Assessing the feasibility and impact of clinical trial trustworthiness checks via an application to Cochrane Reviews: Stage 2 of the INSPECT-SR project

Jack Wilkinson, Calvin Heal, Georgios A. Antoniou, Ella Flemyng, Love Ahnström, Alessandra Alteri, Alison Avenell, Timothy Hugh Barker, David N. Borg, Nicholas JL. Brown, Rob Buhmann, Jose A. Calvache, Rickard Carlsson, Lesley-Anne Carter, Aidan G. Cashin, Sarah Cotterill, Kenneth Färnqvist, Michael C. Ferraro, Steph Grohmann, Lyle C. Gurrin, Jill A. Hayden, Kylie E. Hunter, Natalie Hyltse, Lukas Jung, Ashma Krishan, Silvy Laporte, Toby J. Lasserson, David RT. Laursen, Sarah Lensen, Wentao Li, Tianjing Li, Jianping Liu, Clara Locher, Zewen Lu, Andreas Lundh, Antonia Marsden, Gideon Meyerowitz-Katz, Ben W. Mol, Zachary Munn, Florian Naudet, David Nunan, Neil E. O'Connell, Natasha Olsson, Lisa Parker, Eleftheria Patetsini, Barbara Redman, Sarah Rhodes, Rachel Richardson, Martin Ringsten, Ewelina Rogozińska, Anna Lene Seidler, Kyle Sheldrick, Katie Stocking, Emma Sydenham, Hugh Thomas, Sofia Tsokani, Constant Vinatier, Colby J. Vorland, Rui Wang, Bassel H. Al Wattar, Florencia Weber, Stephanie Weibel, Madelon van Wely, Chang Xu, Lisa Bero, Jamie J. Kirkham

PII: S0895-4356(25)00157-X

DOI: https://doi.org/10.1016/j.jclinepi.2025.111824

Reference: JCE 111824

To appear in: Journal of Clinical Epidemiology

Received Date: 5 December 2024

Revised Date: 24 April 2025

Accepted Date: 5 May 2025

Please cite this article as: Wilkinson J, Heal C, Antoniou GA, Flemyng E, Ahnström L, Alteri A, Avenell A, Barker TH, Borg DN, Brown NJ, Buhmann R, Calvache JA, Carlsson R, Carter L-A, Cashin AG, Cotterill S, Färnqvist K, Ferraro MC, Grohmann S, Gurrin LC, Hayden JA, Hunter KE, Hyltse N, Jung L, Krishan A, Laporte S, Lasserson TJ, Laursen DR, Lensen S, Li W, Li T, Liu J, Locher C, Lu Z, Lundh A, Marsden A, Meyerowitz-Katz G, Mol BW, Munn Z, Naudet F, Nunan D, O'Connell NE, Olsson N, Parker L, Patetsini E, Redman B, Rhodes S, Richardson R, Ringsten M, Rogozińska E, Seidler AL, Sheldrick K, Stocking K, Sydenham E, Thomas H, Tsokani S, Vinatier C, Vorland CJ, Wang R, Al Wattar BH, Weber F, Weibel S, van Wely M, Xu C, Bero L, Kirkham JJ, Assessing the feasibility and impact of clinical trial



trustworthiness checks via an application to Cochrane Reviews: Stage 2 of the INSPECT-SR project, *Journal of Clinical Epidemiology* (2025), doi: https://doi.org/10.1016/j.jclinepi.2025.111824.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2025 The Author(s). Published by Elsevier Inc.

Assessing the feasibility and impact of clinical trial trustworthiness checks via an application to Cochrane Reviews: Stage 2 of the INSPECT-SR project

4

5 Jack Wilkinson^{1*}, Calvin Heal^{1*}, Georgios A Antoniou^{2,3}, Ella Flemyng⁴, Love Ahnström⁵, Alessandra Alteri⁶, Alison Avenell⁷, Timothy Hugh Barker⁸, David N Borg^{9,10}, Nicholas JL Brown¹¹, Rob Buhmann¹², 6 Jose A Calvache^{13,14}, Rickard Carlsson¹¹, Lesley-Anne Carter¹, Aidan G Cashin^{15,16}, Sarah Cotterill¹, 7 Kenneth Färnqvist¹⁷, Michael C Ferraro^{15,16}, Steph Grohmann⁴, Lyle C Gurrin¹⁸, Jill A Hayden¹⁹, Kylie E 8 Hunter²⁰, Natalie Hyltse¹¹, Lukas Jung²¹, Ashma Krishan¹, Silvy Laporte²², Toby J Lasserson⁴, David RT 9 Laursen^{23,24}, Sarah Lensen²⁵, Wentao Li²⁶, Tianjing Li²⁷, Jianping Liu²⁸, Clara Locher²⁹, Zewen Lu¹, 10 11 Andreas Lundh^{23,24,30}, Antonia Marsden¹, Gideon Meyerowitz-Katz³¹, Ben W Mol²⁵, Zachary Munn²⁶, Florian Naudet²⁹, David Nunan³², Neil E O'Connell³³, Natasha Olsson⁵, Lisa Parker³⁴, Eleftheria 12 Patetsini³⁵, Barbara Redman³⁶, Sarah Rhodes¹, Rachel Richardson³⁷, Martin Ringsten³⁸, Ewelina 13 Rogozińska³⁹, Anna Lene Seidler²⁰, Kyle Sheldrick⁴⁰, Katie Stocking¹, Emma Sydenham⁴¹, Hugh 14 Thomas⁴², Sofia Tsokani^{37,43} Constant Vinatier²⁹, Colby J Vorland⁴⁴, Rui Wang²⁰, Bassel H Al Wattar⁴⁵, 15

