

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.



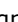
The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email editorial.bmjopen@bmj.com

BMJ Open

Disagreements in risk of bias assessment for randomised controlled trials included in more than one Cochrane systematic reviews

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-028382
Article Type:	Research
Date Submitted by the Author:	05-Dec-2018
Complete List of Authors:	Bertizzolo, Lorenzo; INSERM, U 1153, Equipe Methods Bossuyt, Patrick; Academic Medical Center; University of Amsterdam, Dept. Clinical Epidemiology and Biostatistics Atal, Ignacio; INSERM U1153, Team Methods Ravaud, Philippe; INSERM, U1153, Epidemiology and Biostatistics Sorbonne Paris Cite Research Center (CRESS), Methods of therapeutic evaluation of chronic diseases team (METHODS) Dechartres, Agnes; INSERM U738 H+  pital H+  telDieu Universit+  ParisDescartes, Centre dEpidemiologie Clinique
Keywords:	risk of bias, Cochrane, systematic reviews, interrater agreement, reproducibility

SCHOLARONE™
Manuscripts

1
2
3 **Disagreements in risk of bias assessment for randomised controlled trials included in**
4
5 **more than one Cochrane systematic reviews: a research on research study**
6
7
8
9

10
11
12 Lorenzo Bertizzolo¹, Patrick M Bossuyt², Ignacio Atal^{1, 5}, Philippe Ravaud^{1, 3-6}, Agnès
13 Dechartres^{1, 3-5, 7}
14
15

16
17
18
19 ¹ INSERM, U1153 Epidemiology and Biostatistics Sorbonne Paris Cité Research Center
20 (CRESS), Methods of therapeutic evaluation of chronic diseases Team (METHODS), Paris,
21 F-75004 France; Paris Descartes University, Sorbonne Paris Cité, France.
22
23

24
25
26 ² Department of Clinical Epidemiology, Biostatistics and Bioinformatics, Academic Medical
27 Center, University of Amsterdam, Netherlands.
28
29

30
31 ³ Centre d'Épidémiologie Clinique, Hôpital Hôtel Dieu, AP-HP (Assistance Publique des
32 Hôpitaux de Paris), Paris, France.
33
34

35 ⁴ Faculté de Médecine, Université Paris Descartes, Sorbonne Paris Cité, Paris, France.
36

37 ⁵ Cochrane France, Paris, France
38

39
40 ⁶ Columbia University, Mailman School of Public Health, Department of Epidemiology, New
41 York, USA
42
43

44
45 ⁷ Sorbonne Université, INSERM, Institut Pierre Louis de Santé Publique, Département
46 Biostatistique, Santé Publique et Information Médicale, AP-HP, Hôpitaux Universitaires Pitié
47 Salpêtrière – Charles Foix, Paris, France
48
49

50
51
52
53 Correspondence to: Lorenzo Bertizzolo, M.D., MPH

54
55
56 INSERM, U1153 - Centre d'Épidémiologie Clinique –
57
58

Hôpital Hôtel-Dieu,

1, place du parvis Notre Dame, 75004 Paris, France

Tel: +33 (0)1 42 34 78 25

E-mail: lorenzo.bertizzolo@gmail.com

Key-words: risk of bias, Cochrane, Systematic Reviews, Interrater Agreement,

Reproducibility

Word Count:

Copyright: The Corresponding Author has the right to grant on behalf of all authors and does grant on behalf of all authors, a worldwide licence to the Publishers and its licensees in perpetuity, in all forms, formats and media (whether known now or created in the future), to i) publish, reproduce, distribute, display and store the Contribution, ii) translate the Contribution into other languages, create adaptations, reprints, include within collections and create summaries, extracts and/or, abstracts of the Contribution, iii) create any other derivative work(s) based on the Contribution, iv) to exploit all subsidiary rights in the Contribution, v) the inclusion of electronic links from the Contribution to third party material where-ever it may be located; and, vi) licence any third party to do any or all of the above.

Transparency declaration: The guarantor (Lorenzo Bertizzolo) affirms that the manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

Ethics approval: Not applicable. This is a research on research study.

ABSTRACT (299 words)

Objectives: Assess the frequency and reasons for disagreements in risk of bias assessments for randomised clinical trials (RCTs) included in more than one Cochrane review.

Design: Research on research study, using cross-sectional design.

Data sources: 2,796 Cochrane reviews published between March 2011 and September 2014.

Data selection: RCTs included in more than one review.

Data extraction: Risk of bias assessment and support for judgement for five key risk of bias items.

Data synthesis: For each item, we compared risk of bias assessment made in each review and calculated proportion of agreement. Two reviewers independently analysed 50% of all disagreements by comparing support for each judgement with information from study report to evaluate whether disagreements were related to a difference in information (e.g., contact the study author) or a difference in interpretation (same support for judgement but different interpretation). They also identified main reasons for different interpretation.

Results: 1,604 RCTs were included in more than one review. Proportion of agreement ranged from 57% (770/1,348 trials) for incomplete outcome data to 81% for random sequence generation (1,193/1,466). Most common source of disagreement was difference in interpretation of the same information, ranging from 65% (88/136) for random sequence generation to 90% (56/62) for blinding of participants and personnel. Access to different information explained 32/136 (24%) disagreements for random sequence generation and 38/205 (19%) for allocation concealment. Disagreements related to difference in interpretation were frequently related to incomplete or unclear reporting in the study report (83% of disagreements related to different interpretation for random sequence generation).

1
2
3 **Conclusions:** Risk of bias judgements of RCTs included in more than one Cochrane review
4 differed substantially. Most disagreements were related to a difference in interpretation of an
5 incomplete or unclear description in the study report. A clearer guidance on common causes
6 of incomplete information may improve agreement.
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Article summary

Strengths and limitations of this study

- Use of a very large and comprehensive collection of Cochrane reviews to assess the agreement in risk of bias assessment and to understand reasons of disagreement.
- Analysis of the full-text of study reports to underline what information were available to review authors and how they utilized them while assessing risk of bias.
- Focus on disagreements only. Possible that a proportion of agreements happened “by chance”. For example review authors may express the same risk of bias judgement while using different information or interpreting information differently.
- No evaluation of the potential impact of disagreements in conclusion making at the review level.

INTRODUCTION

Systematic reviews aim to synthesise all existing evidence for a research question by the use of a rigorous and reproducible methodology¹. Because reviews may be affected by bias at the level of individual studies², an assessment of the risk of bias in these studies is a crucial step in conducting a systematic review^{3 4}.

Cochrane has developed a tool to provide a standardised approach to the assessment of the risk of bias in randomised clinical trials (RCTs)⁵. The risk of bias tool is based on specific characteristics related to study design and conduct, selected on theoretical grounds and on empirical evidence from meta-epidemiological studies that these characteristics are associated with differences in treatment effect estimates⁶⁻¹¹. The tool includes seven items (random sequence generation, allocation concealment, blinding of participants and personnel, blinding of outcome assessment, incomplete outcome data, selective reporting, other source of bias), the researchers assess and judge as either “high”, “low” or “unclear” risk of bias^{11 12}. Although Cochrane provides detailed guidance on how to use the tool and recommends consensus between two independent reviewers¹¹, personal judgement is also involved, which may bring variability. Several studies have evaluated the reproducibility of the risk of bias tool, generally shown to be poor¹²⁻¹⁹. However, there is some uncertainty about the main causes of disagreements. For example, some reviewers may search for additional information such as protocols or contact study authors and this difference in available information, rather than a difference in judgement, may explain some of the disagreements.

In this study, we used a large collection of Cochrane reviews to evaluate the reproducibility of risk of bias assessments by identifying randomised controlled trials included in more than one Cochrane review and comparing the assessments. In addition, we examined the likely

1
2
3 reasons for any disagreements. In particular, we evaluated whether disagreements were
4
5 related to differences in information available to reviewers or differences in interpreting the
6
7 same information and what could explain such different interpretation.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

METHODS

This is a research on research study on risk of bias assessment, which used a cross sectional design. We identified RCTs included in more than one reviews included in a large collection of Cochrane reviews. For key risk of bias items, we evaluated agreement between the different systematic reviews; analysed whether disagreements were related to a difference in information available to reviewers or a difference in interpretation of the same information and highlighted the main reasons for disagreements by an in-depth, one-by-one evaluation of disagreements.

Data sources

We obtained data from the 2,796 Cochrane reviews, which corresponds to all reviews available in the Cochrane library between March 2011 and September 2014, including updates (March 2011 corresponds to the last update of the risk of bias tool⁵). Data consisted of one XML file per review, each file containing all data entered by review authors in RevMan, the software used for managing Cochrane reviews²⁰. All individual XML files were merged in a single database by using R v3.2.2²¹ with the XML package²². The vocabulary used for risk of bias items slightly varied across reviews (e.g., some reviews could refer to “allocation concealment” as “allocation masking”). For this reason, two authors independently evaluated all terms used and classified them according to the vocabulary of the tool. Disagreements were resolved by consensus. This standardization was done for a previous publication²³.

Selection of eligible reviews

1
2
3 We excluded withdrawn or “empty” reviews (i.e., systematic reviews not including any
4 study) as well as reviews including observational or non-randomised studies and considered
5 only reviews with an assessment of risk of bias for at least one item of the risk of bias tool.
6
7
8
9

10 11 12 ***Selection of eligible RCTs*** 13

14 To identify single RCTs included and assessed for risk of bias in more than one systematic
15 review, we proceeded as follows. For each RCT, we identified the primary reference(s),
16 which was the reference identified by review authors as the main reference(s) for an included
17 study. Then, we used a matching algorithm²⁴ to identify studies that shared the same primary
18 reference. If several primary references were reported, we considered all of them. We
19 manually checked that the studies sharing the same primary reference in the reviews
20 corresponded to the same RCT.
21
22
23
24
25
26
27
28
29
30

31 32 33 ***Extraction of risk of bias assessment*** 34

35 For each eligible RCT, we extracted the risk of bias assessment and the corresponding
36 support for judgement for each risk of bias item in each review. Whenever a single RCT was
37 included in three or more reviews, we considered only the risk of bias assessment from two
38 reviews chosen at random; this situation concerned less than 10% of our included RCTs and
39 was decided because of workload and to facilitate direct comparison of two assessments. We
40 focused on five risk of bias items: random sequence generation, allocation concealment,
41 blinding of participants and personnel, blinding of outcome assessment and incomplete
42 outcome data. We did not consider selective reporting because it is difficult to evaluate in the
43 absence of the study protocol, which is frequently lacking, especially for older studies^{11 12 14}.
44
45 We also did not consider the item other bias because the definition is very wide (i.e., “any
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 important concerns about bias not covered in the other domains in the tool”¹¹), so
4
5 comparisons across reviews are difficult.
6
7
8
9

10 ***Comparison of risk of bias assessment between reviews***

11
12 For each item, we compared the risk of bias assessment in terms of “high”, “low” or
13
14 “unclear” risk of bias between the two reviews. According to the Cochrane handbook, the
15
16 items blinding of outcome assessment and incomplete outcome data should be assessed for
17
18 each outcome. Therefore, when the reviews reported an assessment of these items at the
19
20 outcome level, we manually checked that outcomes were identical in both reviews and we
21
22 retained for our analysis only the assessments that focused on the same outcomes. For
23
24 blinding, we followed the last version of the Cochrane handbook and we retained only
25
26 assessments of blinding of participants and personnel and blinding of outcome assessment as
27
28 two independent items, excluding different types of assessment (i.e., blinding as a single
29
30 item, blinding of only participants or of only personnel).
31
32
33
34

35 We calculated the percentage agreement for each risk of bias item, as the proportion of
36
37 studies with a concordant assessment in both reviews (e.g. “low” risk of bias AND “low” risk
38
39 of bias). Not all reviews assessed all five key risk of bias items for each RCT included;
40
41 consequently, the number of RCTs evaluated for discrepancies varies depending on the item
42
43 considered.
44
45
46
47
48

49 ***Selection of studies for in-depth analysis of disagreements***

50
51 For workload reasons, we in-depth evaluated the reasons for disagreements for 50% of the
52
53 studies analysed in the previous step. In cases of more than one shared RCTs within a given
54
55
56
57
58
59
60

1
2
3 pair of Cochrane reviews, we selected only one RCT at random. To reach 50% of the total
4
5 sample, we used a simple random selection in the remaining database.
6
7
8
9

10 *Classification of disagreements*

11
12 For the random selection, two reviewers (LB and AD) independently evaluated all
13
14 disagreements in the risk of bias assessment in the two systematic reviews. They first
15
16 scrutinised the support for the judgement in each review and evaluated whether it was the
17
18 same or “conceptually” the same in both reviews (e.g., “randomised, probably done”;
19
20 “randomised, probably not done”; “study only mention randomization, but does not specify
21
22 how randomization was performed; unclear”; “study states it is randomized; low risk”). If the
23
24 support differed, they assessed any other information regarding the study as reported in both
25
26 reviews, systematically searching and evaluating the full-text study report indicated in the
27
28 primary reference. A formalized data extraction process for full texts was not used. Full-texts were
29
30 examined, looking primarily for correspondence between information reported by the reviewers in
31
32 their Support for Judgement and the text.
33
34
35
36

37 They independently classified each case of disagreement as follows:

- 38 • Disagreement related to differences in interpretation:
 - 39 ○ The support for judgement was the same (or “conceptually” the same) in both
40
41 reviews, but the interpretation differed.
 - 42 ○ One review clearly confused one item of the risk of bias tool with a different
43
44 one or the review authors misunderstood the definition of the item (e.g., for
45
46 random sequence generation, support for judgement reports “600 opaque
47
48 envelopes, 1 was drawn every time”).
49
50
51
52
53
54
55
56
57
58
59
60

- Disagreement related to differences in information: the support for the judgement cites information that is not available in the study report; additional sources are cited (e.g., protocol) or the review authors reported that they had contacted the RCT author for additional data.
- Disagreement related to information missed by the review authors: the study report clearly describes the information, but some review authors seemed to have missed this information in the study report.
- Disagreement related to input mistakes: risk of bias assessment in terms of “high”/“low”/“unclear” did not match the support for the judgement (e.g., “Randomization described explicitly”, judgement “Unclear”).
- Unclear: when it was not possible to classify the disagreement because the support for the judgement was empty or because we could not retrieve the full-text study report.

Any disagreements between reviewers were solved by discussion to reach consensus. In the Supplementary Appendix 1, we report a figure synthesizing how the in-depth analysis process was conducted.

Identification of main reasons for different interpretation

For each disagreement related to a difference in interpretation, we evaluated the probable reason for disagreement. For example, the interpretation could differ because of confusion with another risk of bias item (e.g., random sequence generation and allocation concealment) or because the information was unclear or insufficiently detailed in the article. When we were unsure about the reason, we classified the reason as unclear. Two authors (LB and AD) conducted this process in duplicate by using all available information (i.e., support for the

1
2
3 judgement, characteristics of the study reported in the review, full-text article), with
4
5 disagreements resolved by discussion.
6
7
8
9

10 **Statistical analysis**

11
12 Analysis was descriptive with use of frequencies and percentages for qualitative variables.

13
14 Statistical analysis was conducted with Stata 13.1²⁵. We decided to use simple percent
15
16 agreement because other static approaches were problematic. The Kappa statistic requires
17
18 having defined reviewers, which is not the case of our approach. Another statistic, the
19
20 intraclass correlation coefficient (ICC) is not suitable, because it requires assessments to be in
21
22 an ordinal order, which is not our case. There is no continuum between the assessments of
23
24 low, unclear and high risk of bias.
25
26
27
28
29

30 **Patient involvement**

31
32
33 Patients were not involved in any aspect of the study design, conduct, or the development of
34
35 the research question or outcome measures. This is a research-on-research study and
36
37 therefore there was no active patient recruitment for data collection.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

RESULTS

Selection process

Figure 1 shows the selection process. From the 2,796 systematic reviews published between March 2011 and September 2014, 2,291 reviews included RCTs only and reported a risk of bias assessment. Of these, 797 included at least one RCT whose primary reference was shared with another review for which a risk of bias assessment was reported. These 797 reviews included 1,604 single RCTs evaluated for the same risk of bias item in more than one review. The Supplementary Appendix 2 reports the frequency of the different Cochrane groups among those reviews.

Among the 1,604 selected RCTs: 1,603 had duplicate evaluation for allocation concealment, 1,466 for random sequence generation, 375 for blinding of participants and personnel, 583 for blinding of outcome assessment and 1,348 for incomplete outcome data.

Evaluation of agreement and distribution of disagreements

The agreement of risk of bias judgements ranged from 57% (770/1,348 trials) for incomplete outcome data to 81% (1,193/1,466 trials) for random sequence generation (Figure 2). We identified most disagreements for “low” and “unclear” risk of bias judgments, especially for random sequence generation (231/273 trials, 85%). Disagreements between “low” and “high” risk of bias were generally rare, for example 8/273 of disagreements (3%) for random sequence generation, with the exception of incomplete outcome data for which they were more frequent (190/578, 33%). For blinding of participants and personnel, the most frequent disagreement was between “unclear” and “high” risk of bias (50/107, 47%), then “low” versus “unclear” (34/107, 32%), and “low” versus “high” (23/107, 21%) (Figure 2).

Classification of disagreements

The in-depth analysis of disagreements included 802 studies: 799 for allocation concealment, 747 for random sequence generation, 206 for blinding of participants and personnel, 297 for blinding of outcome assessment and 660 for incomplete outcome data. The agreement results of this sample and the distribution of disagreements are reported in the Supplementary Appendix 3.

For all items, the most common source of disagreement was a difference in interpretation, with frequencies ranging from 88/136 (65%) for random sequence generation to 56/62 (90%) for blinding of participants and personnel (Figure 3). The access to additional or different information accounted for disagreements in 32/136 (24%) trials for random sequence generation and 38/205 (19%) for allocation concealment. Access to additional information was less common for the remaining items, with proportions ranging from 2% to 4%. In 80% of the cases, the access to additional information was through the contact of the study author. The other sources of disagreement were less common; input mistake ranged from 1% to 6%, missed information from 1% to 6%. We could not determine the source of disagreement in 5% of our disagreements. For this analysis, we accessed the full text of 216 different trials to help us in the process. The Supplementary Appendix 4 reports some examples of disagreements in which the access to the study report helped us in the classification and the analysis of reasons of disagreement. We could not retrieve or access 19 full-texts we deemed necessary for the categorization of disagreements and this explain the majority of cases where we were unable to categorize the source of disagreement (“unclear” source in Figure 3).

