting free-text boxes for each
23 SQ and each result in each study was time-consuming.

24 Our results relate to one particularly challenging systematic review with network meta-analysis of
25 pragmatic trials of complex interventions that required deliverer and participant involvement.
26 Therefore, the time taken and some of our difficulties may not manifest for other reviews. However,
27 the findings are from a real-world review without special support from the Cochrane Methods
28 Support Unit. The order in which we reviewed studies was not randomly selected and so estimates
29 of the effect of experience are limited in this regard. We did not attempt to clarify uncertainties
30 regarding risk of bias with study authors, as recommended for domain 5 [8]; such action may have
31 affected the time we took as well as our judgments. The approaches we took to implementing RoB2
32 were based on extensive reading of the guidance and discussion among the authors who comprise

1 experienced trialists and systematic reviewers. Nonetheless, they may not fit with the intentions of
2 the tool authors, so we advise caution in following these approaches.

3 We are aware of only one other study that reports upon the implementation of the RoB2 tool in a
4 systematic review: Minozzi et al. analysed the impact an implementation document had on times
5 and inter-rater reliability of their RoB2 assessments in a review of cannabis and cannabinoids for
6 people with multiple sclerosis [11]. The supplementary implementation document details their
7 approach to similar challenges. Like Minozzi et al. we identified an improvement in speed over time,
8 although to a lesser extent. We additionally estimated how number of results and reports per study
9 affected time. Minozzi et al. implemented a 90% cut-point for judging nearly all data being available,
10 while we implemented a 5% rule relating to missingness (SQ 3.1). More specific guidance from the
11 RoB2 tool authors would help to develop consistency between review teams, in the absence of
12 which reviewers are likely to develop their own cut-points. We decided that none of the analyses
13 that attempted to correct for missing outcome data were sufficient to judge a low risk of bias
14 whereas Minozzi et al. accepted these. For bias due to assessment of the outcome being influenced
15 by knowledge of the intervention received we both decided that participant reported outcomes
16 requiring judgement should be treated as some concerns rather than high risk, although their
17 position was informed by relevant evidence.

18 We agree with Minozzi et al. [11] that review authors should develop guidance specific to their
19 review to establish how to assess issues such as those described in this article and assure
20 consistency across assessed results. Review authors can use our findings to plan the time it will take
21 to assess risk of bias, although this may vary substantially. Review users should be aware that RoB2
22 implementation will affect judgement and thus certainty assessments.

23 Our findings suggest refined guidance from the tool developers is warranted, with a focus on
24 operationalising the tool. We would welcome further specific examples directly related to the
25 signalling questions, and more examples of ‘judgement calls’ rather than extremes. For example, the
26 What Works Clearinghouse Standards Handbook provides boundaries for unacceptable risk of bias
27 for combinations of overall and differential attrition under cautious and optimistic assumptions [15],
28 the former being similar to our implementation for SQ 3.4. However, we recognise that it is
29 important to limit the overall volume of guidance. Sometimes the theoretical background included
30 factors that were not to be assessed; these and empirical evidence could be moved to an appendix.
31 Reporting guidelines and critical appraisal tools should also include details of risk-of-bias
32 implementation (e.g. PRISMA 2020 [16], AMSTAR2 [17]).

1 5 Conclusion

2 The RoB2 tool is a positive development from the original risk-of-bias tool, with the addition of
3 signalling questions and an algorithm likely to improve consistency if carefully followed. Assessing
4 individual results is useful for differentiating within-study risk of bias between outcomes and
5 timepoints, particularly in domains such as missing outcome data. This combined with an overall
6 risk-of-bias judgement assist in progressing to an assessment of the certainty of the evidence.
7 Nonetheless, conducting assessments with RoB2 is a substantial and challenging undertaking. We
8 recommend reporting guidelines and critical appraisal tools reflect this, requiring implementation
9 details. Furthermore, the burden on reviewers should be reduced by improving the guidance with
10 greater emphasis on the application of, and location of dividing lines for, each signalling question.

1 Acknowledgements

2 We thank the members of the wider systematic review and network meta-analysis project team who
3 attended Project Management Group meetings: Professor of biostatistics Richard Riley, applied
4 statistics lecturer Ram Bajpai, and medical statistics research assistant Matthew Bond, ageing
5 research assistant Eleftheria Patetsini, ageing research assistant Ridha Ramiz; or who otherwise co-
6 authored the project report on which this is based, including: information specialist Deirdre Andre,
7 aged care researcher Alison Ellwood, rehabilitation research programme manager John Green,
8 geriatric academic fellow Matthew Hale, geriatric medicine doctor Jessica Morgan, cardiovascular
9 ageing clinical lecturer Oliver Todd, and anaesthetics doctor Rebecca Walford. We are grateful to
10 Jasmin Manik for her support with testing the RoB2 tool algorithms.

11 Funding

12 This project is funded by the National Institute for Health Research (NIHR) Health Technology
13 Assessment programme (NIHR128862). The views expressed are those of the author(s) and not
14 necessarily those of the NIHR or the Department of Health and Social Care. The grant applicants
15 designed the overarching systematic review and network meta-analysis; the funders approved the
16 protocol. The funders have not been involved in any aspect of the collection, analysis or
17 interpretation of data; in the writing of the report; or in the decision to submit the article for
18 publication.

19 Data availability

20 An anonymised version of the dataset upon which this article is based will be made available upon
21 reasonable request to the author.