- 16 Florencia Weber⁴⁶, Stephanie Weibel⁴⁶, Madelon van Wely⁴⁷, Chang Xu^{48, 49}, Lisa Bero^{50†}, Jamie J
- 17 Kirkham^{1†}
- 18
- 19 *Joint first authorship
- 20 +Joint senior authorship
- 21
- Centre for Biostatistics, The University of Manchester, Manchester Academic Health Science
 Centre, Manchester, UK
- 24 2. Manchester Vascular Centre, Manchester University NHS Foundation Trust, Manchester, UK
- 25 3. Division of Cardiovascular Sciences, School of Medical Sciences, Manchester Academic Health
- 26 Science Centre, The University of Manchester, Manchester, UK
- 27 4. Evidence Production and Methods Directorate, Cochrane Central Executive, London, UK
- 28 5. Department of Clinical Neuroscience, Karolinska Institutet, Stockholm, Sweden
- 29 6. Obstetrics and Gynaecology Unit, IRCCS San Raffaele Scientific Institute, Milan, Italy
- 30 7. Health Services Research Unit, University of Aberdeen, Aberdeen, UK
- 8. Health Evidence Synthesis, Recommendations and Impact, School of Public Health, The University
 of Adelaide, SA, Australia
- 9. Australian Sports Commission, Australian Institute of Sport, Bruce, Australian Capital Territory,Australia
- 35 10. Queensland University of Technology, School of Exercise Nutrition Sciences, Brisbane,
- 36 Queensland, Australia.

- 37 11. Department of Psychology, Linnaeus University, Växjö, Sweden
- 38 12. University of the Sunshine Coast, School of Health, Queensland, Australia
- 39 13. Department of Anesthesiology, Erasmus MC, The Netherlands
- 40 14. Department of Anesthesiology, Universidad del Cauca, Colombia
- 41 15. Centre for Pain IMPACT, Neuroscience Research Australia, Randwick, Australia
- 42 16. School of Health Sciences, Faculty of Medicine and Health, University of New South Wales43 Sydney, Australia
- 44 17. Department of Molecular medicine and surgery, Karolinska Institute, Stockholm, Sweden
- 45 18. School of Population and Global Health, The University of Melbourne
- 46 19. Department of Community Health & Epidemiology, Dalhousie University, Halifax, Canada
- 47 20. NHMRC Clinical Trials Centre, University of Sydney, Australia
- 48 21. Independent researcher, Heidelberg, Germany
- 49 22. Clinical Pharmacology Unit, INSERM U1059 Sainbiose, University Hospital of Saint-Etienne,
 50 France
- 51 23. Cochrane Denmark & Centre for Evidence-Based Medicine Odense (CEBMO), Department of
- 52 Clinical Research, University of Southern Denmark, Odense, Denmark
- 53 24. Open Patient data Explorative Network (OPEN), Odense University Hospital, Odense, Denmark
- 54 25. Department of Obstetrics, Gynaecology and Newborn Health, Royal Women's Hospital,
 55 University of Melbourne, Melbourne, VIC, Australia
- 26. National Perinatal Epidemiology and Statistics Unit, Centre for Big Data Research in Health and
 School of Women's and Children's Health, The University of New South Wales, Sydney, Australia
- 27. Department of Ophthalmology and Department of Epidemiology, University of ColoradoAnschutz Medical Campus, Colorado, USA
- 28. Centre for Evidence-Based Chinese Medicine, Beijing University of Chinese Medicine, Beijing,China.
- 62 29. Univ Rennes, CHU Rennes, Inserm, EHESP, Irset (Institut de recherche en santé, environnement
- 63 et travail) UMR_S 1085, Centre d'investigation clinique de Rennes (CIC1414), F-35000 Rennes,
- 64 France
- 30. Department of Respiratory Medicine and Infectious Diseases, Copenhagen University Hospital,Bispebjerg and Frederiksberg, Denmark
- 67 31. School of Nursing, University of Wollongong, Australia
- 68 32. Centre for Evidence-Based Medicine, Nuffield Department of Primary Care Health Sciences,
- 69 University of Oxford, Oxford, UK
- 33. Dept of Health Sciences, Centre for Health and Wellbeing across the Lifecourse, Brunel University
 London, LK
- 71 London, UK