Main reasons of disagreements for different interpretation

1
2
3 The main reasons for a difference in interpretation for each item are reported in Table 1.
4
5 Additional examples are provided for each item for the high-low disagreements
6
7 (Supplementary Appendix 5). The most common reason across items was related to
8
9 incomplete or unclear reporting in the RCT. For random sequence generation, disagreements
10
11 in 73/88 (83%) trials were related to lack of a precise description of the randomization
12
13 process with reviewers evaluating “low”, “high” or “unclear” risk of bias the reporting of
14
15 “randomised” in the text. For allocation concealment, the most common reason for
16
17 disagreement was a different interpretation of description of the envelopes used to conceal
18
19 allocation (17%, n=26/149 trials). For the two blinding items, many disagreements occurred
20
21 when the article mentioned only “double blind” in RCTs without an additional description
22
23 (16% of cases, n=9/56 trials for blinding of participants and personnel, 13%, n=9/70 for
24
25 blinding of outcome assessment). For incomplete outcome data, reviewers assessed
26
27 differently the statement from the study report of “no missing data” or “all data reported”
28
29 (10%, 22/220 trials). Another common reason for a difference in interpretation was confusion
30
31 with another item. Allocation concealment was confused with blinding (10%, n=15/149
32
33 trials) but also with random sequence generation (4%, n=6/149). For blinding of participants
34
35 and personnel, the most common cause for disagreement concerned the interpretation of
36
37 cases when blinding was not feasible (36%, n=20/56 trials), assessed at high risk by some
38
39 reviewers and low by others. Another common cause of disagreement for the two blinding
40
41 items related to the assessment of outcomes that should not be affected by blinding (e.g.,
42
43 mortality); it explained 21% (n=12 trials) of disagreements for blinding of participants and
44
45 personnel and 23% (n=16 trials) for blinding of outcome assessment, often low versus high
46
47 disagreements.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 For incomplete outcome data, the use of different cut-offs for the rate of missing data is the
4 most common reason for disagreement (26%, n=57 trials); also common is considering the
5 explanation of reasons for missing data enough to attribute a low risk of bias (13%. n=28
6 trials).
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

DISCUSSION

In this study, we took advantage of a very large sample of Cochrane reviews to explore the sources of disagreements in risk of bias assessment for trials included in several reviews. Our results confirm that the agreement for risk of bias assessments is generally suboptimal, with better agreement for random sequence generation and allocation concealment and less agreement for incomplete outcome data. Access to different sources of information explained why 24% of the trials had disagreements in the assessment of risk of bias for random sequence generation and 19% for allocation concealment. However, the main source of disagreements was a difference in interpretation of the same information, which was frequently related to incomplete or unclear reporting in the study report.

Strengths and weaknesses

Our study goes beyond previous literature on the topic^{3 12-18 26}. As compared with most other studies¹²⁻¹⁷ we used real-world data to explore agreement of risk of bias assessments in real scenarios. We evaluated a very large and comprehensive collection of Cochrane reviews that spanned multiple specialties and topics, including a number of trials about ten times larger than the largest study on the topic¹². We completed our analysis by searching individual study reports to give support to our comments on reasons for disagreements, which, to our knowledge, has not been done in previous, smaller works that used a similar methodology¹⁸. While doing this, we developed a suitable classification scheme for sources of disagreements and conducted, in duplicate, an extensive analysis to understand the risk of bias assessment process and explored the most common reasons for disagreements.

1
2
3 Our study has limitations. Although the classification of disagreements was conducted in
4 duplicate following a formalised process, there remains a component of personal judgement.
5
6 We evaluated only disagreements, but a number of agreements might have occurred “by
7
8 chance”. In our analysis of likely reasons for disagreements, some resulted from confusion
9
10 between risk of bias items. Similar discrepancies might have occurred among agreements;
11
12 indeed, previous literature on the topic demonstrated that reviewers do not accurately follow
13
14 the risk of bias tool²⁷. We did not assess non-Cochrane reviews, even if they often use the
15
16 Cochrane risk of bias tool. Agreement in these reviews is likely worse because of less
17
18 familiarity and training with the tool. We also did not assess the selective reporting item that
19
20 is frequently judged on incomplete information. We did not evaluate whether disagreements
21
22 varied depending on the Cochrane review group or year of publication. Finally, we did not
23
24 evaluate the impact of disagreements and the extent to which the evidence base for making
25
26 conclusions and providing summary statements of effectiveness may have been affected by
27
28 changing the rating.
29
30
31
32
33
34
35
36
37

38 **Comparison with other studies**

39
40 Our findings confirm the importance of issues that were previously identified by Jorgensen et
41
42 al.³ and Savovic et al.²⁶. In particular, Savovic et al.²⁶, surveying users of the risk of bias tool,
43
44 reported on the possibility of confusion between random sequence generation and allocation
45
46 concealment and between allocation concealment and blinding; the uncertainty on how to
47
48 address unfeasibility of blinding; and the difficulties in assessing incomplete outcome data
49
50 especially regarding the acceptable rate of missing data. More recently, Jorgensen et al.³,
51
52 evaluating comments on the use of the risk of bias tool, highlighted how authors complained
53
54 that judgment often originates from incomplete or missing information.
55
56
57
58
59
60

1
2
3 A previous study identified 46 RCTs included in different systematic reviews in the field of
4 fertility and evaluated the percentage agreement in risk of bias assessment. That analysis
5 showed generally worse agreement than in our study, with percentage agreement ranging
6 from 35% to 71%. Differences in sample size and the particular topic may explain these
7 differences. In addition, although the authors had compared supports for judgement between
8 reviews, this evaluation may have been incomplete, because they did not evaluate the primary
9 study reports¹⁸.

20 21 **Implications**

22
23 Our results confirm that the agreement in risk of bias assessment would be enhanced by more
24 detailed guidance in use of the risk of bias tool with particular focus on common causes of
25 disagreements. We showed that in many cases, the unclear reporting from source material
26 allows reviewers ample space for personal judgement and differences in judgement.
27

28
29 The scientific community continues to stress the importance of improving the reporting of
30 trials²⁸⁻³¹, which may limit disagreements when assessing risk of bias. In parallel, we could
31 also work on restricting the space for personal interpretation when assessing risk of bias. A
32 suggestion could be to give clearer instruction on how to evaluate common cases, for
33 example when confronted with nothing more than the term “randomised” or “double blind”
34 in the study report. Similarly, a threshold could be set on the quota for missing data and
35 indications on which imputation methods are appropriate and in which situations.
36
37

38
39 To minimise research waste, it could be interesting to have access to risk of bias assessments
40 from other Cochrane groups and the supports they used, including information from authors
41 or from protocols to help reviewers in their assessments. This process would imply having a
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 unique study identification number across reviews and a central shared repository for all
4
5 studies included in any Cochrane reviews.
6

7
8 Following the suggestions based on the findings and comments provided by Jorgensen et al.³
9
10 and Savovic et al.²⁶, Cochrane has been working on a new version of the risk of bias tool,
11
12 which has recently been released^{32 33}. The new version has a different approach to the risk of
13
14 bias assessment, guiding reviewers through the process with the use of “signalling questions”,
15
16 which might leave less room for subjectivity. In addition, there is more guidance in assessing
17
18 some items. For example, the new tool better clarifies some aspects of the randomization
19
20 process, especially about what to do in some cases of incomplete information (e.g.,
21
22 randomization list created by an external centre with no other indication). The new tool also
23
24 has a different approach to the blinding aspect, oriented to the implications of the masking
25
26 process. However, the new tool does not cover some of our concerns, especially those related
27
28 to incomplete outcome data: quota for missing data that are considered acceptable, and
29
30 whether reviewers should focus more on the reasons for the missing data or their magnitude.
31
32 It also does not address the common case of authors reporting “no missing data”. Research-
33
34 on-research studies are needed to evaluate whether this new version of the tool results in
35
36 improved reproducibility.
37
38
39
40
41
42
43
44

45 **Conclusion**

46
47 This analysis of risk of bias assessment for more than 1,600 trials included in more than one
48
49 reviews showed that agreement remains suboptimal. Most disagreements come from a
50
51 difference in interpretation of an incomplete or unclear description in the study report. In
52
53 some cases, the difference in the assessment was due to some but not all review authors
54
55 obtaining additional information, from a protocol or from contacting study author.
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

REFERENCES

1. Cook DJ, Mulrow CD, Haynes RB. Systematic reviews: synthesis of best evidence for clinical decisions. *Ann Intern Med* 1997;126(5):376-80. [published Online First: 1997/03/01]
2. Hopewell S, Boutron I, Altman DG, et al. Incorporation of assessments of risk of bias of primary studies in systematic reviews of randomised trials: a cross-sectional study. *BMJ Open* 2013;3(8):e003342. doi: 10.1136/bmjopen-2013-003342 [published Online First: 2013/08/27]
3. Jorgensen L, Paludan-Muller AS, Laursen DR, et al. Evaluation of the Cochrane tool for assessing risk of bias in randomized clinical trials: overview of published comments and analysis of user practice in Cochrane and non-Cochrane reviews. *Syst Rev* 2016;5:80. doi: 10.1186/s13643-016-0259-8 [published Online First: 2016/05/11]
4. Hrobjartsson A, Boutron I, Turner L, et al. Assessing risk of bias in randomised clinical trials included in Cochrane Reviews: the why is easy, the how is a challenge. *Cochrane Database Syst Rev* 2013(4):ED000058. doi: 10.1002/14651858.ED000058 [published Online First: 2013/06/04]
5. Higgins JP, Altman DG, Gotzsche PC, et al. The Cochrane Collaboration's tool for assessing risk of bias in randomised trials. *BMJ* 2011;343(oct18 2):d5928. doi: 10.1136/bmj.d5928
6. Schulz KF, Chalmers I, Hayes RJ, et al. Empirical evidence of bias. Dimensions of methodological quality associated with estimates of treatment effects in controlled trials. *JAMA* 1995;273(5):408-12. [published Online First: 1995/02/01]
7. Moher D, Pham B, Jones A, et al. Does quality of reports of randomised trials affect estimates of intervention efficacy reported in meta-analyses? *The Lancet* 1998;352(9128):609-13. doi: 10.1016/s0140-6736(98)01085-x
8. Balk EM, Bonis PA, Moskowitz H, et al. Correlation of quality measures with estimates of treatment effect in meta-analyses of randomized controlled trials. *JAMA* 2002;287(22):2973-82. [published Online First: 2002/06/08]
9. Page MJ, Higgins JP, Clayton G, et al. Empirical Evidence of Study Design Biases in Randomized Trials: Systematic Review of Meta-Epidemiological Studies. *PLoS One* 2016;11(7):e0159267. doi: 10.1371/journal.pone.0159267 [published Online First: 2016/07/12]
10. Dechartres A, Trinquart L, Faber T, et al. Empirical evaluation of which trial characteristics are associated with treatment effect estimates. *J Clin Epidemiol* 2016;77:24-37. doi: 10.1016/j.jclinepi.2016.04.005
11. Higgins JPT, Altman DG, Sterne JAC (editors). Chapter 8: Assessing risk of bias in included studies. In: Higgins JP, Green S, eds. *Cochrane Handbook for Systematic Reviews of Interventions*. Version 5.1.0 (updated March 2011) ed: The Cochrane Collaboration, 2011. Available from www.handbook.cochrane.org.
12. Hartling L, Ospina M, Liang Y, et al. Risk of bias versus quality assessment of randomised controlled trials: cross sectional study. *BMJ* 2009;339(oct19 1):b4012. doi: 10.1136/bmj.b4012
13. Hartling L, Bond K, Vandermeer B, et al. Applying the risk of bias tool in a systematic review of combination long-acting beta-agonists and inhaled corticosteroids for persistent asthma. *PLoS One* 2011;6(2):e17242. doi: 10.1371/journal.pone.0017242
14. Hartling L, Hamm MP, Milne A, et al. Testing the risk of bias tool showed low reliability between individual reviewers and across consensus assessments of reviewer pairs. *J Clin Epidemiol* 2013;66(9):973-81. doi: 10.1016/j.jclinepi.2012.07.005
15. Armijo-Olivo S, Ospina M, da Costa BR, et al. Poor reliability between Cochrane reviewers and blinded external reviewers when applying the Cochrane risk of bias tool in physical therapy trials. *PLoS One* 2014;9(5):e96920. doi: 10.1371/journal.pone.0096920 [published Online First: 2014/05/16]
16. Armijo-Olivo S, Stiles CR, Hagen NA, et al. Assessment of study quality for systematic reviews: a comparison of the Cochrane Collaboration Risk of Bias Tool and the Effective Public Health Practice Project Quality Assessment Tool: methodological research. *J Eval Clin Pract* 2012;18(1):12-8. doi: 10.1111/j.1365-2753.2010.01516.x
17. da Costa BR, Beckett B, Diaz A, et al. Effect of standardized training on the reliability of the Cochrane risk of bias assessment tool: a prospective study. *Syst Rev* 2017;6(1):44. doi: 10.1186/s13643-017-0441-7

18. Jordan VM, Lensen SF, Farquhar CM. There were large discrepancies in risk of bias tool judgments when a randomized controlled trial appeared in more than one systematic review. *J Clin Epidemiol* 2017;81:72-76. doi: 10.1016/j.jclinepi.2016.08.012 [published Online First: 2016/09/14]
19. Wilkins AJ. Risk of bias in assessing Risk of Bias. *Ophthalmic Physiol Opt* 2017;37(1):107-09. doi: 10.1111/opo.12333
20. Review Manager (RevMan) [Computer program] [program]. Version 5.3 version. Copenhagen: The Nordic Cochrane Centre, The Cochrane Collaboration, 2014.
21. R: A language and environment for statistical computing. [program]. Viena: R Foundation for Statistical Computing, 2013.
22. XML: Tools for Parsing and Generating XML Within R and S-Plus. [program]. 3.2.2 version, 2017.
23. Dechartres A, Trinquart L, Atal I, et al. Evolution of poor reporting and inadequate methods over time in 20 920 randomised controlled trials included in Cochrane reviews: research on research study. *BMJ* 2017;357:j2490. doi: 10.1136/bmj.j2490
24. deMelo, VV, Conference: Advances in Logic Based Intelligent Systems; 2005.
25. Stata Statistical Software: Release 13 [program]. College Station, TX: StataCorp LP, 2013.
26. Savovic J, Weeks L, Sterne JA, et al. Evaluation of the Cochrane Collaboration's tool for assessing the risk of bias in randomized trials: focus groups, online survey, proposed recommendations and their implementation. *Syst Rev* 2014;3:37. doi: 10.1186/2046-4053-3-37 [published Online First: 2014/04/16]
27. Propadalo I, Tranfic M, Vuka I, et al. In Cochrane reviews risk of bias assessments for allocation concealment were frequently not in line with Cochrane's Handbook guidance. *J Clin Epidemiol* 2018 doi: 10.1016/j.jclinepi.2018.10.002
28. Shamseer L, Hopewell S, Altman DG, et al. Update on the endorsement of CONSORT by high impact factor journals: a survey of journal "Instructions to Authors" in 2014. *Trials* 2016;17(1):301. doi: 10.1186/s13063-016-1408-z
29. Turner L, Shamseer L, Altman DG, et al. Does use of the CONSORT Statement impact the completeness of reporting of randomised controlled trials published in medical journals? A Cochrane review. *Syst Rev* 2012;1:60. doi: 10.1186/2046-4053-1-60
30. Turner L, Shamseer L, Altman DG, et al. Consolidated standards of reporting trials (CONSORT) and the completeness of reporting of randomised controlled trials (RCTs) published in medical journals. *Cochrane Database Syst Rev* 2012;11:MR000030. doi: 10.1002/14651858.MR000030.pub2
31. Altman DG, Moher D, Schulz KF. Improving the reporting of randomised trials: the CONSORT Statement and beyond. *Stat Med* 2012;31(25):2985-97. doi: 10.1002/sim.5402
32. RoB 2.0 Tool [Available from: <https://sites.google.com/site/riskofbiastool/welcome/rob-2-0-tool> accessed 22/08/2018.
33. Higgins JP, Sterne JA, Savovic J, et al. A revised tool for assessing risk of bias in randomized trials. *Cochrane Database of Systematic Reviews* 2016(10 (Suppl 1)) doi: dx.doi.org/10.1002/14651858.CD201601

ACKNOWLEDGEMENTS:

We thank Camila Olarte Parra for her help during the data management phase and her comments on this manuscript. We thank David Tovey, editor in chief of the Cochrane Library, for agreeing to share data from Cochrane reviews; Javier Mayoral Campos, system administrator; the Cochrane Central Executive for preparing files; and all Cochrane reviewers who collected data. We also thank Laura Smales for English revision of the manuscript.

CONTRIBUTIONS:

Lorenzo Bertizzolo was involved in the study conception, selection of trials, data extraction, data analysis, interpretation of results and drafting the manuscript.

Patrick Bossuyt was involved in the study conception, data analysis, interpretation of results and drafting the manuscript.

Ignacio Atal was involved in the study conception, data extraction, and drafting the manuscript.

Philippe Ravaud was involved in the study conception and drafting the manuscript.

Agnès Dechartres was involved in the study conception, selection of trials, data extraction, data analysis, interpretation of results and drafting the manuscript.

Lorenzo Bertizzolo is the guarantor. He had full access to all the data in the study and takes responsibility for the integrity of the data and the accuracy of the data analysis.

FUNDING:

This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 676207.

COMPETING INTERESTS:

1 The authors declare that they have no competing interests in relation to this study.
2
3
4

5 **DATA SHARING:**
6

7 Raw data and analyses are available on request from the authors.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

List of tables

Table 1: Main reasons for disagreements in cases of a different interpretation of the same information.

List of figures

Figure 1: Flow-chart of the selection process.

Figure 2: Distribution of agreements and disagreements for the different risk of bias items analysed; raw number and percentages of the total. For disagreements, distribution of the different discrepancies.

Figure 3: Classification of disagreements for the different items; raw number and percentage of the total.

List of appendices

Supplementary Appendix 1: explanatory figure of the categorization process for the in-depth analysis of disagreement.

Supplementary Appendix 2: Frequency of the different Cochrane review groups involved in the included reviews.

Supplementary Appendix 3: for the 50% selection of studies for the in-depth analysis; distribution of agreements and disagreements for the different risk of bias items analysed; raw number and percentages of the total. For disagreements, distribution of the different discrepancies.

Supplementary Appendix 4: selected examples where the access to the study report helped us in the categorization of the disagreement and in highlighting the reasons for disagreement.

Supplementary Appendix 5 Reasons for disagreements in cases of different interpretation of the same information; focus on “low” versus “high” disagreements

Table 1| Main reasons for disagreements in cases of a different interpretation of the same information.

<i>Risk of bias</i>	<i>Main reasons for disagreements</i>	<i>N (%)⁺</i>	<i>Examples of support for judgement from the review*</i>
<i>Item</i>			
<i>random sequence generation</i>	Consider differently incomplete or unclear description	73 (83)	“States “cluster randomisation by computer””; Low risk of bias “Cluster randomisation by computer. No further information provided”; Unclear risk of bias
	Confusion with allocation concealment	9 (10)	“allocation was done using sealed envelopes containing name of one of the two groups.”; Low risk of bias
<i>allocation concealment</i>	Consider differently incomplete or unclear description	49 (33)	“Not specified.”; High risk of bias “Method of concealment not described.”; Unclear risk of bias
	Consider differently envelopes description	26 (17)	““Sequentially numbered sealed envelopes”. Does not state if opaque envelopes.”; Unclear risk of bias “Sequentially numbered sealed envelopes.”; Low risk of bias
	Random sequence generated by computer or external centre considered enough for Low risk	21 (14)	“Treatment was allocated based on the computer-generated number list.”; Low risk of bias
	Confusion in the definition of the item	19 (13)	“Researchers attempted to contact all patients seen by physicians during one month”; High risk of bias
	Confusion with blinding	15 (10)	“participants were told to which compound they had been allocated.”; High risk of bias
	Confusion with random sequence generation	6 (4)	“Computer generated randomised lists.”; Low risk of bias

⁺ Number of RCTs disagreeing for this reason; percentage over the total of disagreements for different interpretation.