22 References

- 23 [1] Higgins JPT, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, *et al.* *Cochrane handbook for*
24 *systematic reviews of interventions, Second Edition*: John Wiley & Sons; 2019.
- 25 [2] Guyatt GH, Oxman AD, Vist G, Kunz R, Brozek J, Alonso-Coello P, *et al.* GRADE guidelines: 4.
26 Rating the quality of evidence—study limitations (risk of bias). *J Clin Epidemiol* 2011;**64**:407-
27 15. <https://doi.org/10.1016/j.jclinepi.2010.07.017>

- 1 [3] Higgins JP, Altman DG. Assessing Risk of Bias in Included Studies. In: Higgins JP, Green S,
2 editors. *Cochrane Handbook for Systematic Reviews of Interventions*; 2008:187-241.
3 <https://doi.org/10.1002/9780470712184.ch8>
- 4 [4] Higgins JPT, Altman DG, Gøtzsche PC, Jüni P, Moher D, Oxman AD, *et al.* The Cochrane
5 Collaboration's tool for assessing risk of bias in randomised trials. *BMJ* 2011;**343**:d5928.
6 <https://doi.org/10.1136/bmj.d5928>
- 7 [5] Jørgensen L, Paludan-Müller AS, Laursen DRT, Savović J, Boutron I, Sterne JAC, *et al.*
8 Evaluation of the Cochrane tool for assessing risk of bias in randomized clinical trials:
9 overview of published comments and analysis of user practice in Cochrane and non-
10 Cochrane reviews. *Systematic Reviews* 2016;**5**:80. [https://doi.org/10.1186/s13643-016-](https://doi.org/10.1186/s13643-016-0259-8)
11 [0259-8](https://doi.org/10.1186/s13643-016-0259-8)
- 12 [6] Higgins JPT, Savović J, Page MJ, Elbers RG, Sterne JAC. Chapter 8: Assessing risk of bias in a
13 randomized trial. In: Higgins J, Thomas J, Chandler J, Cumpston M, Li T, Page M, *et al.*,
14 editors. *Cochrane handbook for systematic reviews of interventions* Version 6.3 (updated
15 February 2022). Cochrane; 2022. URL: www.training.cochrane.org/handbook/ (Accessed 12
16 January 2023).
- 17 [7] Sterne JAC, Savović J, Page MJ, Elbers RG, Blencowe NS, Boutron I, *et al.* RoB 2: a revised tool
18 for assessing risk of bias in randomised trials. *BMJ* 2019;**366**:l4898.
19 <https://doi.org/10.1136/bmj.l4898>
- 20 [8] Higgins JPT, Savović J, Page MJ, Sterne JAC, RoB2 Development Group. *Revised Cochrane*
21 *risk-of-bias tool for randomized trials (RoB 2)*; 2019. URL: <https://www.riskofbias.info/>
22 (Accessed 12 January 2023).
- 23 [9] Minozzi S, Gonzalez-Lorenzo M, Cinquini M, Berardinelli D, C C, Ciardullo S, *et al.* Adherence
24 of systematic reviews to Cochrane RoB2 guidance was frequently poor: a meta
25 epidemiological study. *J Clin Epidemiol* 2022;**152**:47-55.
26 <https://doi.org/10.1016/j.jclinepi.2022.09.003>
- 27 [10] Minozzi S, Cinquini M, Gianola S, Gonzalez-Lorenzo M, Banzi R. The revised Cochrane risk of
28 bias tool for randomized trials (RoB 2) showed low interrater reliability and challenges in its
29 application. *J Clin Epidemiol* 2020;**126**:37-44. <https://doi.org/10.1016/j.jclinepi.2020.06.015>

- 1 [11] Minozzi S, Dwan K, Borrelli F, Filippini G. Reliability of the revised Cochrane risk-of-bias tool
2 for randomised trials (RoB2) improved with the use of implementation instruction. *J Clin*
3 *Epidemiol* 2022;**141**:99-105. <https://doi.org/10.1016/j.jclinepi.2021.09.021>
- 4 [12] Crocker TF, Clegg A, Riley RD, Lam N, Bajpai R, Jordão M, *et al.* Community-based complex
5 interventions to sustain independence in older people, stratified by frailty: a protocol for a
6 systematic review and network meta-analysis. *BMJ Open* 2021;**11**:e045637.
7 <https://doi.org/10.1136/bmjopen-2020-045637>
- 8 [13] Eldridge S, Campbell MK, Campbell MJ, Drahota AK, Giraudeau B, Reeves BC, *et al.* Revised
9 Cochrane risk of bias tool for randomized trials (RoB 2): Additional considerations for cluster-
10 randomized trials (RoB 2 CRT) (18 March 2021). 2021.
- 11 [14] Cochrane Training. *RoB 2: Learning Live webinar series*. URL:
12 <https://training.cochrane.org/rob-2-learning-live-webinar-series> (Accessed 12 January
13 2023).
- 14 [15] What Works Clearinghouse. *What Works Clearinghouse Procedures and Standards*
15 *Handbook* Version 5.0 (updated December 2022). Washington, DC: U.S. Department of
16 Education, Institute of Education Sciences, National Center for Education Evaluation; 2022.
17 URL: <https://ies.ed.gov/ncee/wwc/handbooks> (Accessed 30 May 2023).
- 18 [16] Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, *et al.* The PRISMA
19 2020 statement: an updated guideline for reporting systematic reviews. *BMJ* 2021;**372**:n71.
20 <https://doi.org/10.1136/bmj.n71>
- 21 [17] Shea BJ, Reeves BC, Wells G, Thuku M, Hamel C, Moran J, *et al.* AMSTAR 2: a critical appraisal
22 tool for systematic reviews that include randomised or non-randomised studies of
23 healthcare interventions, or both. *BMJ* 2017;**358**:j4008. <https://doi.org/10.1136/bmj.j4008>