- 72 34. Faculty of Medicine & Health, Charles Perkins Centre, University of Sydney, NSW, Australia
- 73 35. York Trials Unit, Department of Health Sciences, University of York, York, UK
- 74 36. New York University Grossman School of Medicine, New York, USA
- 75 37. Methods Support Unit, Cochrane, UK
- 76 38. Cochrane Sweden, Skåne University Hospital, Lund University, Lund, Sweden
- 77 39. Meta-Analysis Group, Institute of Clinical Trials and Methodology, MRC Clinical Trials Unit at UCL,
- 78 London, UK
- 79 40. Faculty of Medicine, University of New South Wales, Australia
- 80 41. Cochrane Central Editorial Service, London, UK
- 81 42. The Lancet Gastroenterology & Hepatology, London, UK
- 43. Laboratory of Hygiene, Social & Preventive Medicine and Medical Statistics, School of Medicine,
- 83 Aristotle University of Thessaloniki, Thessaloniki, Greece
- 84 44. Department of Epidemiology and Biostatistics, Indiana University School of Public Health-
- 85 Bloomington, USA
- 86 45. Clinical Trials Unit, Anglia Ruskin University, Chelmsford, UK
- 46. University Hospital Würzburg, Department of Anaesthesiology, Intensive Care, Emergency andPain Medicine, Würzburg, Germany
- 47. Centre for Reproductive Medicine, Department of Obstetrics and Gynaecology, AmsterdamUniversity Medical Center, Netherlands
- 91 48. The Third Department of Hepatic Surgery, Eastern Hepatobiliary Surgery Hospital, Third Affiliated
- 92 Hospital, Second Military Medical University, Naval Medical University, Shanghai, China
- 93 49. Proof of Concept Center, Eastern Hepatobiliary Surgery Hospital, Third Affiliated Hospital, Second
- 94 Military Medical University, Naval Medical University, Shanghai, China
- 95 50. University of Colorado Anschutz Medical Campus, Colorado, USA
- 96

97 What is new

- 98 An extensive list of potential checks for assessing study trustworthiness was assessed via an • 99 application to 95 randomised controlled trials (RCTs) in 50 Cochrane Reviews. 100 Following application of the checks, assessors had concerns about the authenticity of 32% of • 101 the RCTs. 102 If these RCTs were excluded, 22% of meta-analyses would have no remaining RCTs. 103 However, the study showed that some checks were frequently infeasible, and others could be • 104 easily misunderstood or misinterpreted.
- The study restricted assessment to meta-analyses including five or fewer RCTs, which might distort the impact of applying the checks.
- 107

108

109 Abstract

110 Background

111 The aim of the INSPECT-SR project is to develop a tool to identify problematic RCTs in 112 systematic reviews. In Stage 1 of the project, a list of potential trustworthiness checks was 113 created. The checks on this list must be evaluated to determine which should be included in 114 the INSPECT-SR tool.

115 Methods

116 We attempted to apply 72 trustworthiness checks to RCTs in 50 Cochrane Reviews. For each, 117 we recorded whether the check was passed, failed or possibly failed, or whether it was not

118 feasible to complete the check. Following application of the checks, we recorded whether we

119 had concerns about the authenticity of each RCT. We repeated each meta-analysis after

removing RCTs flagged by each check, and again after removing RCTs where we had concerns

121 about authenticity, to estimate the impact of trustworthiness assessment. Trustworthiness

122 assessments were compared to Risk of Bias and GRADE assessments in the reviews.

123

124 Results

95 RCTs were assessed. Following application of the checks, assessors had some or serious 125 concerns about the authenticity of 25% and 6% of the RCTs, respectively. Removing RCTs with 126 127 either some or serious concerns resulted in 22% of meta-analyses having no remaining RCTs. However, many checks proved difficult to understand or implement, which may have led to 128 129 unwarranted scepticism in some instances. Furthermore, we restricted assessment to metaanalyses with no more than 5 RCTs (54% contained only 1 RCT), which will distort the impact 130 131 on results. No relationship was identified between trustworthiness assessment and Risk of 132 Bias or GRADE.

133

134 Conclusions

This study supports the case for routine trustworthiness assessment in systematic reviews, as problematic studies do not appear to be flagged by Risk of Bias assessment. The study produced evidence on the feasibility and impact of trustworthiness checks. These results will be used, in conjunction with those from a subsequent Delphi process, to determine which checks should be included in the INSPECT-SR tool.

140

141 Plain language summary

142 Systematic reviews collate evidence from randomised controlled trials (RCTs) to find out 143 whether health interventions are safe and effective. However, it is now recognised that the

findings of some RCTs are not genuine, and some of these studies appear to have been 144 fabricated. Various checks for these "problematic" RCTs have been proposed, but it is 145 necessary to evaluate these checks to find out which are useful and which are feasible. We 146 applied a comprehensive list of "trustworthiness checks" to 95 RCTs in 50 systematic reviews 147 to learn more about them, and to see how often performing the checks would lead us to 148 classify RCTs as being potentially inauthentic. We found that applying the checks led to 149 concerns about the authenticity of around 1 in 3 RCTs. However, we found that many of the 150 checks were difficult to perform and could have been misinterpreted. This might have led us 151 to be overly sceptical in some cases. The findings from this study will be used, alongside other 152 evidence, to decide which of these checks should be performed routinely to try to identify 153 problematic RCTs, to stop them from being mistaken for genuine studies and potentially being 154 155 used to inform healthcare decisions.