^{*} When two extracts are reported, they refer to the same study.

Peer review only

<i>Risk of bias Item</i>	<i>Main reasons for disagreements</i>	<i>N (%)⁺</i>	<i>Examples of support for judgement from the review*</i>
<i>blinding of participants and personnel</i>	Assess risk differently if blinding was not feasible because of the type of intervention	20 (36)	"Not possible to blind participants"; Low risk of bias "Participants were not blinded for provided treatment. This is inherent to study design"; High risk of bias
	outcome considered not influenced by blinding	12 (21)	"No information given about whether patients were blind to physician allocation but treatment outcomes judged unlikely to be affected by lack of blinding"; Low risk of bias
	Consider differently information of "double blind"	9 (16)	"Quote: ". . . patients were randomised in double-blind conditions . . ."Comment: probably done"; Low risk of bias "Quote: "double blind conditions". No further details."; Unclear risk of bias
	Consider differently incomplete or unclear description	7 (12)	"Researchers were blind until after the baseline assessment. participants were not blinded."; Unclear risk of bias "Not possible to blind participants to intervention. Insufficient information to make a judgement about blinding of therapists"; High risk of bias
	Confusion in the definition of the item	5 (9)	"Described as an "open-label" pilot study."; Low risk of bias
<i>blinding of outcome assessment</i>	Consider differently incomplete or unclear description	24 (34)	"Not explicitly discussed in the publish study, it was assumed to be open label"; High risk of bias "Not described in published study"; Unclear risk of bias
	outcome considered not influenced by blinding	16 (23)	"Not stated, but it was unlikely that the outcome was influenced by lack of blinding"; Low risk of bias
	Consider differently patient-reported outcomes when patients are blinded or not to the intervention	9 (13)	"Comment: depression assessed by patient self-report"; High risk of bias "Insufficient information available to assess"; Low risk of bias
	Consider differently information of "double blind"	9 (13)	"Quote: ". . . double blind" Comment: probably done"; Low risk of bias "Quote: "double blind conditions". No further details."; Unclear risk of bias
	Assess risk differently if blinding was not feasible because of the type of intervention	6 (9)	"blinding not possible due to intervention"; High risk of bias "Unclear blinding of outcome assessment"; Low risk of bias

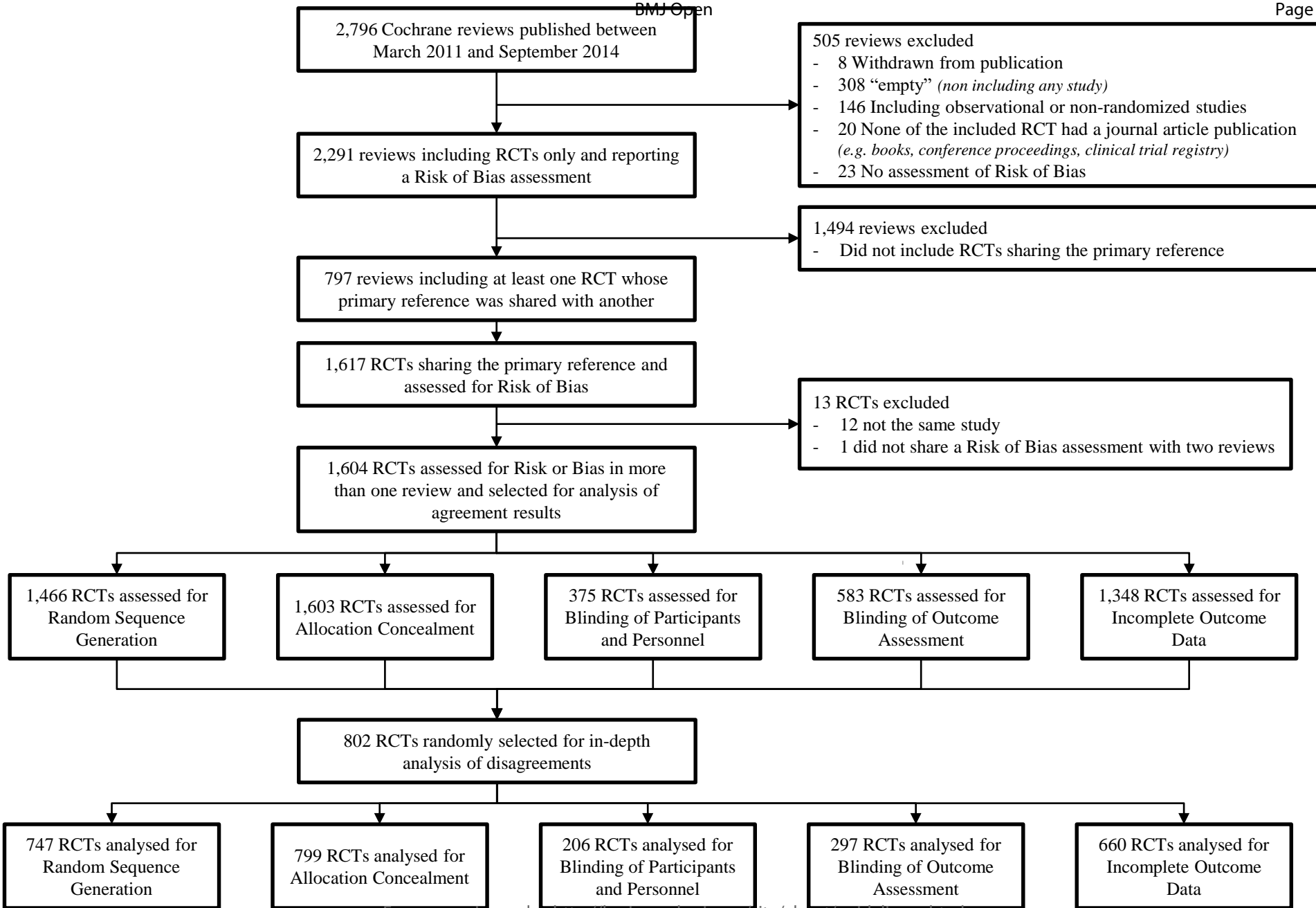
⁺ Number of RCTs disagreeing for this reason; percentage over the total of disagreements for different interpretation.

^{*} When two extracts are reported, they refer to the same study.

<i>Risk of bias Item</i>	Main reasons for disagreements	N (%)⁺	Examples of support for judgement from the review*
<i>incomplete outcome data</i>	Use different cut-off for the rate of missing data	57 (26)	“11 withdrawals (10%).”; Low risk of bias “Comment: there were post-randomisation drop-outs”; High risk of bias
	Focus on number vs reasons/precise report of missing data	28 (13)	“20 drop-outs (27.2%) with 4 deaths (3 males, 1 female) from cardiovascular events”; High risk of bias “Numbers and reasons for dropouts and withdrawals in all intervention groups were described.”; Low risk of bias
	Consider differently incomplete or unclear description	27 (12)	“Women who were untraceable or unsuitable for follow-up were excluded, other losses included as smokers”; Low risk of bias “167/1287 (12.9%) (C = 83, I = 84) excluded from analysis due to moving away, being untraceable or deemed unsuitable for follow-up (e.g. miscarriage). 1120 in sample. 51/1287 non-responders were included as continuing smokers.” High risk of bias
	Consider differently intention-to-treat analysis	25 (11)	“147 randomised; 4 in the letrozole group and 3 in the LOD dropped out of the trial, all for non-compliance. However, ITT analysis was not conducted.”; Unclear risk of bias “7 women lost to follow up, but similar (3 vs 4) in both groups; losses due to noncompliance”; Low risk of bias
	Consider differently report of “no missing data”	22 (10)	“Did not report number of withdrawals. Comment: all patients who were randomised were included in the final analysis. ITT analysis was conducted.”; Unclear risk of bias “It does not appear that there were any withdrawals or dropouts” Low risk of bias
	Consider differently imputation of missing data	20 (9)	“Imputation method not described”; Unclear risk of bias “Dropout rate was not significant”; Low risk of bias
	Use different cut-off for difference in the rate missing data between different arms/comparisons	13 (6)	“Dropout higher in placebo group (35% vs 25% in budesonide group). ITT used.”; High risk of bias “Similar rates of withdrawal between arms. Withdrawals: 36 BUD, 51 placebo”; Low risk of bias

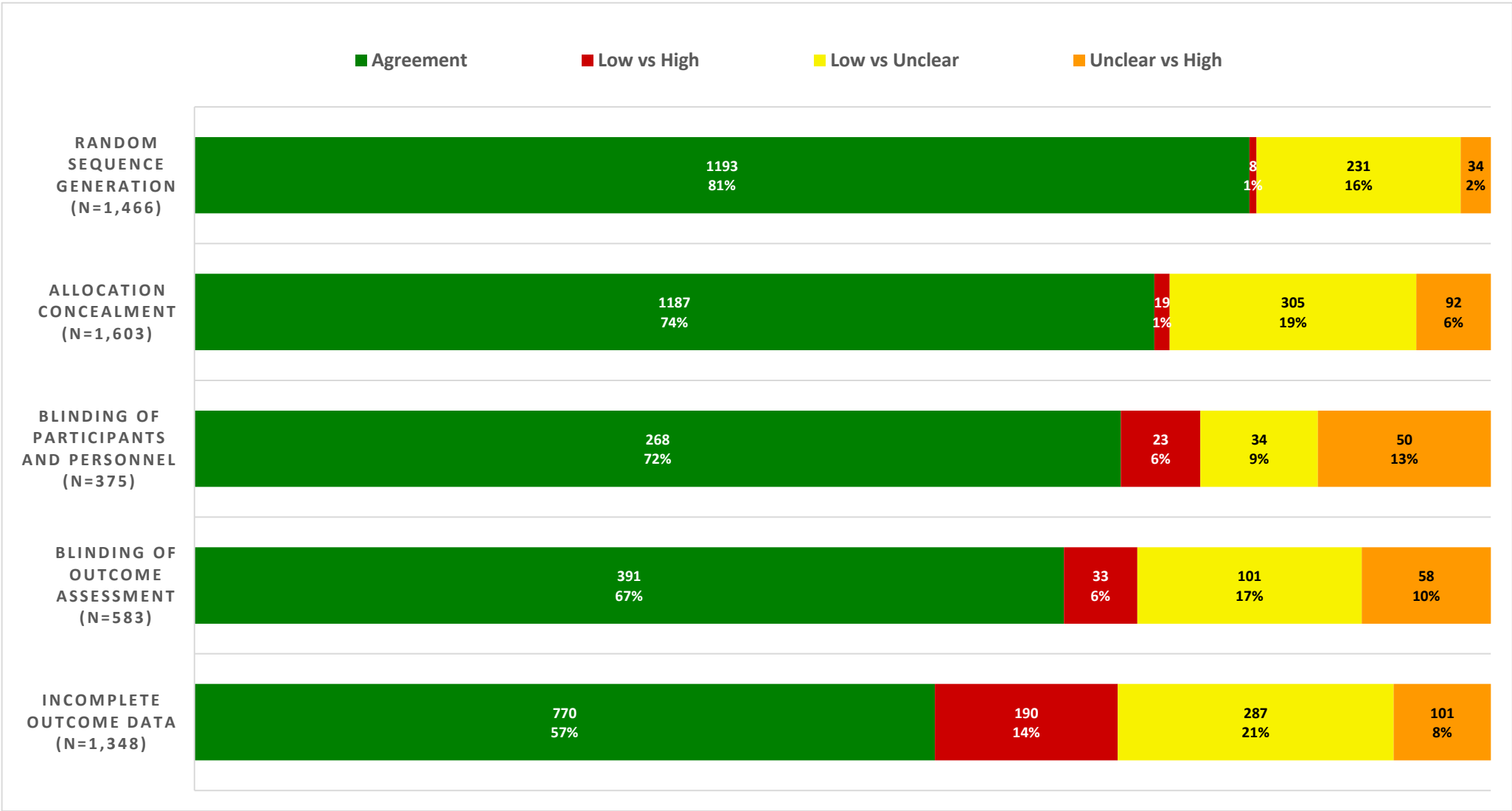
⁺ Number of RCTs disagreeing for this reason; percentage over the total of disagreements for different interpretation.

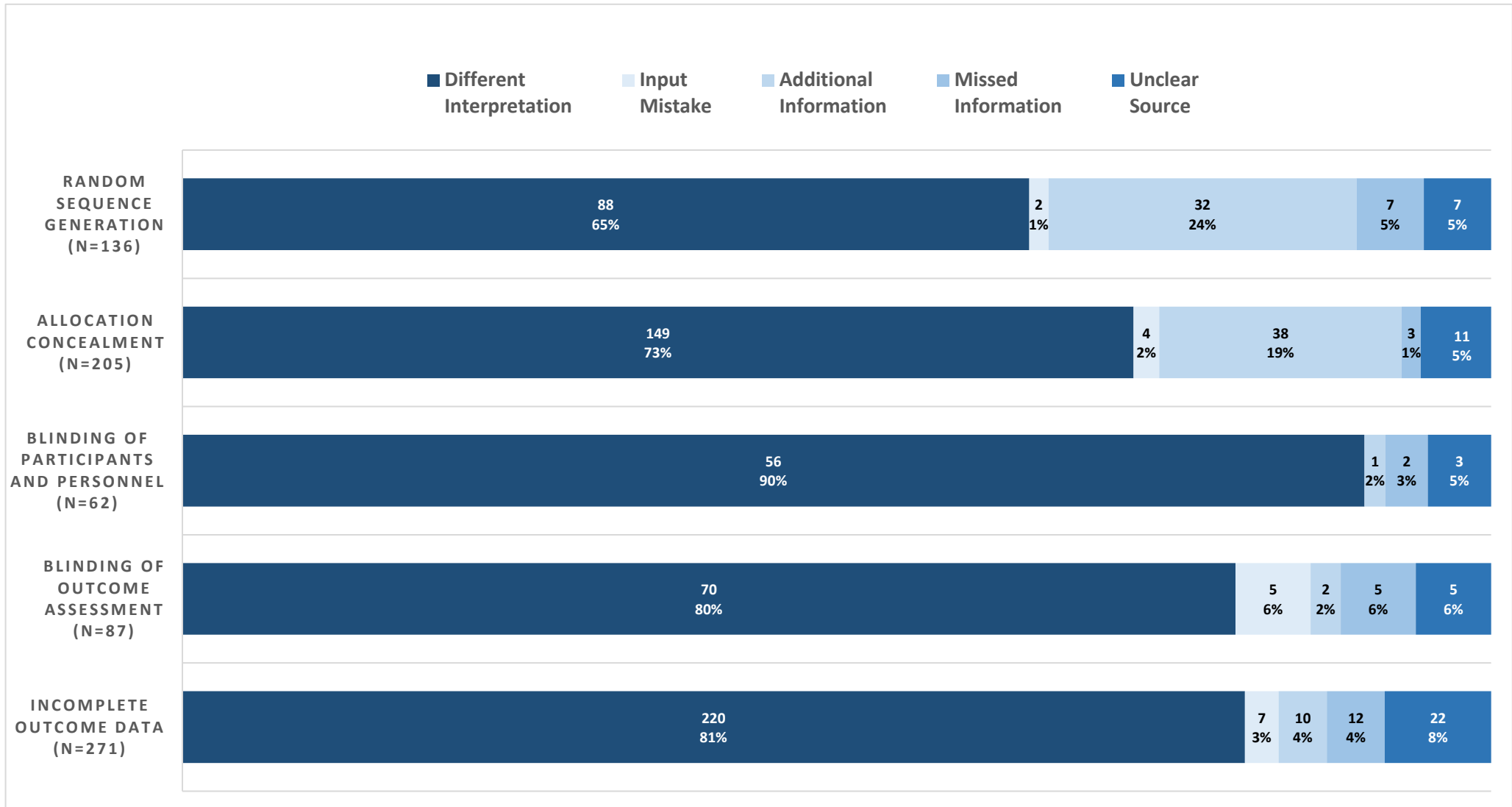
* When two extracts are reported, they refer to the same study.