156

157 MAIN TEXT

158 Background

Systematic reviews of randomised controlled trials (RCTs) aim to include all trials that address 159 the review question and meet the prespecified eligibility criteria. There is an understanding 160 that RCTs included in a systematic review should be scrutinised for their internal validity, for 161 162 example, using Risk of Bias tools (1, 2). These assessments require that the reviewer can trust what is written in a trial report to be an authentic account of what took place. However, this 163 no longer appears to be tenable as a default assumption, as recent large-scale assessments 164 165 have cast doubt on the veracity of many RCTs submitted to journals (3) or published in 166 systematic reviews (4). Recent examples, such as ivermectin for COVID-19, illustrate how the 167 failure to routinely interrogate the authenticity of eligible RCTs in systematic reviews allows 168 fake studies to influence patient care (5).

Cochrane defines 'problematic studies' as studies where there are 'serious questions about 169 170 the trustworthiness of the data or findings' (6). Problematic studies could represent instances 171 of academic misconduct such as research fraud, or could be the result of critical errors in trial 172 processes. Cochrane policy, introduced in 2021, states that potentially problematic RCTs 173 should not be included in a systematic review (6, 7). This prompts the question of what criteria could be used to identify problematic studies, which may appear to be high-quality on the 174 175 basis of traditional Risk of Bias assessment (8). Cochrane's implementation guidance recognises that a number of methods for identifying problematic studies have been proposed, 176 but does not recommend a method at this time. 177

The aim of the INSPECT-SR (INveStigating ProblEmatic Clinical Trials in Systematic Reviews) project is to develop a tool that can be used by systematic reviewers to assess the trustworthiness of RCTs (9). Several tools have recently been proposed for this purpose (10-14). However, none of these have involved a comprehensive evaluation and subsequent selection of potential trustworthiness checks. In Stage 1 of the development process, we

identified an extensive list of potential trustworthiness checks (15). A tool including all of 183 these checks would not be practicable, and we anticipate that many of the checks will turn 184 185 out to be infeasible or otherwise not useful. In Stages 2 (application to Cochrane Reviews) 186 and 3 (Delphi survey), the checks on this list will be evaluated to determine which should be included in the final tool. These results will then feed into a series of consensus meetings 187 (Stage 4) which will be used to develop a draft version of the INSPECT-SR tool. The draft tool 188 will then be tested in the assessment of RCTs (Stage 5). Feedback from Stage 5 will be used to 189 finalise the tool. The current study describes Stage 2 of the project, in which the identified 190 191 checks were applied to RCTs included in a sample of Cochrane Reviews, in order to evaluate their feasibility and impact on review results, and to evaluate how often assessors had 192 193 concerns about the authenticity of RCTs after applying the checks.

194

195 Methods

A protocol describing the INSPECT-SR project methods has previously been published (9). We undertook a large, collaborative project in which assessors applied a series of 72 trustworthiness checks to RCTs included in 50 Cochrane Reviews. The University of Manchester Ethics tool was used to determine that ethical approval was not required for this study (30th Sept 2022).

201

202 Description of trustworthiness checks

Prior to this exercise, a list of trustworthiness checks was assembled using a scoping review 203 204 (16), qualitative study (17) and survey of experts (15). This list contained 116 checks arranged into five domains: Inspecting results in the paper, Inspecting the research team and their work, 205 206 Inspecting conduct, governance and transparency, Inspecting text and publication details, and 207 Inspecting individual participant data. In the current study we only considered the first four 208 domains, as individual participant data are not generally available during systematic reviews 209 and meta-analyses based on aggregate data; nor were they available to us. An extension to the INSPECT-SR tool based on the checks in the fifth domain, which can be applied when 210 individual participant data are available, 'INSPECT-IPD', has been funded for development 211 (Reference: NIHR30355). The first four domains included 76 checks (Tables 1 and 2). We made 212 some modifications to the list in preparation for the current study, in consultation with the 213 project expert advisory panel. This included refining the language of some items to improve 214 215 clarity. To assist assessors in applying the checks, we drafted brief explanations for each check 216 (S Tables 1 to 4). Four checks (checks 45, 66, 67, 72, Tables 1 and 2) were not assessed as they 217 were not considered practicable in the context of the present study. Consequently, 72 checks 218 were assessed.

219

220 Description of assessors

The INSPECT-SR working group includes a core management group and an expert advisory 221 panel. Members of both were invited to act as assessors for the current study. We also invited 222 223 additional collaborators who had expressed an interest in contributing to the development 224 process. Collaborators were identified from a variety of sources. We invited attendees at presentations relating to the project to contact us to express an interest, and also invited 225 individuals who had expressed an interest in the topic to JW using personalised emails and 226 via social media. All assessors were considered to have sufficient expertise in research 227 methods (specifically, to evaluate RCTs) to enable them to undertake the assessment. We did 228 229 not require assessors to hold any particular qualification however. We did not require assessors to have specialist expertise relating to research integrity (for example, use of 230 231 forensic statistical methods or investigation of misconduct cases), as a key objective was to 232 learn about the usefulness and feasibility of the checks when applied by potential users of the 233 INSPECT-SR tool (i.e. systematic reviewers, researchers, peer reviewers) who would not be 234 expected to possess this specialist knowledge. Assessors who were considered to have made 235 a substantial contribution to data acquisition and critical review of manuscript drafts, were 236 given the option to co-author the manuscript.

237

238 Selection of Cochrane Reviews and RCTs

The sample size of 50 Cochrane reviews represented a number that was considered feasible 239 240 to complete, while facilitating the evaluation of feasibility and impact of applying the checks across different topic areas. A preliminary pilot was conducted on a small number of RCTs to 241 confirm this. The 50 reviews were purposefully selected from the Cochrane Library. To be 242 243 eligible, a review could not be authored or co-authored by the assessor, and could not contain 244 RCTs authored or co-authored by the assessor, to prevent any conflict of interest in 245 conducting the assessment. As a feasibility requirement, we also required that the review 246 contained at least one (meta-) analysis containing one to five RCTs. For brevity, we use the 247 term 'meta-analysis' in this article to describe an analysis which produces a pooled average 248 estimate and confidence interval for a treatment effect on an outcome based on the included 249 studies, recognising that, when there is only one study, this involves reporting the estimate and confidence interval from that study. The RCTs in the first eligible meta-analysis in the 250 251 review were subjected to the trustworthiness assessment, as a feasibility constraint. We also 252 required that the review had not already undergone a trustworthiness assessment as part of 253 the review process, since this could have resulted in the prior removal of problematic studies, 254 distorting our assessment. Assessors were asked to suggest a topic with which they were 255 broadly familiar. We attempted to match assessors to review topics, to replicate a typical scenario in which INSPECT-SR would be used (a systematic review would often be undertaken 256 257 by someone with some relevant subject-matter knowledge). We then selected the most 258 recent Cochrane Review relating to the topic suggested by the assessor that met the eligibility 259 criteria. Assessors did not always have subject-matter knowledge relating to the review(s) 260 they assessed, however. For example, some assessors were primarily methodologists, with limited clinical knowledge of the subject matter. For these people, we attempted to select 261

review topics to cover a broad range of health areas. We asked each assessor to record their familiarity with the review topic during data extraction (little or no familiarity, some familiarity, or high familiarity).

265

266 Data extraction and trustworthiness assessment

A bespoke data extraction form was produced, and was revised following piloting on a small 267 number of RCTs, and can be accessed at <u>https://osf.io/9pyw2/</u>. Assessors were informed of 268 software that could be used to implement some of the statistical checks. Examples include 269 270 the scrutiny package in R (18), online applications created to implement some checks e.g. 271 applications for performing GRIMMER (Granularity-Related Inconsistency of Means Mapped to Error Repeats) and SPRITE (Sample Parameter Reconstruction via Iterative Techniques)(19, 272 273 20), or Microsoft Excel (21) for basic statistical checks, but it was not a requirement to use any particular software to undertake the assessment. For each Cochrane Review, the assessor 274 275 extracted data and applied the list of checks to each RCT in the meta-analysis. An exception 276 was check 26 - Is there heterogeneity across studies in degree of imbalance in baseline 277 characteristics (in meta-analysis), which was assessed only once per review. The assessor 278 extracted the year of publication for each RCT, the summary data entered in the meta-279 analysis, and Risk of Bias and GRADE (Grading of Recommendations Assessment, Development and Evaluation) (22) assessments as presented in the Cochrane Review. 280

A second assessor performed a quality check of accuracy and completeness of this 281 information following extraction. Any disagreements were resolved by discussion between 282 283 assessors and a third team member (JW). The assessor attempted to apply each of the 72 checks to the trial, selecting one of four response options: not feasible; passed the check; 284 possible fail; fail. For each check, assessors were asked to supply free text to explain their 285 assessment. The country or countries in which the RCT was conducted was also recorded. 286 287 After applying the checks, assessors recorded their answer to the question "Do you have concerns about the authenticity of this study?" using one of four response options: no; some 288 289 concerns; serious concerns; don't know. Assessors were asked if they had performed any 290 additional checks (not included on the list) and if so, to describe both the checks and the results of applying them. There was space for the assessor to add any additional information, 291 292 and to provide an estimate of how many hours it took them to assess the RCT. The intention had been for one assessor to assess all of the RCTs in the review, before checking by a second 293 assessor. However, some assessors failed to complete the assessment of all RCTs in their 294 295 allocated review, and so for several reviews the RCTs were split between two assessors, before being checked by a third assessor. 296

297

298 Statistical analysis

We summarised trial and Cochrane Review characteristics, and the responses for each check. We calculated how often assessors had concerns about study authenticity. We evaluated the impact of applying each check by comparing the analysis/ meta-analysis including all trials as per the review to a version in which any RCTs flagged by the check were removed, in terms of the numbers of trials, sample size, change in effect estimate, 95% confidence interval width, heterogeneity, and change in inference.