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46





1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Same support for judgement

Risk of bias assessment does not match the support for judgement
(e.g. "Randomization described explicitly", judgement "Unclear")

Input mistake

One review confuses one item with a different one / misunderstanding of definition of the item
(e.g. for Random Sequence Generation "600 opaque envelopes, I was drawn every time")

Differences in interpretation

Different support for judgement

Access to the study report

Information in the report is incomplete or unclear

Study report clearly describes the information, but one review seemed to have missed it

Missed information from the study report

Study report does not describe the information reported by one author

No access to study report

Mention of access to additional information in the review

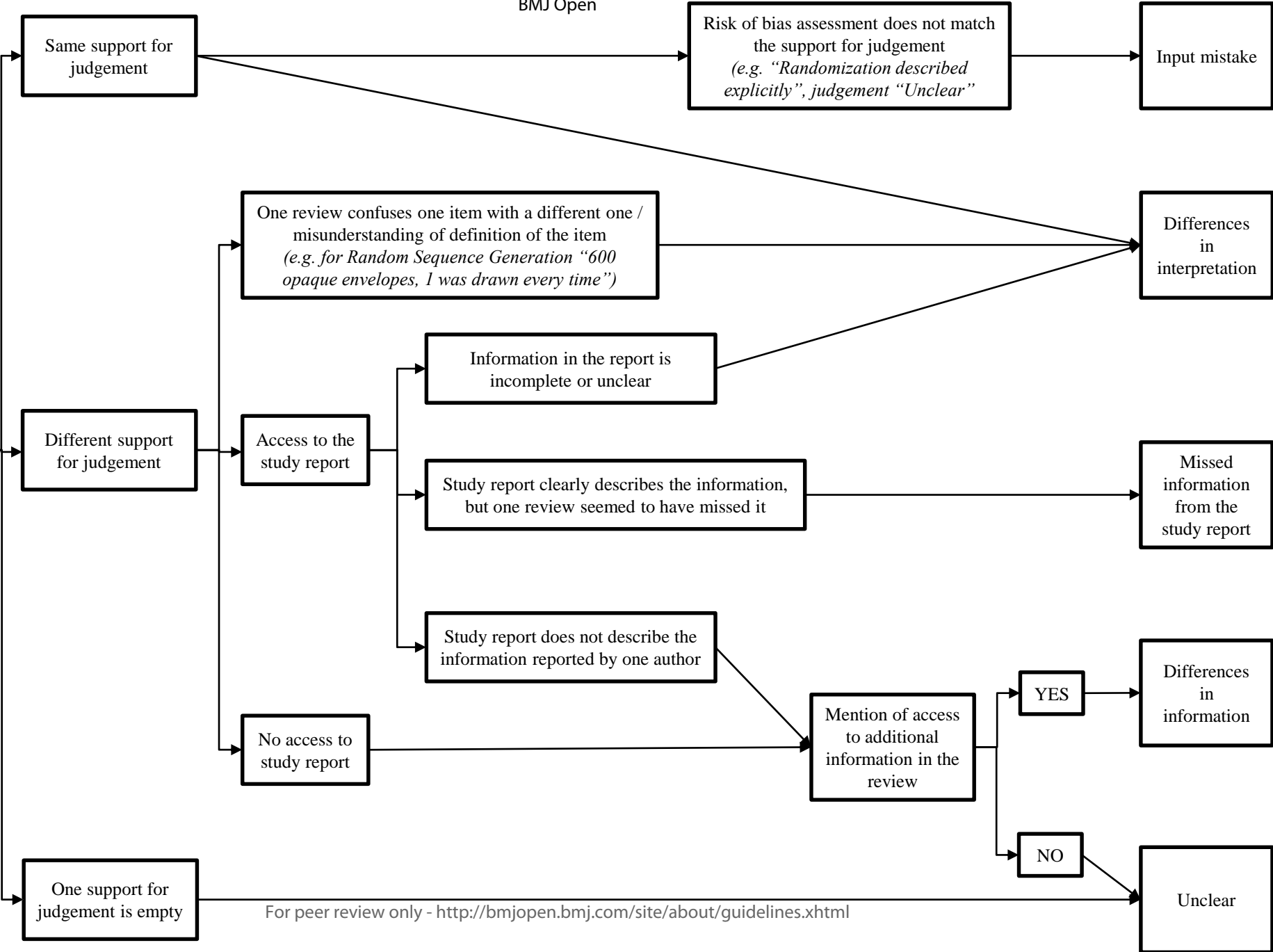
YES

Differences in information

NO

Unclear

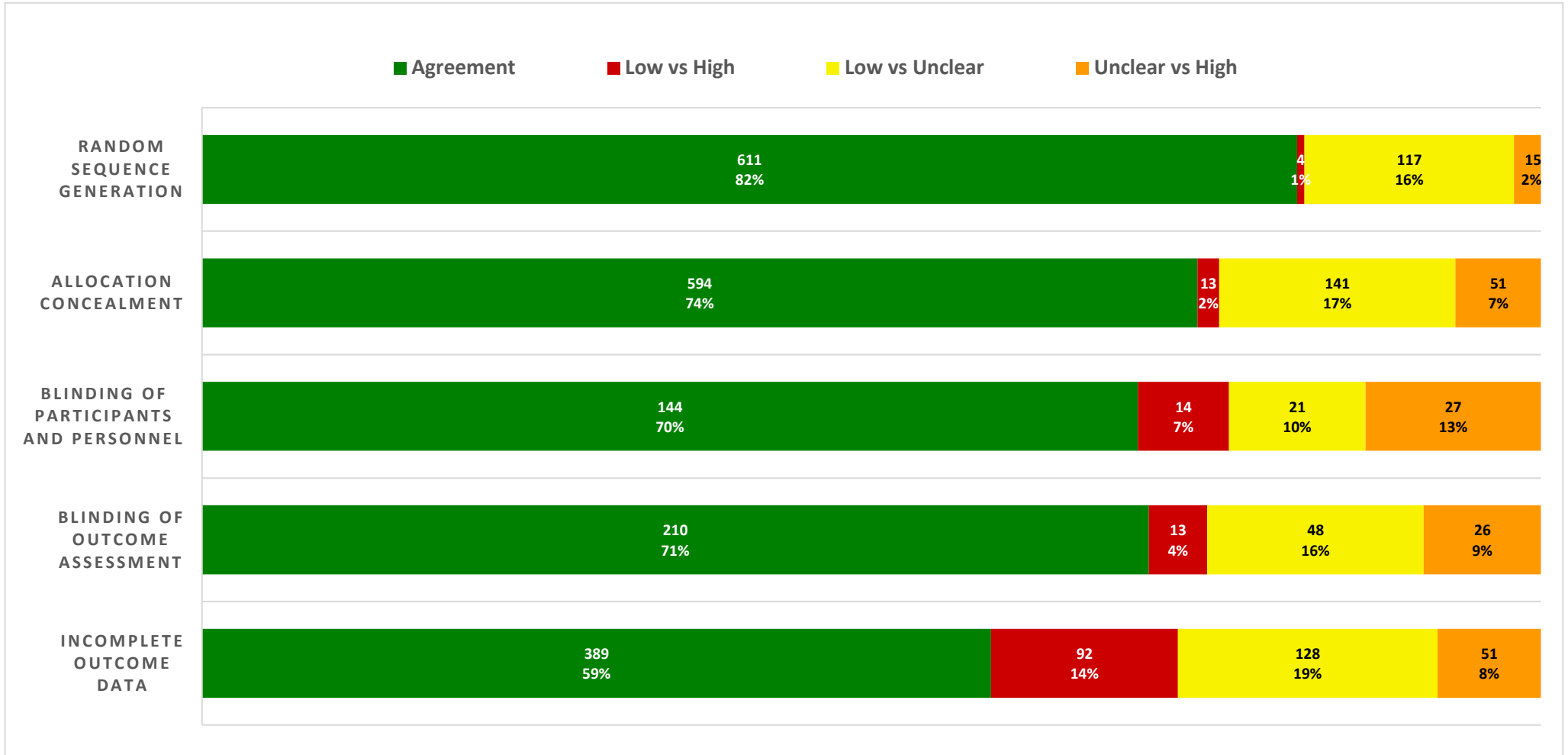
One support for judgement is empty



	Cochrane Group	Number of reviews	% on the total
	Pregnancy and Childbirth	93	11.7%
	Airways	48	6.0%
	Pain, Palliative and Supportive Care Group	42	5.3%
	Acute Respiratory Infections	37	4.6%
	Gynaecology and Fertility	29	3.6%
	Neonatal	29	3.6%
	Tobacco Addiction	27	3.4%
	Stroke	25	3.1%
	Gynaecological, Neuro-oncology and Orphan Cancer Group	23	2.9%
	Wounds	23	2.9%
	Hepato-Biliary	22	2.8%
	Cystic Fibrosis and Genetic Disorders	21	2.6%
	Anaesthesia	20	2.5%
	Drugs and Alcohol	20	2.5%
	Neuromuscular	19	2.4%
	Common Mental Disorders	18	2.3%
	Fertility Regulation	18	2.3%
	Heart	17	2.1%
	Developmental, Psychosocial and Learning Problems	16	2.0%
	Incontinence	16	2.0%
	Kidney disease	16	2.0%
	Schizophrenia	16	2.0%
	Infectious Diseases	15	1.9%
	Oral Health	14	1.8%
	Vascular	14	1.8%
	Dementia and Cognitive Improvement	13	1.6%
	Musculoskeletal	12	1.5%
	Consumers and Communication	10	1.3%
	Epilepsy	10	1.3%
	Eyes and Vision	9	1.1%
	Metabolic and Endocrine Disorders	9	1.1%
	Back and Neck	8	1.0%
	Hypertension	8	1.0%
	Multiple Sclerosis	8	1.0%
	Effective Practice and Organisation of Care	7	0.9%
	HIV/AIDS	7	0.9%
	Inflammatory Bowel Disease	7	0.9%
	Injuries	7	0.9%
	Bone, Joint and Muscle Trauma Group	6	0.8%
	ENT	6	0.8%
	Haematological Malignancies	6	0.8%

Cochrane Group	Number of reviews	% on the total
Breast Cancer	5	0.6%
Colorectal Cancer	5	0.6%
Lung Cancer	3	0.4%
Movement Disorders	3	0.4%
Skin	3	0.4%
Occupational Health	2	0.3%
Sexually Transmitted Infections	2	0.3%
Public Health	1	0.1%
Upper GI and Pancreatic Diseases	1	0.1%
Urology	1	0.1%
Total	797	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46



Supplementary Appendix 4| Examples of in-depth analysis of disagreements conducted with the support of the study report

Risk of bias item	Study Name	Support for judgement*	Information in the study report^	Category of disagreement	Reason of disagreement
Random sequence generation	ABCD 2004	Review 4136: Generated the randomisation list using SAS, stratified by sex and SCr; Low Risk	<i>“The (...) statistician generated the randomization list using SAS (...) stratified by sex and baseline serum creatinine concentration (...).”</i>	Missed information from the study report	
		Review 8277: Method not reported; Unclear Risk			
	Cho 2006	Review 7566: Stated that it is a quasi-randomized study but details not given; High Risk	<i>“... using a quasi-experimental design with a non-equivalent control group.”</i> <i>“They were randomly assigned to participate in the experimental group (...) or a waiting-list control group (...).”</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 9553: Participants randomly allocated to treatment or control group; Unclear Risk			
	Petersen 2005	Review 9132: Quote: “[P]atients were randomly assigned...”Quote: “We used an adaptive allocation scheme for the treatment assignment, with the MMSE score, age and APOE e4 status as balancing covariates”; Low Risk	<i>“We used an adaptive allocation scheme for the treatment assignment, with the MMSE score, age, and APOE e 4 status as balancing covariates.”</i>	Different interpretation	Confusion or misknowledge
		Review 7176: The trial is described as randomised, but the method of sequence generation was not specified. Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report^	Category of disagreement	Reason of disagreement
Allocation concealment	Burge 2000	Review 2991: Participants were randomly assigned sequentially from a list comprising treatment numbers only; Low Risk	<i>"We used a computer generated allocation schedule stratified by centre (block size of six). Patients were randomised sequentially from a list comprising treatment numbers only".</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 10115: Information not available; Unclear Risk			
	McMurdo 1993	Review 4294: Quote: "Randomisation was by opening sealed envelopes supplied in sequence by the study co-ordinator; Low Risk	<i>"Randomization was by opening sealed envelopes supplied in sequence by the study co-ordinator (...), and prepared from a computer-generated random numbers table."</i>	Different interpretation	Consider differently envelopes description
		Review 4963: Unclear, insufficient reporting to permit judgement; Unclear Risk			
	Draper 2007	Review 8179: "... and alternating between treatment or wait list control groups."; High Risk	<i>"On each occasion that a least eight patients had been recruited, their names were selected at random by a blinded investigator to be allocated alternately to the immediate treatment group or a wait-list control group."</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 1919: "Reported as concealed but specific method for concealment not reported"; Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report^	Category of disagreement	Reason of disagreement
Blinding of participants and personnel	Nielsen 2006	Review 9672: "Double-blind"; Low Risk	<p><i>"This study was a randomized, placebo-controlled, double blind, Danish, multi-center (two centers) study."</i></p> <p><i>"The treatment was applied by a nasal spray with one puff in each nostril every day either in the morning or evening."</i></p>	Different information	One review accessed additional data through another study report
		Review 4143: Although "All treatments were supplied as identical intranasal sprays..." the 2004 publication describes a higher rate of withdrawal due to adverse effects in the intervention groups [11.7% in the placebo group, 21.7% in the 150 gm group and 28.7% in the 300 gm group} which may have affected blinding status; Unclear Risk			
	Gersel 1979	Review 10562: Described as double-blind [presumed participants and personnel/investigators]; Low Risk	<p><i>"A double-blind experimental design was used, employing each patient as his own control."</i></p>	Different interpretation	Consider differently information of "double blind"
		Review 6968: Not mentioned and no information to suggest this was done.; Unclear Risk			
	Stein 2011	Review 7025: Not possible to blind participants to intervention. Insufficient information to make a judgement about blinding of therapists; High Risk	<p><i>"Follow-up assessment was made 3 months after release (research staff conducting assessments were blind to treatment assignment)."</i></p> <p><i>"Randomization was accomplished via random numbers table in advance and placed in an envelope by the project coordinator. Following baseline assessment, research staff opened the envelope to learn of intervention assignment."</i></p>	Different interpretation	Consider differently incomplete or unclear description
		Review 10901: Researchers were blind until after the baseline assessment. Participants were not blinded.; Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report^	Category of disagreement	Reason of disagreement
Blinding of outcome assessment	Schoen 2007	Review 3603: Outcome assessor was not blinded.; High risk	<p><i>"In total 72 patients were screened by a maxillofacial surgeon (PJS) and prosthodontist (HR)."</i></p> <p><i>"All clinical assessments were performed by the investigator (PJS) who was not involved in treatment of the patients."</i></p>	Different interpretation	Consider differently incomplete or unclear description
		Review 5005: Outcome assessor may have been unaware of allocation: "All clinical assessments were performed by the investigator (PJS) who was not involved in treatment of the patients."; Low risk			
	Geroin 2011	Review 6185: Not done; High risk	<p><i>"All patients were evaluated by the same examiner (an experienced internal coworker) who was not aware of the treatment received by the patients"; Low Risk</i></p>	Missed information from the study report	
		Review 9645: Quote: "All patients were evaluated by the same examiner (an experienced internal coworker) who was not aware of the treatment received by the patients"; Low Risk			
	McCambridge 2004	Review 8969: As one interventionist was the study PI, a second independent interviewer who was blind to study condition was employed to conduct 3 month follow-ups, and an additional interviewer who was blind to initial group allocation was employed for 12 months follow-ups; Low Risk	<p><i>"further area of possible bias was that intervention recipients might report more favourable outcome data to the researcher who had delivered the intervention (J.M.). To study any such bias, a second independent interviewer who was blind to study condition, was employed to interview a sample of participants."</i></p>	Different interpretation	Consider differently incomplete or unclear description
		Review 7025: A second independent interviewer who was blind to study condition was employed to interview a sample of participants, though not all participants; Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report [^]	Category of disagreement	Reason of disagreement
Incomplete outcome data	Altmaier 1992	Review 1822: All subjects recorded follow-up data; Low Risk	<i>[From table] “The n = 21 for control group and n = 24 for psychological group on all process measures.”</i> <i>[From table] The n = 21 for each group at each assessment.]</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 7407: Inadequately reported; High Risk			
	Killen 1984	Review 146: 11/75 recruited dropped out before full treatment, and are excluded from analyses.; Low Risk	<i>“The first 75 were accepted into the study. Seven failed to attend (...) two dropped (...). The final sample (N = 64).”</i>	Missed information from the study report	
		Review 3999: Losses to follow-up not reported, all participants included; Unclear Risk			
	Creager 2008	Review 986: There was a huge loss to follow up (only 50% completed the 6 month follow up) in this study and therefore there is a high risk of attrition bias; High Risk	<i>“The remaining 525 patients met the inclusion criteria (...) The remaining 430 patients met their criteria for randomization (...). The ITT population consisted of 370 randomized patients (...). The per-protocol patient population consisted of 214 randomized patients”</i>	Different interpretation	Consider differently intention-to-treat analysis
		Review 5262: Unclear why of patients stopped medication, unclear whether data presented represents intention-to-treat or per-protocol analysis			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

[^] Information that were highlighted in the study report to support the analysis process.

Supplementary Appendix 5| Reasons for disagreements in cases of a different interpretation of the same information; focus on “low” versus “high” disagreements.

Risk of bias item	Main reasons for disagreements	Examples
random sequence generation	Consider differently incomplete or unclear description	<p>“The names of communities within each group of three were written on individual cards, mixed and selected randomly: the first from each group was assigned to arm A (IEC alone), the second to arm B (IEC and STI management) and the third to arm C”; Low risk of bias</p> <p>“Names of communities within each triplet were written on separate cards and shuffled.”; High risk of bias</p>
allocation concealment	Consider differently envelopes description	<p>“Closed envelopes”; Low risk of bias</p> <p>“Closed envelopes, although not opaque.”; High risk of bias</p>
	Confusion in the definition of the item	<p>“pg. 2 - Methods - randomisation was done centrally to preserve allocation concealment”; Low risk of bias</p> <p>“904 patients were eligible for the study. 446 patients were randomised (49%). Due to the number of patients declining screening, there is an increased risk of inclusion bias.”; High risk of bias</p>
	Confusion with blinding	<p>“States used “preprogrammed laptop computer”. Remote site”; Low risk of bias</p> <p>“participants were told to which compound they had been allocated.”; High risk of bias</p>
blinding of participants and personnel	Assess risk differently if blinding was not feasible because of the type of intervention	<p>“Not possible to blind participants”; Low risk of bias</p> <p>“Participants were not blinded for provided treatment. This is inherent to study design”; High risk of bias</p>
	Outcome considered not influenced by blinding	<p>“Not possible to blind but most of the outcomes not likely to be influenced by lack of blinding.”; Low risk of bias</p> <p>“Not blinded due to nature of intervention.”; High risk of bias</p>
	Confusion with allocation concealment	<p>“participants were randomly allocated to either intervention or control group by an independent party”; Low risk of bias</p> <p>“Control group did not receive the comparable non-exercise related attention to the intervention group”; High risk of bias</p>

Risk of bias Item	Main reasons for disagreements	Examples
blinding of outcome assessment	<p>outcome considered not influenced by blinding</p> <hr/> <p>Consider differently patient reported outcomes when patients are blinded or not to the intervention</p> <hr/> <p>Assess risk differently if blinding was not feasible because of the type of intervention</p>	<p>“No information given about whether patients or assessors were blind to physician allocation but primary outcomes (treatment outcome and patient reported physician cultural competency) judged unlikely to be affected by lack of blinding”; Low risk of bias</p> <p>“Unblinded.”; High risk of bias</p> <hr/> <p>“Insufficient information available to assess”; Low risk of bias</p> <p>“Comment: depression assessed by patient self-report”; High risk of bias</p> <hr/> <p>“Unclear blinding of outcome assessment”; Low risk of bias</p> <p>“blinding not possible due to intervention”; High risk of bias</p>
incomplete outcome data	<p>Use different cut-off for the rate of missing data</p> <hr/> <p>Focus on number vs reasons/precise report of missing data</p> <hr/> <p>Consider differently incomplete or unclear description</p> <hr/> <p>Use different cut-off for difference in the rate missing data between different arms/comparisons</p>	<p>“11 withdrawals (10%).”; Low risk of bias</p> <p>“Comment: there were post-randomisation drop-outs”; High risk of bias</p> <hr/> <p>“Numbers and reasons for dropouts and withdrawals in all intervention groups were described.”; Low risk of bias</p> <p>“20 drop-outs (27.2%) with 4 deaths (3 males, 1 female) from cardiovascular events”; High risk of bias</p> <hr/> <p>“Women who were untraceable or unsuitable for follow-up were excluded, other losses included as smokers”; Low risk of bias</p> <p>“167/1287 (12.9%) (C = 83, I = 84) excluded from analysis due to moving away, being untraceable or deemed unsuitable for follow-up (e.g. miscarriage). 1120 in sample. 51/1287 non-responders were included as continuing smokers.” High risk of bias</p> <hr/> <p>“Similar rates of withdrawal between arms. Withdrawals: 36 BUD, 51 placebo”; Low risk of bias</p> <p>“Dropout higher in placebo group (35% vs 25% in budesonide group). ITT used.”; High risk of bias</p>

BMJ Open

Disagreements in risk of bias assessment for randomised controlled trials included in more than one Cochrane systematic reviews: a research on research study using cross-sectional design