The first two of these metrics were assessed over all reviews, while the remainder were assessed separately for binary and continuous outcomes. We used the metafor package (23) in R to perform all meta-analyses, using odds ratios to summarise treatment effects with binary outcomes, and standardised mean differences to summarise treatment effects with continuous outcomes. Random effects meta-analyses using the DerSimonian and Laird (24) method were performed, as the most typical method employed in systematic reviews (25, 26).

We assessed potential redundancies among the checks by plotting the responses for each 312 check for each RCT in an array. We made the post-hoc decision to undertake a hierarchical 313 314 cluster analysis, using complete agglomeration based on Gower dissimilarity, as implemented in the cluster package in R (27). The purpose of this analysis was to identify possible clusters 315 of checks that could potentially be combined. We used multinomial regression to assess the 316 relationship between trustworthiness assessment and each Risk of Bias domain, and ordinal 317 318 regression (proportional odds logistic regression) to consider the relationship between the 319 GRADE assessment and the number of trials flagged for concerns. We used likelihood ratio 320 tests for inference following regression model fits. We conducted an additional analysis which 321 had not been specified in the protocol, where we evaluated the relationship between the assessment for each check and the overall assessment of the trial using the N-1 chi-squared 322 323 test (28), to determine which checks were influential in reaching an overall assessment. The N-1 chi-squared test was used in anticipation of small expected counts (29). This analysis was 324 performed in trials where the check was considered to be feasible, and the assessments were 325 analysed as 'passed' vs 'fail or possible fail'. We used a post-hoc significance threshold of 1% 326 to highlight checks associated with the overall assessment, creating contingency tables 327 328 (outcome of check vs overall assessment) for these checks to determine whether failing the 329 check was associated with an assessor having overall concerns. We categorised the free-text responses to the question asking how long it took to complete the assessment in a post-hoc 330 331 fashion (less than 90 minutes, 90 minutes to 3 hours, more than 3 hours). The dataset and analysis code for this study are available at https://osf.io/9pyw2/ . 332

333

334 Results

We included a total of 95 RCTs from 50 Cochrane Reviews. The reviews were from 24 different Cochrane Groups (S Table 5). Assessors considered themselves to have high familiarity with the review topic for 7/50 (14%) reviews, some familiarity for 20/50 (40%) of reviews, and little or no familiarity for 23/50 (46%). The characteristics of included Cochrane Reviews are shown

- in Table 3. The median (IQR) number of RCTs per review was 1 (1 to 1.9). 27 (54%) contained 339 only 1 RCT. The median (IQR) number of participants in the assessed RCTs was 71 (40 to 174). 340 Fifteen of 95 (16%) were conducted in multiple countries, with the remaining 80 taking place 341 342 in one of 21 different countries (S Table 6). Twenty-four (26%) RCTs took less than 90 minutes 343 to assess, 29 (31%) took between 90 minutes and 3 hours, and 40 (42%) took more than 3 hours.
- 344
- 345

Responses to individual trustworthiness checks 346

Figure 1 and S Table 7 summarise the responses for each check, and S Figure 1 shows the 347 348 study-level responses for each check. S Figure 2 shows how the checks are clustered in the dataset. Missing data for trustworthiness checks were infrequent, with only one check having 349 350 missing data for as many as five RCTs (check 42). Check 26 is 'missing' for 10 RCTs, as it was only assessed once per review. A number of checks were considered to have 'failed' or 351 'possibly failed' often. The five checks most often receiving an assessment of 'failed' or 352 353 'possibly failed' were check 61 - Are the data publically available? (81%), check 30 - Are contributorship statements present? (69%), check 31 - Are contributorship statements 354 complete? (57%), check 64 - Has the study been prospectively registered? (56%), check 49 – Is 355 356 a funding source reported? (40%). Some statistical checks frequently resulted in responses of 'failed' or 'possibly failed'. Examples include check 12 - Are differences in variances in baseline 357 358 variables between randomised groups plausible? (28%), check 11 – Are statistical test [results] 359 of outcomes correct? (21%).

A number of checks were considered to be infeasible in most cases. The checks most 360 frequently considered infeasible were check 40 - Is the standard deviation of summary 361 statistics in multiple studies by same authors plausible (when compared to simulated or 362 bootstrapped data?) (99%), check 62 - Are additional patient data recorded in patient case 363 records beyond what is reported in the paper? (98%), check 38 - Is the distribution of non-first 364 365 digits in manuscripts from one author compatible with a genuine measurement process? 366 (90%), check 35 - Is any duplicate reporting acknowledged or explained? (89%), and check 29 - Are withdrawal and loss to follow-up in multiple trials by the same author consistent with 367 368 the expected (random) binomial distribution? (84%).

369

Overall assessment and relationship to individual checks 370

Overall, responses to the question "Do you have concerns about the authenticity of this 371 372 study?" were: no (60/95, 64%); some concerns (24/95, 25%); serious concerns (6/95, 6%). Pvalues from chi-squared tests looking at the outcome of each check against the overall 373 trustworthiness assessment of the study are shown in S Table 8. Noting that these analyses 374 were post-hoc and exploratory, 19 checks were associated with overall assessment using a 375 376 1% significance level. Contingency tables were inspected to examine the nature of these

associations (specifically to confirm that failing or possibly failing, rather than passing, a checkwas associated with the presence of concerns).

379 Of these 19, there were 11 checks for which failing (as opposed to passing) the check appeared to correlate with an assessment of overall concern: 1. Are any baseline data 380 implausible with respect to magnitude, frequency, or variance? (p=0.00001); 2. Is the number 381 of participant withdrawals compatible with the disease, age and timeline? (p=0.005); 8. Are 382 there any discrepancies between data reported in figures, tables and text? (p = 0.00006), 9. 383 Are any outcome data, including estimated treatment effects, implausible? (p = 0.000002), 19. 384 385 Are results internally consistent? (p=0.00008), 37. Does the statistics methods section use 386 generic language, suggesting lack of expert statistical input? (p=0.003), 51. Is the reported staffing adequate for the study conduct as reported? (p=0.009), 52. Is the recruitment of 387 participants plausible within the stated time frame for the research? (p=0.0005), 53. Is the 388 389 recruitment of participants plausible considering the epidemiology of the disease in the area of the study location? (p=0.0004), 56. Are there any concerns about unethical practice? 390 (p=0.001), 64. Has the study been prospectively registered? (p=0.004). 391

392

393 Impact of applying the trustworthiness assessments on systematic review results

S Table 9 and S Table 10 show the impact of removing RCTs flagged by each check (considered 394 395 individually) from meta-analysis, for binary and continuous outcomes respectively. In 396 continuous outcome meta-analyses, removal of RCTs flagged by a check resulted in a median 397 of 4% (IQR 0% to 12.5%, range 0% to 67%) of meta-analyses having no remaining trials. In binary outcome meta-analysis, the corresponding values were 4% (IQR 0% to 8%, range 0% 398 399 to 73%). The sample size of reviews would be reduced to a median (of means) of 93% (IQR 400 87% to 97%, range 27% to 100%) of the original size. The median (of means) number of trials 401 that would be removed from meta-analysis was 0.14 (IQR 0.06 to 0.24, range 0 to 1.52).

When RCTs were removed on the basis of the overall assessment (some or serious concerns), 402 33% of continuous outcome meta-analyses and 12% of binary outcome meta-analyses had no 403 404 remaining trials. Amongst meta-analyses with at least one RCT remaining, for binary outcome 405 meta-analyses, the mean ROR was 0.98; SE increased by 19% on average; none changed in 406 terms of statistical significance (using a 5% significance threshold); and the mean ratio of 407 confidence interval widths (width expressed as the ratio of upper to lower limit on OR scale) 408 was 4.52. For continuous outcome meta-analyses with at least one RCT remaining, the 409 average change in estimate was -0.02 SDs; SE (and, equivalently, CI width) increased by a mean of 5%; and none of the meta-analyses changed in terms of statistical significance. 410

411

412 Relationship between trustworthiness assessments, Risk of Bias and GRADE

We only investigated the relationship between overall trustworthiness assessment and risk of bias for reviews using the first version of the Cochrane RoB tool since there were only 10 reviews applying RoB 2. Multinomial regression did not indicate associations between any risk of bias domain and overall concern, with the exception of allocation concealment. However, 417 this was not in the expected direction, with concerns expressed more often for studies with

418 unclear or low bias assessment compared to high bias assessment (p=0.01). The estimated 419 relationship between number of trials flagged for concerns and GRADE assessment was 420 imprecise (OR = 0.68, 95% CI = (0.39 to 1.17)).

421

422 New checks used by assessors

Assessors described eight checks which they used and which they felt were additional to the 423 list of checks assessed in the study. Two of these – checking for trial registration, and checking 424 425 the author list – were already covered by the primary list. Three others were variations of existing checks - checking the certification status of the ethical committee or institutional 426 427 review board, looking at a related publication of a subgroup, and checking for consistency 428 with the main article. Three were new: looking to see whether the authors exclusively worked 429 together, checking whether the first author's department had participated in other RCTs, and 430 looking into the reported funder.