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-028382.R1
Article Type:	Research
Date Submitted by the Author:	14-Feb-2019
Complete List of Authors:	Bertizzolo, Lorenzo; INSERM, U 1153, Equipe Methods Bossuyt, Patrick; Academic Medical Center; University of Amsterdam, Dept. Clinical Epidemiology and Biostatistics Atal, Ignacio; INSERM U1153, Team Methods Ravaud, Philippe; INSERM, U1153, Epidemiology and Biostatistics Sorbonne Paris Cite Research Center (CRESS), Methods of therapeutic evaluation of chronic diseases team (METHODS) Dechartres, Agnes; INSERM U738 H+pital H+telDieu Universit+ParisDescartes, Centre dEpidemiologie Clinique
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Evidence based practice, Public health, Qualitative research, Research methods
Keywords:	risk of bias, Cochrane, systematic reviews, interrater agreement, reproducibility, Public health < INFECTIOUS DISEASES

SCHOLARONE™
Manuscripts

1
2
3 **Disagreements in risk of bias assessment for randomised controlled trials included in**
4 **more than one Cochrane systematic reviews: a research on research study using cross-**
5 **sectional design**
6
7
8
9
10
11
12
13
14

15 Lorenzo Bertizzolo¹, Patrick M Bossuyt², Ignacio Atal^{1, 5}, Philippe Ravaud^{1, 3-6}, Agnès
16 Dechartres⁷
17
18
19
20
21

22 ¹ INSERM, U1153 Epidemiology and Biostatistics Sorbonne Paris Cité Research Center
23 (CRESS), Methods of therapeutic evaluation of chronic diseases Team (METHODS), Paris,
24 F-75004 France; Paris Descartes University, Sorbonne Paris Cité, France.
25
26

27 ² Department of Clinical Epidemiology, Biostatistics and Bioinformatics, Academic Medical
28 Center, University of Amsterdam, Netherlands.
29
30

31 ³ Centre d'Épidémiologie Clinique, Hôpital Hôtel Dieu, AP-HP (Assistance Publique des
32 Hôpitaux de Paris), Paris, France.
33
34

35 ⁴ Faculté de Médecine, Université Paris Descartes, Sorbonne Paris Cité, Paris, France.
36
37

38 ⁵ Cochrane France, Paris, France
39
40

41 ⁶ Columbia University, Mailman School of Public Health, Department of Epidemiology, New
42 York, USA
43
44

45 ⁷ Sorbonne Université, INSERM, Institut Pierre Louis de Santé Publique, Département
46 Biostatistique, Santé Publique et Information Médicale, AP-HP, Hôpitaux Universitaires Pitié
47 Salpêtrière – Charles Foix, Paris, France
48
49
50
51

52 Correspondence to: Lorenzo Bertizzolo, M.D., MPH
53
54
55
56
57
58
59
60

INSERM, U1153 - Centre d'Epidémiologie Clinique –

Hôpital Hôtel-Dieu,

1, place du parvis Notre Dame, 75004 Paris, France

Tel: +33 (0)1 42 34 78 25

E-mail: lorenzo.bertizzolo@gmail.com

Key-words: risk of bias, Cochrane, Systematic Reviews, Interrater Agreement, Reproducibility

Word Count: 3800

Copyright: The Corresponding Author has the right to grant on behalf of all authors and does grant on behalf of all authors, a worldwide licence to the Publishers and its licensees in perpetuity, in all forms, formats and media (whether known now or created in the future), to i) publish, reproduce, distribute, display and store the Contribution, ii) translate the Contribution into other languages, create adaptations, reprints, include within collections and create summaries, extracts and/or, abstracts of the Contribution, iii) create any other derivative work(s) based on the Contribution, iv) to exploit all subsidiary rights in the Contribution, v) the inclusion of electronic links from the Contribution to third party material where-ever it may be located; and, vi) licence any third party to do any or all of the above.

Transparency declaration: The guarantor (Lorenzo Bertizzolo) affirms that the manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

Ethics approval: Not applicable. This is a research on research study.

ABSTRACT (299 words)

Objectives: Assess the frequency and reasons for disagreements in risk of bias assessments for randomised clinical trials (RCTs) included in more than one Cochrane review.

Design: Research on research study, using cross-sectional design.

Data sources: 2,796 Cochrane reviews published between March 2011 and September 2014.

Data selection: RCTs included in more than one review.

Data extraction: Risk of bias assessment and support for judgement for five key risk of bias items.

Data synthesis: For each item, we compared risk of bias assessment made in each review and calculated proportion of agreement. Two reviewers independently analysed 50% of all disagreements by comparing support for each judgement with information from study report to evaluate whether disagreements were related to a difference in information (e.g., contact the study author) or a difference in interpretation (same support for judgement but different interpretation). They also identified main reasons for different interpretation.

Results: 1,604 RCTs were included in more than one review. Proportion of agreement ranged from 57% (770/1,348 trials) for incomplete outcome data to 81% for random sequence generation (1,193/1,466). Most common source of disagreement was difference in interpretation of the same information, ranging from 65% (88/136) for random sequence generation to 90% (56/62) for blinding of participants and personnel. Access to different information explained 32/136 (24%) disagreements for random sequence generation and 38/205 (19%) for allocation concealment. Disagreements related to difference in interpretation were frequently related to incomplete or unclear reporting in the study report (83% of disagreements related to different interpretation for random sequence generation).

1
2
3 **Conclusions:** Risk of bias judgements of RCTs included in more than one Cochrane review
4 differed substantially. Most disagreements were related to a difference in interpretation of an
5 incomplete or unclear description in the study report. A clearer guidance on common causes
6 of incomplete information may improve agreement.
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Article summary

Strengths and limitations of this study

- Use of a very large and comprehensive collection of Cochrane reviews to assess the agreement in risk of bias assessment and to understand reasons of disagreement.
- Analysis of the full-text of study reports to underline what information were available to review authors and how they utilized them while assessing risk of bias.
- Focus on disagreements only. Possible that a proportion of agreements happened “by chance”. For example review authors may express the same risk of bias judgement while using different information or interpreting information differently.
- No evaluation of the potential impact of disagreements in conclusion making at the review level.

INTRODUCTION

Systematic reviews aim to synthesise all existing evidence for a research question by the use of a rigorous and reproducible methodology¹. Because reviews may be affected by bias at the level of individual studies², an assessment of the risk of bias in these studies is a crucial step in conducting a systematic review^{3 4}.

Cochrane has developed a tool to provide a standardised approach to the assessment of the risk of bias in randomised clinical trials (RCTs)⁵. The risk of bias tool is based on specific characteristics related to study design and conduct, selected on theoretical grounds and on empirical evidence from meta-epidemiological studies that these characteristics are associated with differences in treatment effect estimates⁶⁻¹¹. The tool includes seven items (random sequence generation, allocation concealment, blinding of participants and personnel, blinding of outcome assessment, incomplete outcome data, selective reporting, other source of bias), the researchers assess and judge as either “high”, “low” or “unclear” risk of bias^{11 12}. Although Cochrane provides detailed guidance on how to use the tool and recommends consensus between two independent reviewers¹¹, personal judgement is also involved, which may bring variability. Several studies have evaluated the reproducibility of the risk of bias tool, generally shown to be poor¹²⁻¹⁹. However, there is some uncertainty about the main causes of disagreements. For example, some reviewers may search for additional information such as protocols or contact study authors and this difference in available information, rather than a difference in judgement, may explain some of the disagreements.

In this study, we used a large collection of Cochrane reviews to evaluate the reproducibility of risk of bias assessments by identifying randomised controlled trials included in more than one Cochrane review and comparing the assessments. In addition, we examined the likely

1
2
3 reasons for any disagreements. In particular, we evaluated whether disagreements were
4
5 related to differences in information available to reviewers or differences in interpreting the
6
7 same information and what could explain such different interpretation.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

METHODS

This is a research on research study on risk of bias assessment, which used a cross sectional design. We identified RCTs included in more than one reviews included in a large collection of Cochrane reviews. For key risk of bias items, we evaluated agreement between the different systematic reviews; analysed whether disagreements were related to a difference in information available to reviewers or a difference in interpretation of the same information and highlighted the main reasons for disagreements by an in-depth, one-by-one evaluation of disagreements.

Data sources

We obtained data from the 2,796 Cochrane reviews, which corresponds to all reviews available in the Cochrane library between March 2011 and September 2014, including updates (March 2011 corresponds to the last update of the risk of bias tool⁵). Data consisted of one XML file per review, each file containing all data entered by review authors in RevMan, the software used for managing Cochrane reviews²⁰. All individual XML files were merged in a single database by using R v3.2.2²¹ with the XML package²². The vocabulary used for risk of bias items slightly varied across reviews (e.g., some reviews could refer to “allocation concealment” as “allocation masking”). For this reason, two authors independently evaluated all terms used and classified them according to the vocabulary of the tool. Disagreements were resolved by consensus. This standardization was done for a previous publication²³.

Selection of eligible reviews

1
2
3 We excluded withdrawn or “empty” reviews (i.e., systematic reviews not including any
4 study) as well as reviews including observational or non-randomised studies and considered
5 only reviews with an assessment of risk of bias for at least one item of the risk of bias tool.
6
7
8
9

10 11 12 ***Selection of eligible RCTs*** 13

14 To identify single RCTs included and assessed for risk of bias in more than one systematic
15 review, we proceeded as follows. For each RCT, we identified the primary reference(s),
16 which was the reference identified by review authors as the main reference(s) for an included
17 study. Then, we used a matching algorithm²⁴ to identify studies that shared the same primary
18 reference. If several primary references were reported, we considered all of them. We
19 manually checked that the studies sharing the same primary reference in the reviews
20 corresponded to the same RCT.
21
22
23
24
25
26
27
28
29
30

31 32 33 ***Extraction of risk of bias assessment*** 34

35 For each eligible RCT, we extracted the risk of bias assessment and the corresponding
36 support for judgement for each risk of bias item in each review. Whenever a single RCT was
37 included in three or more reviews, we considered only the risk of bias assessment from two
38 reviews chosen at random; this situation concerned less than 10% of our included RCTs and
39 was decided because of workload and to facilitate direct comparison of two assessments. We
40 focused on five risk of bias items: random sequence generation, allocation concealment,
41 blinding of participants and personnel, blinding of outcome assessment and incomplete
42 outcome data. We did not consider selective reporting because it is difficult to evaluate in the
43 absence of the study protocol, which is frequently lacking, especially for older studies^{11 12 14}.
44
45 We also did not consider the item other bias because the definition is very wide (i.e., “any
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 important concerns about bias not covered in the other domains in the tool”¹¹), so
4
5 comparisons across reviews are difficult.
6
7
8
9

10 ***Comparison of risk of bias assessment between reviews***

11
12 For each item, we compared the risk of bias assessment in terms of “high”, “low” or
13
14 “unclear” risk of bias between the two reviews. According to the Cochrane handbook, the
15
16 items blinding of outcome assessment and incomplete outcome data should be assessed for
17
18 each outcome. Therefore, when the reviews reported an assessment of these items at the
19
20 outcome level, we manually checked that outcomes were identical in both reviews and we
21
22 retained for our analysis only the assessments that focused on the same outcomes. For
23
24 blinding, we followed the last version of the Cochrane handbook and we retained only
25
26 assessments of blinding of participants and personnel and blinding of outcome assessment as
27
28 two independent items, excluding different types of assessment (i.e., blinding as a single
29
30 item, blinding of only participants or of only personnel).
31
32
33

34
35 We calculated the percentage agreement for each risk of bias item, as the proportion of
36
37 studies with a concordant assessment in both reviews (e.g. “low” risk of bias AND “low” risk
38
39 of bias). Not all reviews assessed all five key risk of bias items for each RCT included;
40
41 consequently, the number of RCTs evaluated for discrepancies varies depending on the item
42
43 considered.
44
45
46
47
48

49 ***Selection of studies for in-depth analysis of disagreements***

50
51 For workload reasons, we in-depth evaluated the reasons for disagreements for 50% of the
52
53 studies analysed in the previous step. In cases of more than one shared RCTs within a given
54
55
56
57
58
59
60

1
2
3 pair of Cochrane reviews, we selected only one RCT at random. To reach 50% of the total
4
5 sample, we used a simple random selection in the remaining database.
6
7
8
9

10 *Classification of disagreements*

11
12 For the random selection, two reviewers (LB and AD) independently evaluated all
13
14 disagreements in the risk of bias assessment in the two systematic reviews. They first
15
16 scrutinised the support for the judgement in each review and evaluated whether it was the
17
18 same or “conceptually” the same in both reviews (e.g., “randomised, probably done”;
19
20 “randomised, probably not done”; “study only mention randomization, but does not specify
21
22 how randomization was performed; unclear”; “study states it is randomized; low risk”). If the
23
24 support differed, they assessed any other information regarding the study as reported in both
25
26 reviews, systematically searching and evaluating the full-text study report indicated in the
27
28 primary reference. A formalized data extraction process for full texts was not used. Full-texts were
29
30 examined, looking primarily for correspondence between information reported by the reviewers in
31
32 their Support for Judgement and the text.
33
34
35
36

37 They independently classified each case of disagreement as follows:

- 38 • Disagreement related to differences in interpretation:
 - 39 ○ The support for judgement was the same (or “conceptually” the same) in both
40
41 reviews, but the interpretation differed.
 - 42 ○ One review clearly confused one item of the risk of bias tool with a different
43
44 one or the review authors misunderstood the definition of the item (e.g., for
45
46 random sequence generation, support for judgement reports “600 opaque
47
48 envelopes, 1 was drawn every time”).
49
50
51
52
53
54
55
56
57
58
59
60

- Disagreement related to differences in information: the support for the judgement cites information that is not available in the study report; additional sources are cited (e.g., protocol) or the review authors reported that they had contacted the RCT author for additional data.
- Disagreement related to information missed by the review authors: the study report clearly describes the information, but some review authors seemed to have missed this information in the study report.
- Disagreement related to input mistakes: risk of bias assessment in terms of “high”/“low”/“unclear” did not match the support for the judgement (e.g., “Randomization described explicitly”, judgement “Unclear”).
- Unclear: when it was not possible to classify the disagreement because the support for the judgement was empty or because we could not retrieve the full-text study report.

Any disagreements between reviewers were solved by discussion to reach consensus. In the Supplementary Appendix 1, we report a figure synthesizing how the in-depth analysis process was conducted.

Identification of main reasons for different interpretation

For each disagreement related to a difference in interpretation, we evaluated the probable reason for disagreement. For example, the interpretation could differ because of confusion with another risk of bias item (e.g., random sequence generation and allocation concealment) or because the information was unclear or insufficiently detailed in the article. When we were unsure about the reason, we classified the reason as unclear. Two authors (LB and AD) conducted this process in duplicate by using all available information (i.e., support for the

1
2
3 judgement, characteristics of the study reported in the review, full-text article), with
4
5 disagreements resolved by discussion.
6
7
8
9

10 **Statistical analysis**

11
12 Analysis was descriptive with use of frequencies and percentages for qualitative variables.

13
14 Statistical analysis was conducted with Stata 13.1²⁵. We decided to use simple percent
15
16 agreement because other static approaches were problematic. The Kappa statistic requires
17
18 having defined reviewers, which is not the case of our approach. Another statistic, the
19
20 intraclass correlation coefficient (ICC) is not suitable, because it requires assessments to be in
21
22 an ordinal order, which is not our case. There is no continuum between the assessments of
23
24 low, unclear and high risk of bias.
25
26
27
28
29

30 **Patient involvement**

31
32 Patients were not involved in any aspect of the study design, conduct, or the development of
33
34 the research question or outcome measures. This is a research-on-research study and
35
36 therefore there was no active patient recruitment for data collection.
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

RESULTS

Selection process

Figure 1 shows the selection process. From the 2,796 systematic reviews published between March 2011 and September 2014, 2,291 reviews included RCTs only and reported a risk of bias assessment. Of these, 797 included at least one RCT whose primary reference was shared with another review for which a risk of bias assessment was reported. These 797 reviews included 1,604 single RCTs evaluated for the same risk of bias item in more than one review. The Supplementary Appendix 2 reports the frequency of the different Cochrane groups among those reviews.

Among the 1,604 selected RCTs: 1,603 had duplicate evaluation for allocation concealment, 1,466 for random sequence generation, 375 for blinding of participants and personnel, 583 for blinding of outcome assessment and 1,348 for incomplete outcome data.

Evaluation of agreement and distribution of disagreements

The agreement of risk of bias judgements ranged from 57% (770/1,348 trials) for incomplete outcome data to 81% (1,193/1,466 trials) for random sequence generation (Figure 2). We identified most disagreements for “low” and “unclear” risk of bias judgments, especially for random sequence generation (231/273 trials, 85%). Disagreements between “low” and “high” risk of bias were generally rare, for example 8/273 of disagreements (3%) for random sequence generation, with the exception of incomplete outcome data for which they were more frequent (190/578, 33%). For blinding of participants and personnel, the most frequent disagreement was between “unclear” and “high” risk of bias (50/107, 47%), then “low” versus “unclear” (34/107, 32%), and “low” versus “high” (23/107, 21%) (Figure 2).

Classification of disagreements

The in-depth analysis of disagreements included 802 studies: 799 for allocation concealment, 747 for random sequence generation, 206 for blinding of participants and personnel, 297 for blinding of outcome assessment and 660 for incomplete outcome data. The agreement results of this sample and the distribution of disagreements are reported in the Supplementary Appendix 3.

For all items, the most common source of disagreement was a difference in interpretation, with frequencies ranging from 88/136 (65%) for random sequence generation to 56/62 (90%) for blinding of participants and personnel (Figure 3). The access to additional or different information accounted for disagreements in 32/136 (24%) trials for random sequence generation and 38/205 (19%) for allocation concealment. Access to additional information was less common for the remaining items, with proportions ranging from 2% to 4%. In 80% of the cases, the access to additional information was through the contact of the study author. The other sources of disagreement were less common; input mistake ranged from 1% to 6%, missed information from 1% to 6%. We could not determine the source of disagreement in 5% of our disagreements. For this analysis, we accessed the full text of 216 different trials to help us in the process. The Supplementary Appendix 4 reports some examples of disagreements in which the access to the study report helped us in the classification and the analysis of reasons of disagreement. We could not retrieve or access 19 full-texts we deemed necessary for the categorization of disagreements and this explain the majority of cases where we were unable to categorize the source of disagreement (“unclear” source in Figure 3).

Main reasons of disagreements for different interpretation

1
2
3 The main reasons for a difference in interpretation for each item are reported in Table 1.
4
5 Additional examples are provided for each item for the high-low disagreements
6
7 (Supplementary Appendix 5). The most common reason across items was related to
8
9 incomplete or unclear reporting in the RCT. For random sequence generation, disagreements
10
11 in 73/88 (83%) trials were related to lack of a precise description of the randomization
12
13 process with reviewers evaluating “low”, “high” or “unclear” risk of bias the reporting of
14
15 “randomised” in the text. For allocation concealment, the most common reason for
16
17 disagreement was a different interpretation of description of the envelopes used to conceal
18
19 allocation (17%, n=26/149 trials). For the two blinding items, many disagreements occurred
20
21 when the article mentioned only “double blind” in RCTs without an additional description
22
23 (16% of cases, n=9/56 trials for blinding of participants and personnel, 13%, n=9/70 for
24
25 blinding of outcome assessment). For incomplete outcome data, reviewers assessed
26
27 differently the statement from the study report of “no missing data” or “all data reported”
28
29 (10%, 22/220 trials). Another common reason for a difference in interpretation was confusion
30
31 with another item. Allocation concealment was confused with blinding (10%, n=15/149
32
33 trials) but also with random sequence generation (4%, n=6/149). For blinding of participants
34
35 and personnel, the most common cause for disagreement concerned the interpretation of
36
37 cases when blinding was not feasible (36%, n=20/56 trials), assessed at high risk by some
38
39 reviewers and low by others. Another common cause of disagreement for the two blinding
40
41 items related to the assessment of outcomes that should not be affected by blinding (e.g.,
42
43 mortality); it explained 21% (n=12 trials) of disagreements for blinding of participants and
44
45 personnel and 23% (n=16 trials) for blinding of outcome assessment, often low versus high
46
47 disagreements.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 For incomplete outcome data, the use of different cut-offs for the rate of missing data is the
4 most common reason for disagreement (26%, n=57 trials); also common is considering the
5 explanation of reasons for missing data enough to attribute a low risk of bias (13%. n=28
6 trials).
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

DISCUSSION

In this study, we took advantage of a very large sample of Cochrane reviews to explore the sources of disagreements in risk of bias assessment for trials included in several reviews. We decided to focus on Cochrane reviews because as these reviews are produced within a single organization, therefore we expected results and procedures to be more appropriately comparable. Authors compiling Cochrane reviews are members of the organization and, in most cases, they underwent a similar training for assessing risk of bias. Our results confirm that the agreement for risk of bias assessments is generally suboptimal, with better agreement for random sequence generation and allocation concealment and less agreement for incomplete outcome data. Access to different sources of information explained why 24% of the trials had disagreements in the assessment of risk of bias for random sequence generation and 19% for allocation concealment. However, the main source of disagreements was a difference in interpretation of the same information, which was frequently related to incomplete or unclear reporting in the study report.

Strengths and weaknesses

Our study goes beyond previous literature on the topic^{3 12-18 26}. As compared with most other studies¹²⁻¹⁷ we used real-world data to explore agreement of risk of bias assessments in real scenarios. We evaluated a very large and comprehensive collection of Cochrane reviews that spanned multiple specialties and topics, including a number of trials about ten times larger than the largest study on the topic¹². We completed our analysis by searching individual study reports to give support to our comments on reasons for disagreements, which, to our knowledge, has not been done in previous, smaller works that used a similar methodology¹⁸.

1
2
3 While doing this, we developed a suitable classification scheme for sources of disagreements
4 and conducted, in duplicate, an extensive analysis to understand the risk of bias assessment
5 process and explored the most common reasons for disagreements.
6
7

8
9
10 Our study has limitations. Whenever a single RCT was included in three reviews or more, we
11 considered only the risk of bias assessment from two reviews chosen at random.
12

13
14 Nevertheless, we cannot exclude that different combinations of two chosen evaluations could
15 have produced slightly different results. Although the classification of disagreements was
16 conducted in duplicate following a formalised process, there remains a component of
17 personal judgement. We evaluated only disagreements, but a number of agreements might
18 have occurred “by chance”. In our analysis of likely reasons for disagreements, some resulted
19 from confusion between risk of bias items. Similar discrepancies might have occurred among
20 agreements; indeed, previous literature on the topic demonstrated that reviewers do not
21 accurately follow the risk of bias tool²⁷. We also did not assess the selective reporting item
22 that is frequently judged on incomplete information. We did not evaluate whether
23 disagreements varied depending on the Cochrane review group or year of publication.
24
25 Finally, we did not evaluate the impact of disagreements and the extent to which the evidence
26 base for making conclusions and providing summary statements of effectiveness may have
27 been affected by changing the rating.
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

47 **Comparison with other studies**

48
49 Our findings confirm the importance of issues that were previously identified by Jorgensen et
50 al.³ and Savovic et al.²⁶. In particular, Savovic et al.²⁶, surveying users of the risk of bias tool,
51 reported on the possibility of confusion between random sequence generation and allocation
52 concealment and between allocation concealment and blinding; the uncertainty on how to
53
54
55
56
57
58
59
60

1
2
3 address unfeasibility of blinding; and the difficulties in assessing incomplete outcome data
4 especially regarding the acceptable rate of missing data. More recently, Jorgensen et al.³,
5
6 evaluating comments on the use of the risk of bias tool, highlighted how authors complained
7
8 that judgment often originates from incomplete or missing information.
9
10

11
12 A previous study identified 46 RCTs included in different systematic reviews in the field of
13
14 fertility and evaluated the percentage agreement in risk of bias assessment. That analysis
15
16 showed generally worse agreement than in our study, with percentage agreement ranging
17
18 from 35% to 71%. Differences in sample size and the particular topic may explain these
19
20 differences. In addition, although the authors had compared supports for judgement between
21
22 reviews, this evaluation may have been incomplete, because they did not evaluate the primary
23
24 study reports¹⁸.
25
26
27
28
29

30 31 **Implications**

32
33 Our results confirm that the agreement in risk of bias assessment would be enhanced by more
34
35 detailed guidance in use of the risk of bias tool with particular focus on common causes of
36
37 disagreements. We showed that in many cases, the unclear reporting from source material
38
39 allows reviewers ample space for personal judgement and differences in judgement.
40
41

42
43 The scientific community continues to stress the importance of improving the reporting of
44
45 trials²⁸⁻³¹, which may limit disagreements when assessing risk of bias. In parallel, we could
46
47 also work on restricting the space for personal interpretation when assessing risk of bias. A
48
49 suggestion could be to give clearer instruction on how to evaluate common cases, for
50
51 example when confronted with nothing more than the term “randomised” or “double blind”
52
53 in the study report. Similarly, a threshold could be set on the quota for missing data and
54
55 indications on which imputation methods are appropriate and in which situations.
56
57
58
59
60

1
2
3 To minimise research waste, it could be interesting to have access to risk of bias assessments
4 from other Cochrane groups and the supports they used, including information from authors
5 or from protocols to help reviewers in their assessments. This process would imply having a
6 unique study identification number across reviews and a central shared repository for all
7 studies included in any Cochrane reviews.
8
9

10
11
12
13
14 Following the suggestions based on the findings and comments provided by Jorgensen et al.³
15 and Savovic et al.²⁶, Cochrane has been working on a new version of the risk of bias tool,
16 which has recently been released^{32 33}. The new version has a different approach to the risk of
17 bias assessment, guiding reviewers through the process with the use of “signalling questions”,
18 which might leave less room for subjectivity. In addition, there is more guidance in assessing
19 some items. For example, the new tool better clarifies some aspects of the randomization
20 process, especially about what to do in some cases of incomplete information (e.g.,
21 randomization list created by an external centre with no other indication). The new tool also
22 has a different approach to the blinding aspect, oriented to the implications of the masking
23 process. However, the new tool does not cover some of our concerns, especially those related
24 to incomplete outcome data: quota for missing data that are considered acceptable, and
25 whether reviewers should focus more on the reasons for the missing data or their magnitude.
26 It also does not address the common case of authors reporting “no missing data”. Research-
27 on-research studies are needed to evaluate whether this new version of the tool results in
28 improved reproducibility.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 **Conclusion**

52
53 This analysis of risk of bias assessment for more than 1,600 trials included in more than one
54 reviews showed that agreement remains suboptimal. Most disagreements come from a
55
56
57
58
59
60

1
2
3 difference in interpretation of an incomplete or unclear description in the study report. In
4
5 some cases, the difference in the assessment was due to some but not all review authors
6
7 obtaining additional information, from a protocol or from contacting study author.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

REFERENCES

1. Cook DJ, Mulrow CD, Haynes RB. Systematic reviews: synthesis of best evidence for clinical decisions. *Ann Intern Med* 1997;126(5):376-80. [published Online First: 1997/03/01]
2. Hopewell S, Boutron I, Altman DG, et al. Incorporation of assessments of risk of bias of primary studies in systematic reviews of randomised trials: a cross-sectional study. *BMJ Open* 2013;3(8):e003342. doi: 10.1136/bmjopen-2013-003342 [published Online First: 2013/08/27]
3. Jorgensen L, Paludan-Muller AS, Laursen DR, et al. Evaluation of the Cochrane tool for assessing risk of bias in randomized clinical trials: overview of published comments and analysis of user practice in Cochrane and non-Cochrane reviews. *Syst Rev* 2016;5:80. doi: 10.1186/s13643-016-0259-8 [published Online First: 2016/05/11]
4. Hrobjartsson A, Boutron I, Turner L, et al. Assessing risk of bias in randomised clinical trials included in Cochrane Reviews: the why is easy, the how is a challenge. *Cochrane Database Syst Rev* 2013(4):ED000058. doi: 10.1002/14651858.ED000058 [published Online First: 2013/06/04]
5. Higgins JP, Altman DG, Gotzsche PC, et al. The Cochrane Collaboration's tool for assessing risk of bias in randomised trials. *BMJ* 2011;343(oct18 2):d5928. doi: 10.1136/bmj.d5928
6. Schulz KF, Chalmers I, Hayes RJ, et al. Empirical evidence of bias. Dimensions of methodological quality associated with estimates of treatment effects in controlled trials. *JAMA* 1995;273(5):408-12. [published Online First: 1995/02/01]
7. Moher D, Pham B, Jones A, et al. Does quality of reports of randomised trials affect estimates of intervention efficacy reported in meta-analyses? *The Lancet* 1998;352(9128):609-13. doi: 10.1016/s0140-6736(98)01085-x
8. Balk EM, Bonis PA, Moskowitz H, et al. Correlation of quality measures with estimates of treatment effect in meta-analyses of randomized controlled trials. *JAMA* 2002;287(22):2973-82. [published Online First: 2002/06/08]
9. Page MJ, Higgins JP, Clayton G, et al. Empirical Evidence of Study Design Biases in Randomized Trials: Systematic Review of Meta-Epidemiological Studies. *PLoS One* 2016;11(7):e0159267. doi: 10.1371/journal.pone.0159267 [published Online First: 2016/07/12]
10. Dechartres A, Trinquart L, Faber T, et al. Empirical evaluation of which trial characteristics are associated with treatment effect estimates. *J Clin Epidemiol* 2016;77:24-37. doi: 10.1016/j.jclinepi.2016.04.005
11. Higgins JPT, Altman DG, Sterne JAC (editors). Chapter 8: Assessing risk of bias in included studies. In: Higgins JP, Green S, eds. *Cochrane Handbook for Systematic Reviews of Interventions*. Version 5.1.0 (updated March 2011) ed: The Cochrane Collaboration, 2011. Available from www.handbook.cochrane.org.
12. Hartling L, Ospina M, Liang Y, et al. Risk of bias versus quality assessment of randomised controlled trials: cross sectional study. *BMJ* 2009;339(oct19 1):b4012. doi: 10.1136/bmj.b4012
13. Hartling L, Bond K, Vandermeer B, et al. Applying the risk of bias tool in a systematic review of combination long-acting beta-agonists and inhaled corticosteroids for persistent asthma. *PLoS One* 2011;6(2):e17242. doi: 10.1371/journal.pone.0017242
14. Hartling L, Hamm MP, Milne A, et al. Testing the risk of bias tool showed low reliability between individual reviewers and across consensus assessments of reviewer pairs. *J Clin Epidemiol* 2013;66(9):973-81. doi: 10.1016/j.jclinepi.2012.07.005
15. Armijo-Olivo S, Ospina M, da Costa BR, et al. Poor reliability between Cochrane reviewers and blinded external reviewers when applying the Cochrane risk of bias tool in physical therapy trials. *PLoS One* 2014;9(5):e96920. doi: 10.1371/journal.pone.0096920 [published Online First: 2014/05/16]
16. Armijo-Olivo S, Stiles CR, Hagen NA, et al. Assessment of study quality for systematic reviews: a comparison of the Cochrane Collaboration Risk of Bias Tool and the Effective Public Health Practice Project Quality Assessment Tool: methodological research. *J Eval Clin Pract* 2012;18(1):12-8. doi: 10.1111/j.1365-2753.2010.01516.x
17. da Costa BR, Beckett B, Diaz A, et al. Effect of standardized training on the reliability of the Cochrane risk of bias assessment tool: a prospective study. *Syst Rev* 2017;6(1):44. doi: 10.1186/s13643-017-0441-7

18. Jordan VM, Lensen SF, Farquhar CM. There were large discrepancies in risk of bias tool judgments when a randomized controlled trial appeared in more than one systematic review. *J Clin Epidemiol* 2017;81:72-76. doi: 10.1016/j.jclinepi.2016.08.012 [published Online First: 2016/09/14]
19. Wilkins AJ. Risk of bias in assessing Risk of Bias. *Ophthalmic Physiol Opt* 2017;37(1):107-09. doi: 10.1111/opo.12333
20. Review Manager (RevMan) [Computer program] [program]. Version 5.3 version. Copenhagen: The Nordic Cochrane Centre, The Cochrane Collaboration, 2014.
21. R: A language and environment for statistical computing. [program]. Viena: R Foundation for Statistical Computing, 2013.
22. XML: Tools for Parsing and Generating XML Within R and S-Plus. [program]. 3.2.2 version, 2017.
23. Dechartres A, Trinquart L, Atal I, et al. Evolution of poor reporting and inadequate methods over time in 20 920 randomised controlled trials included in Cochrane reviews: research on research study. *BMJ* 2017;357:j2490. doi: 10.1136/bmj.j2490
24. deMelo, VV, Conference: Advances in Logic Based Intelligent Systems; 2005.
25. Stata Statistical Software: Release 13 [program]. College Station, TX: StataCorp LP, 2013.
26. Savovic J, Weeks L, Sterne JA, et al. Evaluation of the Cochrane Collaboration's tool for assessing the risk of bias in randomized trials: focus groups, online survey, proposed recommendations and their implementation. *Syst Rev* 2014;3:37. doi: 10.1186/2046-4053-3-37 [published Online First: 2014/04/16]
27. Propadalo I, Tranfic M, Vuka I, et al. In Cochrane reviews risk of bias assessments for allocation concealment were frequently not in line with Cochrane's Handbook guidance. *J Clin Epidemiol* 2018 doi: 10.1016/j.jclinepi.2018.10.002
28. Shamseer L, Hopewell S, Altman DG, et al. Update on the endorsement of CONSORT by high impact factor journals: a survey of journal "Instructions to Authors" in 2014. *Trials* 2016;17(1):301. doi: 10.1186/s13063-016-1408-z
29. Turner L, Shamseer L, Altman DG, et al. Does use of the CONSORT Statement impact the completeness of reporting of randomised controlled trials published in medical journals? A Cochrane review. *Syst Rev* 2012;1:60. doi: 10.1186/2046-4053-1-60
30. Turner L, Shamseer L, Altman DG, et al. Consolidated standards of reporting trials (CONSORT) and the completeness of reporting of randomised controlled trials (RCTs) published in medical journals. *Cochrane Database Syst Rev* 2012;11:MR000030. doi: 10.1002/14651858.MR000030.pub2
31. Altman DG, Moher D, Schulz KF. Improving the reporting of randomised trials: the CONSORT Statement and beyond. *Stat Med* 2012;31(25):2985-97. doi: 10.1002/sim.5402
32. RoB 2.0 Tool [Available from: <https://sites.google.com/site/riskofbiastool/welcome/rob-2-0-tool> accessed 22/08/2018.
33. Higgins JP, Sterne JA, Savovic J, et al. A revised tool for assessing risk of bias in randomized trials. *Cochrane Database of Systematic Reviews* 2016(10 (Suppl 1)) doi: dx.doi.org/10.1002/14651858.CD201601

ACKNOWLEDGEMENTS:

We thank Camila Olarte Parra for her help during the data management phase and her comments on this manuscript. We thank David Tovey, editor in chief of the Cochrane Library, for agreeing to share data from Cochrane reviews; Javier Mayoral Campos, system administrator; the Cochrane Central Executive for preparing files; and all Cochrane reviewers who collected data. We also thank Laura Smales for English revision of the manuscript.

CONTRIBUTIONS:

Lorenzo Bertizzolo was involved in the study conception, selection of trials, data extraction, data analysis, interpretation of results and drafting the manuscript.

Patrick Bossuyt was involved in the study conception, data analysis, interpretation of results and drafting the manuscript.

Ignacio Atal was involved in the study conception, data extraction, and drafting the manuscript.

Philippe Ravaud was involved in the study conception and drafting the manuscript.

Agnès Dechartres was involved in the study conception, selection of trials, data extraction, data analysis, interpretation of results and drafting the manuscript.

Lorenzo Bertizzolo is the guarantor. He had full access to all the data in the study and takes responsibility for the integrity of the data and the accuracy of the data analysis.

FUNDING:

This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 676207.

COMPETING INTERESTS:

1 The authors declare that they have no competing interests in relation to this study.
2
3

4
5 **DATA SHARING:**
6

7 Raw data and analyses are available on request from the authors.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

List of tables

Table 1: Main reasons for disagreements in cases of a different interpretation of the same information.

List of figures

Figure 1: Flow-chart of the selection process.

Figure 2: Distribution of agreements and disagreements for the different risk of bias items analysed; raw number and percentages of the total. For disagreements, distribution of the different discrepancies.

Figure 3: Classification of disagreements for the different items; raw number and percentage of the total.

List of appendices

Supplementary Appendix 1: explanatory figure of the categorization process for the in-depth analysis of disagreement.

Supplementary Appendix 2: Frequency of the different Cochrane review groups involved in the included reviews.

Supplementary Appendix 3: for the 50% selection of studies for the in-depth analysis; distribution of agreements and disagreements for the different risk of bias items analysed; raw number and percentages of the total. For disagreements, distribution of the different discrepancies.

Supplementary Appendix 4: selected examples where the access to the study report helped us in the categorization of the disagreement and in highlighting the reasons for disagreement.

Supplementary Appendix 5 Reasons for disagreements in cases of different interpretation of the same information; focus on “low” versus “high” disagreements

Table 1| Main reasons for disagreements in cases of a different interpretation of the same information.

Risk of bias	Main reasons for disagreements	N (%)⁺	Examples of support for judgement from the review*
Item			
random sequence generation	Consider differently incomplete or unclear description	73 (83)	“States “cluster randomisation by computer””; Low risk of bias “Cluster randomisation by computer. No further information provided”; Unclear risk of bias
	Confusion with allocation concealment	9 (10)	“allocation was done using sealed envelopes containing name of one of the two groups.”; Low risk of bias
allocation concealment	Consider differently incomplete or unclear description	49 (33)	“Not specified.”; High risk of bias “Method of concealment not described.”; Unclear risk of bias
	Consider differently envelopes description	26 (17)	““Sequentially numbered sealed envelopes”. Does not state if opaque envelopes.”; Unclear risk of bias “Sequentially numbered sealed envelopes.”; Low risk of bias
	Random sequence generated by computer or external centre considered enough for Low risk	21 (14)	“Treatment was allocated based on the computer-generated number list.”; Low risk of bias
	Confusion in the definition of the item	19 (13)	“Researchers attempted to contact all patients seen by physicians during one month”; High risk of bias
	Confusion with blinding	15 (10)	“participants were told to which compound they had been allocated.”; High risk of bias
	Confusion with random sequence generation	6 (4)	“Computer generated randomised lists.”; Low risk of bias

⁺ Number of RCTs disagreeing for this reason; percentage over the total of disagreements for different interpretation.

* When two extracts are reported, they refer to the same study.

<i>Risk of bias Item</i>	<i>Main reasons for disagreements</i>	<i>N (%)⁺</i>	<i>Examples of support for judgement from the review*</i>
<i>blinding of participants and personnel</i>	Assess risk differently if blinding was not feasible because of the type of intervention	20 (36)	“Not possible to blind participants”; Low risk of bias “Participants were not blinded for provided treatment. This is inherent to study design”; High risk of bias
	outcome considered not influenced by blinding	12 (21)	“No information given about whether patients were blind to physician allocation but treatment outcomes judged unlikely to be affected by lack of blinding”; Low risk of bias
	Consider differently information of “double blind”	9 (16)	“Quote: “. . . patients were randomised in double-blind conditions . . . ”Comment: probably done”; Low risk of bias “Quote: “double blind conditions”. No further details.”; Unclear risk of bias
	Consider differently incomplete or unclear description	7 (12)	“Researchers were blind until after the baseline assessment. participants were not blinded.”; Unclear risk of bias “Not possible to blind participants to intervention. Insufficient information to make a judgement about blinding of therapists”; High risk of bias
	Confusion in the definition of the item	5 (9)	“Described as an “open-label” pilot study.”; Low risk of bias
<i>blinding of outcome assessment</i>	Consider differently incomplete or unclear description	24 (34)	“Not explicitly discussed in the publish study, it was assumed to be open label”; High risk of bias “Not described in published study”; Unclear risk of bias
	outcome considered not influenced by blinding	16 (23)	“Not stated, but it was unlikely that the outcome was influenced by lack of blinding”; Low risk of bias
	Consider differently patient-reported outcomes when patients are blinded or not to the intervention	9 (13)	“Comment: depression assessed by patient self-report”; High risk of bias “Insufficient information available to assess”; Low risk of bias
	Consider differently information of “double blind”	9 (13)	“Quote: “. . . double blind” Comment: probably done”; Low risk of bias “Quote: “double blind conditions”. No further details.”; Unclear risk of bias
	Assess risk differently if blinding was not feasible because of the type of intervention	6 (9)	“blinding not possible due to intervention”; High risk of bias “Unclear blinding of outcome assessment”; Low risk of bias

⁺ Number of RCTs disagreeing for this reason; percentage over the total of disagreements for different interpretation.

^{*} When two extracts are reported, they refer to the same study.

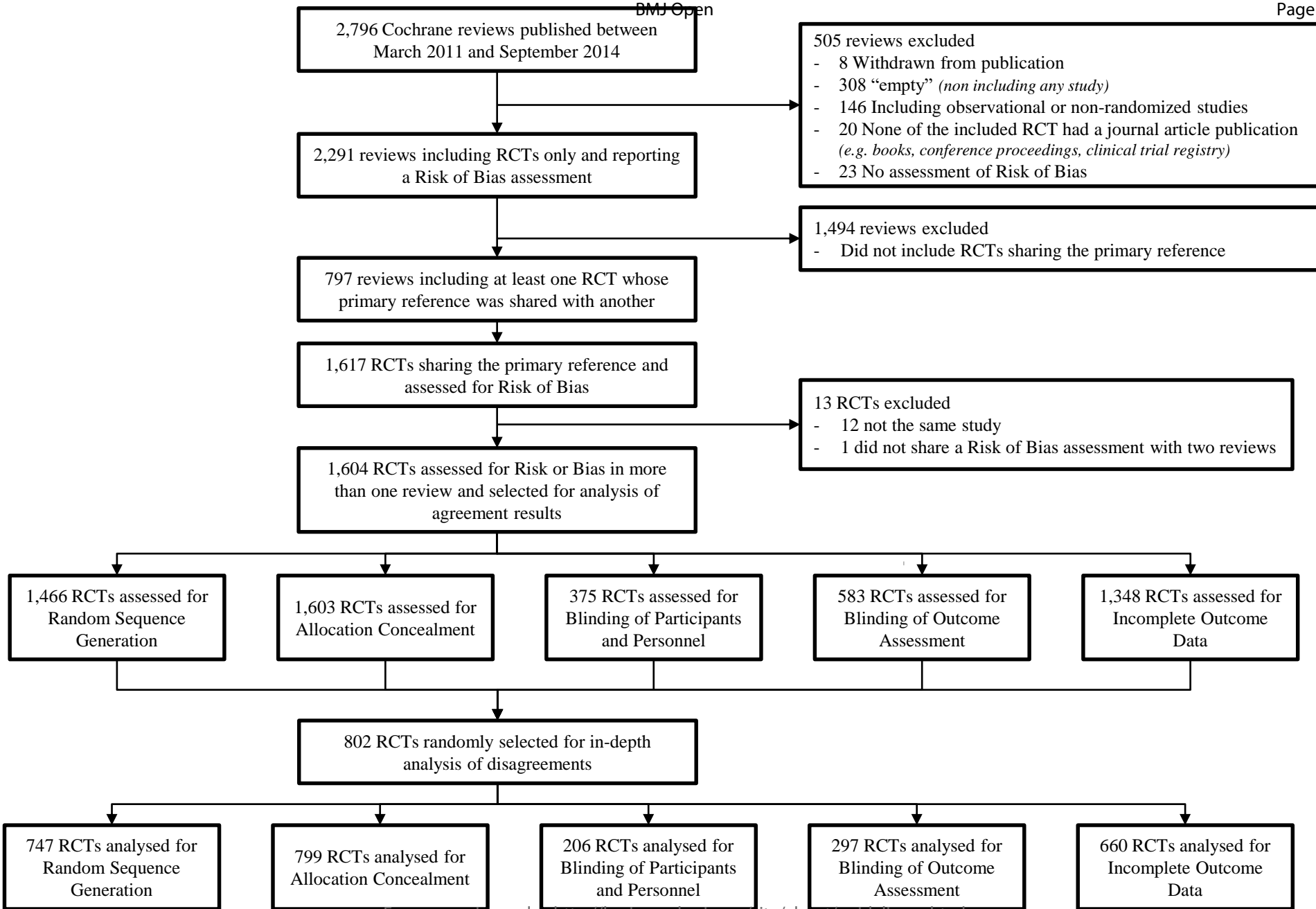
<i>Risk of bias Item</i>	Main reasons for disagreements	N (%)⁺	Examples of support for judgement from the review*
<i>incomplete outcome data</i>	Use different cut-off for the rate of missing data	57 (26)	“11 withdrawals (10%).”; Low risk of bias “Comment: there were post-randomisation drop-outs”; High risk of bias
	Focus on number vs reasons/precise report of missing data	28 (13)	“20 drop-outs (27.2%) with 4 deaths (3 males, 1 female) from cardiovascular events”; High risk of bias “Numbers and reasons for dropouts and withdrawals in all intervention groups were described.”; Low risk of bias
	Consider differently incomplete or unclear description	27 (12)	“Women who were untraceable or unsuitable for follow-up were excluded, other losses included as smokers”; Low risk of bias “167/1287 (12.9%) (C = 83, I = 84) excluded from analysis due to moving away, being untraceable or deemed unsuitable for follow-up (e.g. miscarriage). 1120 in sample. 51/1287 non-responders were included as continuing smokers.” High risk of bias
	Consider differently intention-to-treat analysis	25 (11)	“147 randomised; 4 in the letrozole group and 3 in the LOD dropped out of the trial, all for non-compliance. However, ITT analysis was not conducted.”; Unclear risk of bias “7 women lost to follow up, but similar (3 vs 4) in both groups; losses due to noncompliance”; Low risk of bias
	Consider differently report of “no missing data”	22 (10)	“Did not report number of withdrawals. Comment: all patients who were randomised were included in the final analysis. ITT analysis was conducted.”; Unclear risk of bias “It does not appear that there were any withdrawals or dropouts” Low risk of bias
	Consider differently imputation of missing data	20 (9)	“Imputation method not described”; Unclear risk of bias “Dropout rate was not significant”; Low risk of bias
	Use different cut-off for difference in the rate missing data between different arms/comparisons	13 (6)	“Dropout higher in placebo group (35% vs 25% in budesonide group). ITT used.”; High risk of bias “Similar rates of withdrawal between arms. Withdrawals: 36 BUD, 51 placebo”; Low risk of bias

⁺ Number of RCTs disagreeing for this reason; percentage over the total of disagreements for different interpretation.

* When two extracts are reported, they refer to the same study.

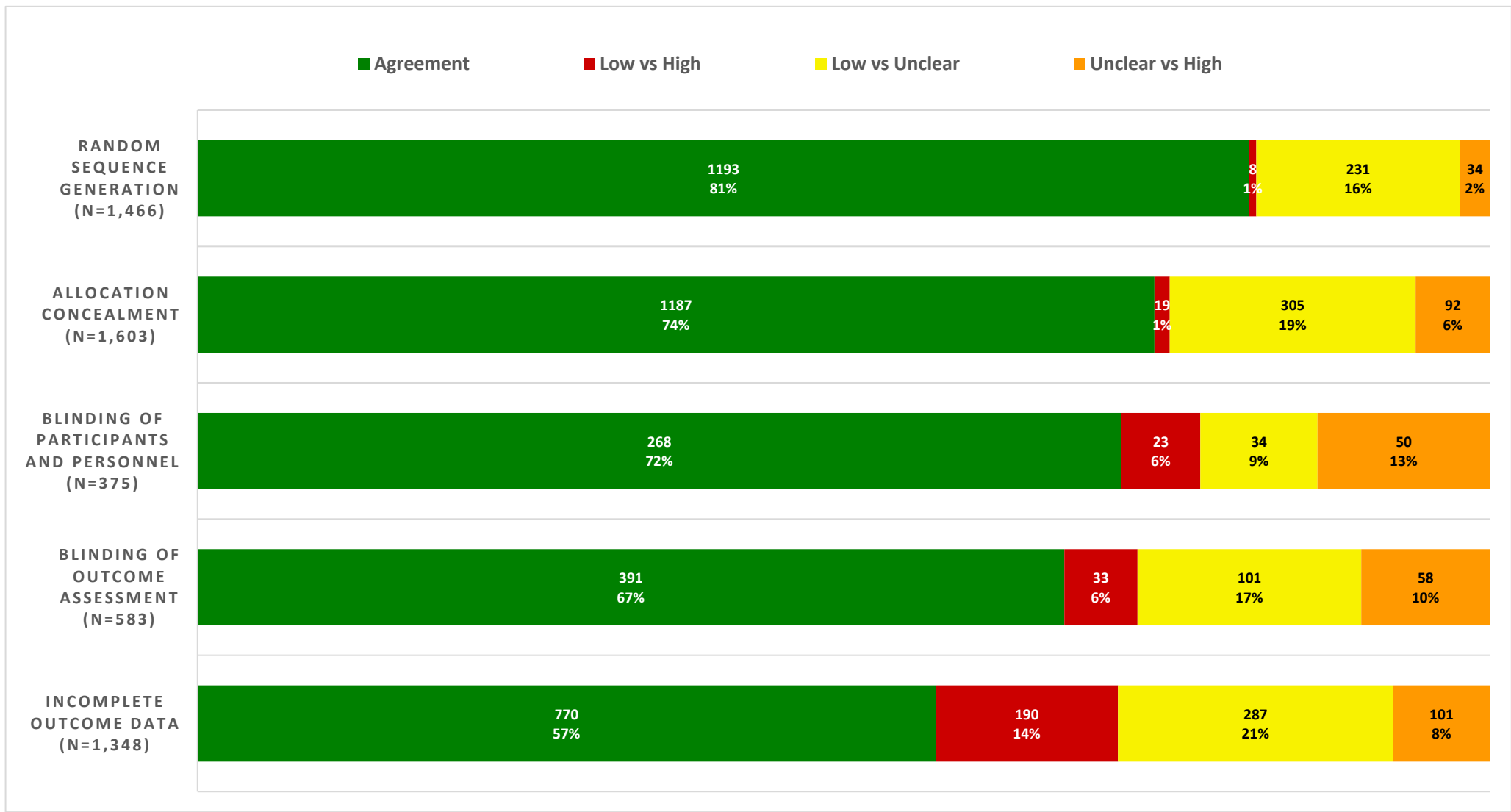
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

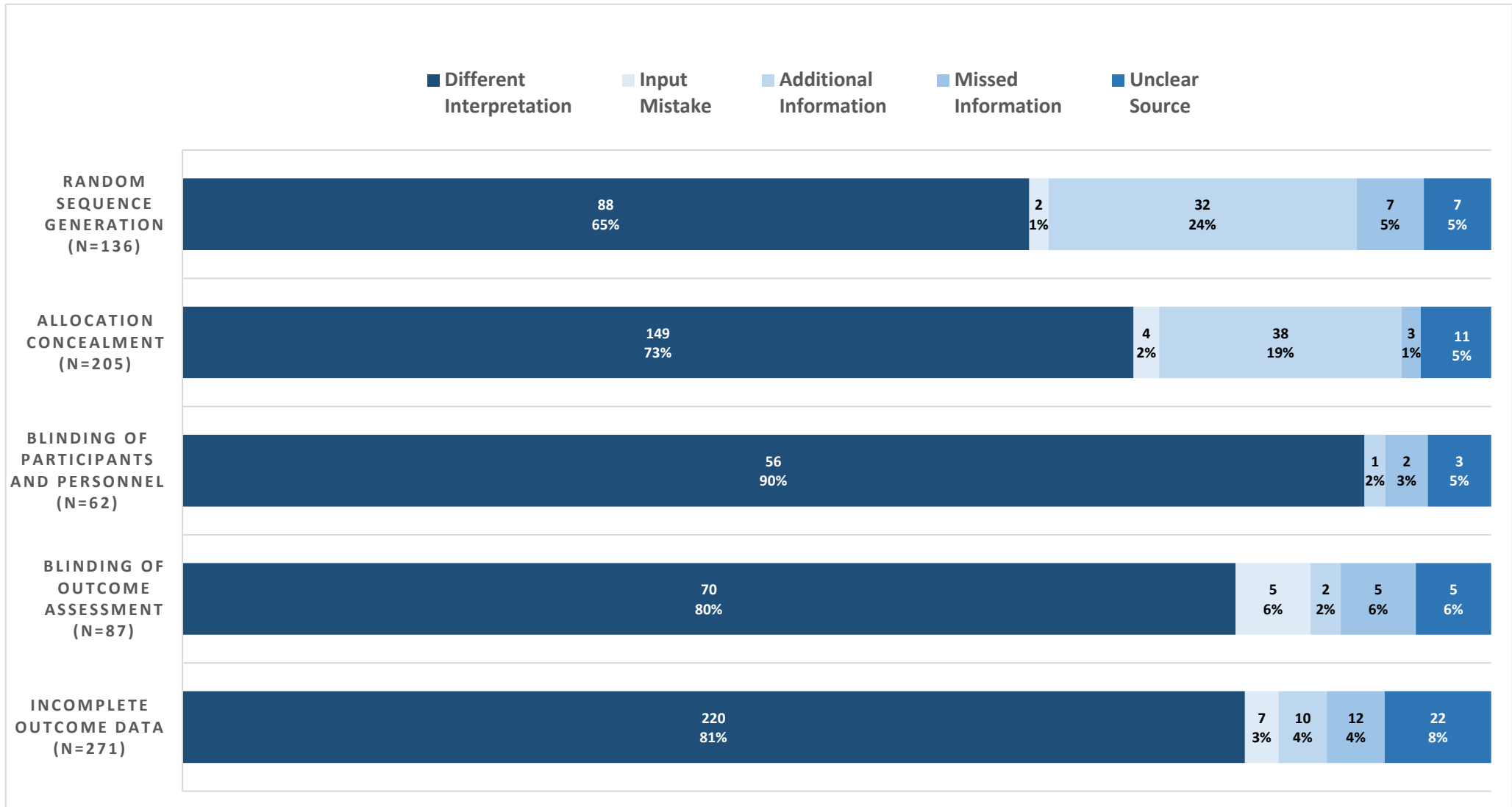
For peer review only



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46





1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Same support for judgement

Risk of bias assessment does not match the support for judgement
(e.g. "Randomization described explicitly", judgement "Unclear")

Input mistake

One review confuses one item with a different one / misunderstanding of definition of the item
(e.g. for Random Sequence Generation "600 opaque envelopes, I was drawn every time")

Differences in interpretation

Assess support for judgement in the two different reviews

Different support for judgement

Access to the study report

Information in the report is incomplete or unclear

Study report clearly describes the information, but one review seemed to have missed it

Missed information from the study report

Study report does not describe the information reported by one author

No access to study report

Mention of access to additional information in the review

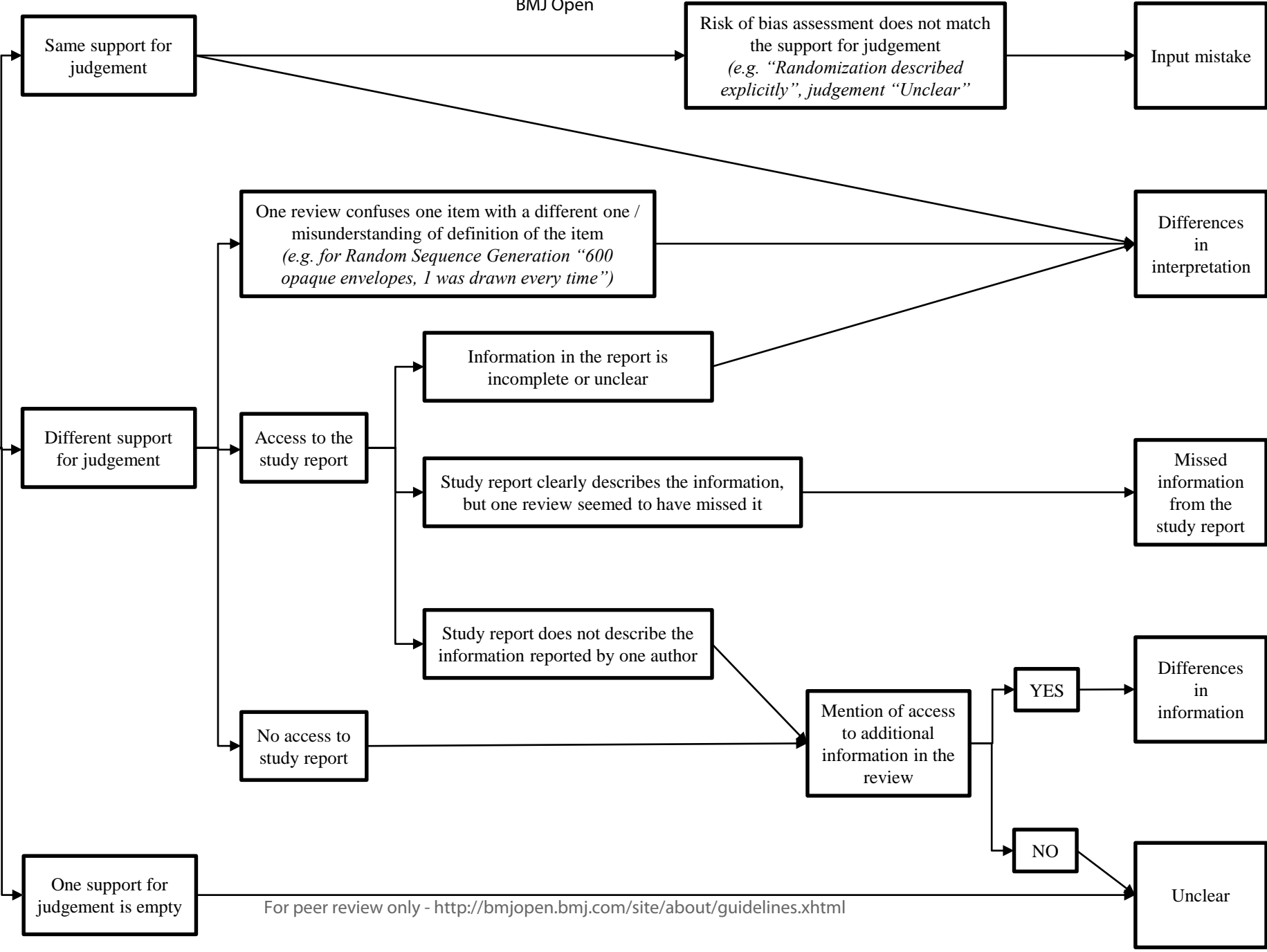
YES

Differences in information

NO

Unclear

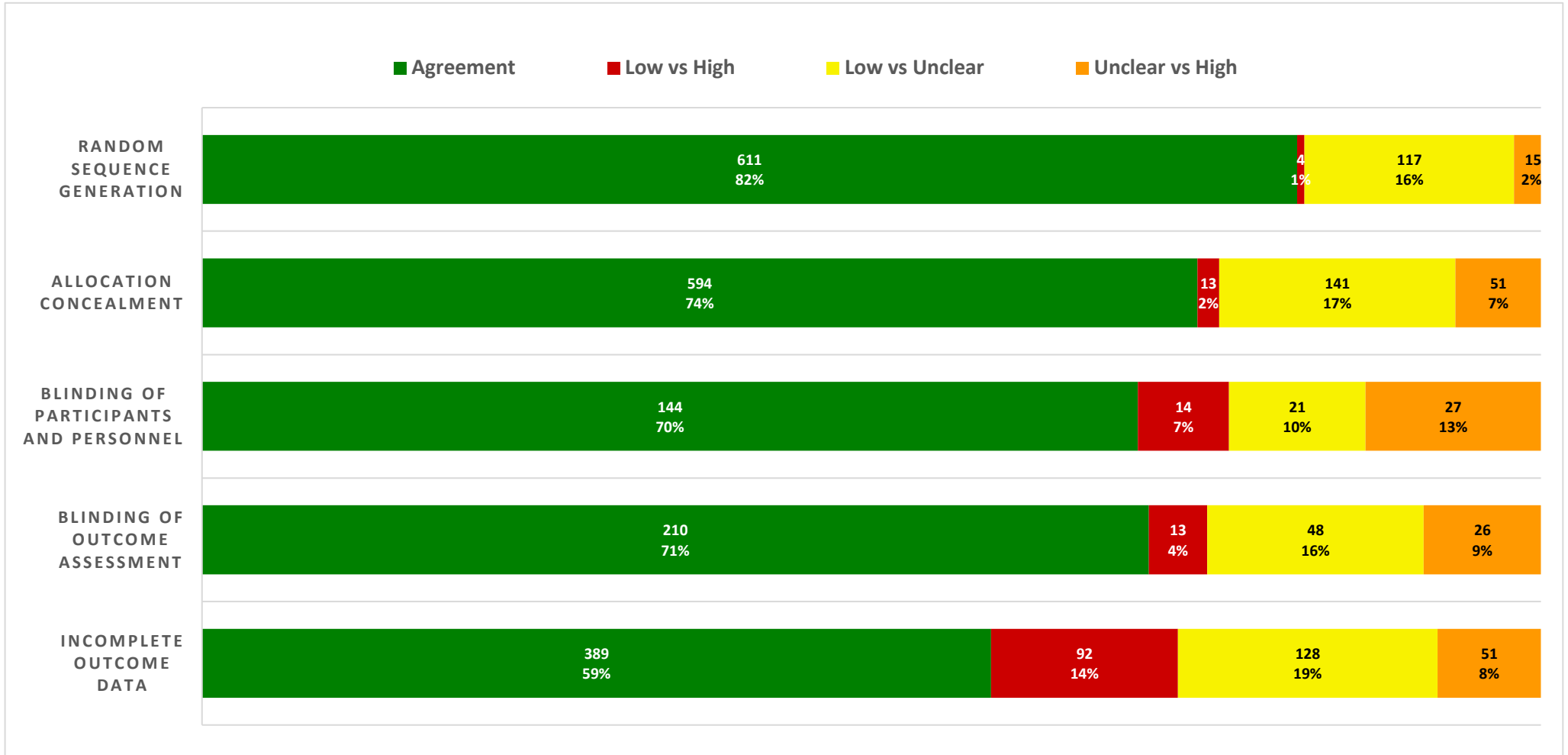
One support for judgement is empty



	Cochrane Group	Number of reviews	% on the total
	Pregnancy and Childbirth	93	11.7%
	Airways	48	6.0%
	Pain, Palliative and Supportive Care Group	42	5.3%
	Acute Respiratory Infections	37	4.6%
	Gynaecology and Fertility	29	3.6%
	Neonatal	29	3.6%
	Tobacco Addiction	27	3.4%
	Stroke	25	3.1%
	Gynaecological, Neuro-oncology and Orphan Cancer Group	23	2.9%
	Wounds	23	2.9%
	Hepato-Biliary	22	2.8%
	Cystic Fibrosis and Genetic Disorders	21	2.6%
	Anaesthesia	20	2.5%
	Drugs and Alcohol	20	2.5%
	Neuromuscular	19	2.4%
	Common Mental Disorders	18	2.3%
	Fertility Regulation	18	2.3%
	Heart	17	2.1%
	Developmental, Psychosocial and Learning Problems	16	2.0%
	Incontinence	16	2.0%
	Kidney disease	16	2.0%
	Schizophrenia	16	2.0%
	Infectious Diseases	15	1.9%
	Oral Health	14	1.8%
	Vascular	14	1.8%
	Dementia and Cognitive Improvement	13	1.6%
	Musculoskeletal	12	1.5%
	Consumers and Communication	10	1.3%
	Epilepsy	10	1.3%
	Eyes and Vision	9	1.1%
	Metabolic and Endocrine Disorders	9	1.1%
	Back and Neck	8	1.0%
	Hypertension	8	1.0%
	Multiple Sclerosis	8	1.0%
	Effective Practice and Organisation of Care	7	0.9%
	HIV/AIDS	7	0.9%
	Inflammatory Bowel Disease	7	0.9%
	Injuries	7	0.9%
	Bone, Joint and Muscle Trauma Group	6	0.8%
	ENT	6	0.8%
	Haematological Malignancies	6	0.8%

Cochrane Group	Number of reviews	% on the total
Breast Cancer	5	0.6%
Colorectal Cancer	5	0.6%
Lung Cancer	3	0.4%
Movement Disorders	3	0.4%
Skin	3	0.4%
Occupational Health	2	0.3%
Sexually Transmitted Infections	2	0.3%
Public Health	1	0.1%
Upper GI and Pancreatic Diseases	1	0.1%
Urology	1	0.1%
Total	797	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46



Supplementary Appendix 4| Examples of in-depth analysis of disagreements conducted with the support of the study report

Risk of bias item	Study Name	Support for judgement*	Information in the study report^	Category of disagreement	Reason of disagreement
Random sequence generation	ABCD 2004	Review 4136: Generated the randomisation list using SAS, stratified by sex and SCr; Low Risk	<i>“The (...) statistician generated the randomization list using SAS (...) stratified by sex and baseline serum creatinine concentration (...).”</i>	Missed information from the study report	
		Review 8277: Method not reported; Unclear Risk			
	Cho 2006	Review 7566: Stated that it is a quasi-randomized study but details not given; High Risk	<i>“... using a quasi-experimental design with a non-equivalent control group.”</i> <i>“They were randomly assigned to participate in the experimental group (...) or a waiting-list control group (...).”</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 9553: Participants randomly allocated to treatment or control group; Unclear Risk			
	Petersen 2005	Review 9132: Quote: “[P]atients were randomly assigned...”Quote: “We used an adaptive allocation scheme for the treatment assignment, with the MMSE score, age and APOE e4 status as balancing covariates”; Low Risk	<i>“We used an adaptive allocation scheme for the treatment assignment, with the MMSE score, age, and APOE e 4 status as balancing covariates.”</i>	Different interpretation	Confusion or misknowledge
		Review 7176: The trial is described as randomised, but the method of sequence generation was not specified. Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report^	Category of disagreement	Reason of disagreement
Allocation concealment	Burge 2000	Review 2991: Participants were randomly assigned sequentially from a list comprising treatment numbers only; Low Risk	<i>"We used a computer generated allocation schedule stratified by centre (block size of six). Patients were randomised sequentially from a list comprising treatment numbers only".</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 10115: Information not available; Unclear Risk			
	McMurdo 1993	Review 4294: Quote: "Randomisation was by opening sealed envelopes supplied in sequence by the study co-ordinator; Low Risk	<i>"Randomization was by opening sealed envelopes supplied in sequence by the study co-ordinator (...), and prepared from a computer-generated random numbers table."</i>	Different interpretation	Consider differently envelopes description
		Review 4963: Unclear, insufficient reporting to permit judgement; Unclear Risk			
	Draper 2007	Review 8179: "... and alternating between treatment or wait list control groups."; High Risk	<i>"On each occasion that a least eight patients had been recruited, their names were selected at random by a blinded investigator to be allocated alternately to the immediate treatment group or a wait-list control group."</i>	Different interpretation	Consider differently incomplete or unclear description
		Review 1919: "Reported as concealed but specific method for concealment not reported"; Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report^	Category of disagreement	Reason of disagreement
Blinding of participants and personnel	Nielsen 2006	Review 9672: "Double-blind"; Low Risk	<p><i>"This study was a randomized, placebo-controlled, double blind, Danish, multi-center (two centers) study."</i></p> <p><i>"The treatment was applied by a nasal spray with one puff in each nostril every day either in the morning or evening."</i></p>	Different information	One review accessed additional data through another study report
		Review 4143: Although "All treatments were supplied as identical intranasal sprays..." the 2004 publication describes a higher rate of withdrawal due to adverse effects in the intervention groups [11.7% in the placebo group, 21.7% in the 150 gm group and 28.7% in the 300 gm group} which may have affected blinding status; Unclear Risk			
	Gersel 1979	Review 10562: Described as double-blind [presumed participants and personnel/investigators]; Low Risk	<p><i>"A double-blind experimental design was used, employing each patient as his own control."</i></p>	Different interpretation	Consider differently information of "double blind"
		Review 6968: Not mentioned and no information to suggest this was done.; Unclear Risk			
	Stein 2011	Review 7025: Not possible to blind participants to intervention. Insufficient information to make a judgement about blinding of therapists; High Risk	<p><i>"Follow-up assessment was made 3 months after release (research staff conducting assessments were blind to treatment assignment)."</i></p> <p><i>"Randomization was accomplished via random numbers table in advance and placed in an envelope by the project coordinator. Following baseline assessment, research staff opened the envelope to learn of intervention assignment."</i></p>	Different interpretation	Consider differently incomplete or unclear description
		Review 10901: Researchers were blind until after the baseline assessment. Participants were not blinded.; Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report [^]	Category of disagreement	Reason of disagreement
Blinding of outcome assessment	Schoen 2007	Review 3603: Outcome assessor was not blinded.; High risk	<p><i>"In total 72 patients were screened by a maxillofacial surgeon (PJS) and prosthodontist (HR)."</i></p> <p><i>"All clinical assessments were performed by the investigator (PJS) who was not involved in treatment of the patients."</i></p>	Different interpretation	Consider differently incomplete or unclear description
		Review 5005: Outcome assessor may have been unaware of allocation: "All clinical assessments were performed by the investigator (PJS) who was not involved in treatment of the patients."; Low risk			
	Geroin 2011	Review 6185: Not done; High risk	<p><i>"All patients were evaluated by the same examiner (an experienced internal coworker) who was not aware of the treatment received by the patients"; Low Risk</i></p>	Missed information from the study report	
		Review 9645: Quote: "All patients were evaluated by the same examiner (an experienced internal coworker) who was not aware of the treatment received by the patients"; Low Risk			
	McCambridge 2004	Review 8969: As one interventionist was the study PI, a second independent interviewer who was blind to study condition was employed to conduct 3 month follow-ups, and an additional interviewer who was blind to initial group allocation was employed for 12 months follow-ups; Low Risk	<p><i>"further area of possible bias was that intervention recipients might report more favourable outcome data to the researcher who had delivered the intervention (J.M.). To study any such bias, a second independent interviewer who was blind to study condition, was employed to interview a sample of participants."</i></p>	Different interpretation	Consider differently incomplete or unclear description
		Review 7025: A second independent interviewer who was blind to study condition was employed to interview a sample of participants, though not all participants; Unclear Risk			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

[^] Information that were highlighted in the study report to support the analysis process.

Risk of bias item	Study Name	Support for judgement*	What is reported in the study report^	Category of disagreement	Reason of disagreement
Incomplete outcome data	Altmaier 1992	Review 1822: All subjects recorded follow-up data; Low Risk	[From table] “The n = 21 for control group and n = 24 for psychological group on all process measures.” [From table] The n = 21 for each group at each assessment.]	Different interpretation	Consider differently incomplete or unclear description
		Review 7407: Inadequately reported; High Risk			
	Killen 1984	Review 146: 11/75 recruited dropped out before full treatment, and are excluded from analyses.; Low Risk	“The first 75 were accepted into the study. Seven failed to attend (...) two dropped (...). The final sample (N = 64).”	Missed information from the study report	
		Review 3999: Losses to follow-up not reported, all participants included; Unclear Risk			
	Creager 2008	Review 986: There was a huge loss to follow up (only 50% completed the 6 month follow up) in this study and therefore there is a high risk of attrition bias; High Risk	“The remaining 525 patients met the inclusion criteria (...) The remaining 430 patients met their criteria for randomization (...).The ITT population consisted of 370 randomized patients (...). The per-protocol patient population consisted of 214 randomized patients”	Different interpretation	Consider differently intention-to-treat analysis
		Review 5262: Unclear why of patients stopped medication, unclear whether data presented represents intention-to-treat or per-protocol analysis			

*Supports for judgement and risk of bias assessments for the two reviews compared. The number of the review corresponds to the last 4 digits of the DOI.

^ Information that were highlighted in the study report to support the analysis process.

Supplementary Appendix 5| Reasons for disagreements in cases of a different interpretation of the same information; focus on “low” versus “high” disagreements.

Risk of bias item	Main reasons for disagreements	Examples
random sequence generation	Consider differently incomplete or unclear description	<p>“The names of communities within each group of three were written on individual cards, mixed and selected randomly: the first from each group was assigned to arm A (IEC alone), the second to arm B (IEC and STI management) and the third to arm C”; Low risk of bias</p> <p>“Names of communities within each triplet were written on separate cards and shuffled.”; High risk of bias</p>
allocation concealment	Consider differently envelopes description	<p>“Closed envelopes”; Low risk of bias</p> <p>“Closed envelopes, although not opaque.”; High risk of bias</p>
	Confusion in the definition of the item	<p>“pg. 2 - Methods - randomisation was done centrally to preserve allocation concealment”; Low risk of bias</p> <p>“904 patients were eligible for the study. 446 patients were randomised (49%). Due to the number of patients declining screening, there is an increased risk of inclusion bias.”; High risk of bias</p>
	Confusion with blinding	<p>“States used “preprogrammed laptop computer”. Remote site”; Low risk of bias</p> <p>“participants were told to which compound they had been allocated.”; High risk of bias</p>
blinding of participants and personnel	Assess risk differently if blinding was not feasible because of the type of intervention	<p>“Not possible to blind participants”; Low risk of bias</p> <p>“Participants were not blinded for provided treatment. This is inherent to study design”; High risk of bias</p>
	Outcome considered not influenced by blinding	<p>“Not possible to blind but most of the outcomes not likely to be influenced by lack of blinding.”; Low risk of bias</p> <p>“Not blinded due to nature of intervention.”; High risk of bias</p>
	Confusion with allocation concealment	<p>“participants were randomly allocated to either intervention or control group by an independent party”; Low risk of bias</p> <p>“Control group did not receive the comparable non-exercise related attention to the intervention group”; High risk of bias</p>

Risk of bias Item	Main reasons for disagreements	Examples
blinding of outcome assessment	outcome considered not influenced by blinding	“No information given about whether patients or assessors were blind to physician allocation but primary outcomes (treatment outcome and patient reported physician cultural competency) judged unlikely to be affected by lack of blinding”; Low risk of bias “Unblinded.”; High risk of bias
	Consider differently patient reported outcomes when patients are blinded or not to the intervention	“Insufficient information available to assess”; Low risk of bias “Comment: depression assessed by patient self-report”; High risk of bias
	Assess risk differently if blinding was not feasible because of the type of intervention	“Unclear blinding of outcome assessment”; Low risk of bias “blinding not possible due to intervention”; High risk of bias
incomplete outcome data	Use different cut-off for the rate of missing data	“11 withdrawals (10%).”; Low risk of bias “Comment: there were post-randomisation drop-outs”; High risk of bias
	Focus on number vs reasons/precise report of missing data	“Numbers and reasons for dropouts and withdrawals in all intervention groups were described.”; Low risk of bias “20 drop-outs (27.2%) with 4 deaths (3 males, 1 female) from cardiovascular events”; High risk of bias
	Consider differently incomplete or unclear description	“Women who were untraceable or unsuitable for follow-up were excluded, other losses included as smokers”; Low risk of bias “167/1287 (12.9%) (C = 83, I = 84) excluded from analysis due to moving away, being untraceable or deemed unsuitable for follow-up (e.g. miscarriage). 1120 in sample. 51/1287 non-responders were included as continuing smokers.” High risk of bias
	Use different cut-off for difference in the rate missing data between different arms/comparisons	“Similar rates of withdrawal between arms. Withdrawals: 36 BUD, 51 placebo”; Low risk of bias “Dropout higher in placebo group (35% vs 25% in budesonide group). ITT used.”; High risk of bias