431

432 Discussion

An extensive list of trustworthiness checks was assessed for their feasibility and impact by 433 application to 95 RCTs in a sample of 50 Cochrane Reviews. The study allowed us to estimate 434 435 how often each of the checks would be considered infeasible for routine use in systematic reviews, how often each would fail, and what the impact of applying the check would be on 436 the estimates from meta-analysis. We found that, in the context of conducting a systematic 437 review, the checks can be applied to identify problematic studies. Furthermore, the findings 438 suggest that a substantial portion of meta-analyses would be left with no remaining RCTs if 439 failed checks were used to identify and exclude problematic studies. Amongst those with 440 remaining RCTs, there was a larger impact on precision than on the magnitude of effect 441 442 estimates. The study also found that, following application of the checks, assessors frequently 443 had concerns about the RCTs included in Cochrane Reviews, with "some concerns" being reported for 25% of studies, and "serious concerns" for a further 6%. 444

445 Feasibility of the checks

A number of checks were deemed generally infeasible. For example, assessments which 446 447 involved taking an author-wide view have been successfully implemented in particular cases 448 (e.g. (8, 30, 31)), but were not considered feasible by assessors in the context of the current 449 study. One possible reason is that these checks require additional data collection to find out more about the authors of a study, their research team, or their other publications. In a 450 451 previous survey of experts, the need for a trustworthiness tool to be practical and not too 452 burdensome was emphasised (15), and therefore checks which require the identification and comparison of additional studies are unlikely to be palatable. Other checks that were deemed 453 infeasible include checking for evidence of copied work, including copied sample 454 455 characteristics and results tables. Unless the copying is identified between RCTs that both 456 happen to be included in the review, it is difficult to see how this sort of check would be

457 practicable in the absence of automated solutions. Assessing the plausibility of various RCT

458 features is likely to be difficult without recourse to domain expertise. Clearly, the results459 indicate that it would not be feasible to apply such a long list of checks routinely, as this took

460 more than 3 hours for 42% of the trials.

461

462 Identification of problematic studies

Failed checks are potential indicators of a problematic trial. Some checks failed for most RCTs. 463 For example, the study agreed with previous work suggesting that many (in the present study, 464 465 most) RCTs are not prospectively registered (32-34), and few make the underlying data available (35). In relation to registration, assessors were much more likely to have concerns 466 467 about authenticity for studies that were not prospectively registered compared to those that 468 were. This could indicate that lack of registration was influential in reaching an overall judgement, or rather that studies with other problematic features were less likely to be 469 prospectively registered. Prospective registration is routinely considered in relation to 470 471 reporting bias, and an important question to be resolved in the INSPECT-SR development 472 process is whether there is additional value in considering prospective registration in the 473 assessment of trustworthiness.

Our findings also indicate that some checks may be prone to misinterpretation or 474 misapplication, which was suggested by high failure rates. In particular, several statistical 475 476 checks proved challenging. For example, 20% of RCTs were considered to have 'failed' or 477 'possibly failed' a check looking to see whether results of statistical tests of outcomes were 478 correct. Some of these failures might be attributable to the rounding of continuous variables 479 in published articles; p-values obtained from rounded summary statistics can differ from 480 those obtained from analysis of the underlying data, meaning the question of assessing 481 consistency cannot just be assessed reproducing the test and looking for an exact match (36). 482 Another example was checking differences in baseline variance between groups, which "failed" or "possibly failed" for 28% of RCTs. Assessors were directed to use an F test here. 483 However, this test has an inflated type 1 error rate for skewed variables (37), such that 484 485 rejection of the test assumptions may have been frequently mistaken for rejection of the hypothesis of equality of variances. Instances such as these may have led to unwarranted 486 skepticism about a study's authenticity in some instances. Although we did not detect 487 associations between failure of these checks and concerns in post-hoc analyses, it remains 488 489 possible that errors of this nature did sometimes occur, and may have influenced the overall 490 assessment of a trial's authenticity.

491

492 Impact of identifying problematic studies on systematic reviews

The impact of removing RCTs flagged by these checks from meta-analyses may appear alarming; for example, removing RCTs for which assessors expressed concerns would result in 11 of 50 meta-analyses with no RCTs remaining. However, caution is needed here due to limitations introduced by our study design. We only assessed meta-analyses containing five

497 or fewer trials in this study. Consequently, more than half contained only a single RCT, 498 although small numbers appear to be typical (38). This will exaggerate the number of reviews 499 with no remaining studies following trustworthiness assessment. Moreover, several other 500 metrics, such as the change in point estimate and associated uncertainty, could only be 501 evaluated in meta-analyses with at least one trial remaining following assessment. Due to the fact that many meta-analyses only included one trial initially, this subset will omit many of 502 the meta-analyses with any trials flagged at all, causing the impact of the checks on these 503 metrics to be understated. 504

505 In line with our expectations, there did not appear to be a clear association between Risk of Bias domains and overall trustworthiness assessment, reinforcing the premise that these 506 frameworks are evaluating different aspects of trials. Many problematic studies appear to 507 frequently describe perfectly sound methods (8). We were unable to ascertain whether there 508 509 is any link between GRADE and trustworthiness assessment, as our estimate of the relationship was too imprecise. We suggest that trustworthiness assessment should be 510 performed prior to Risk of Bias and GRADE assessments, because the value of assessing the 511 internal or external validity of a problematic study is doubtful. 512

- 513
- 514

515 Implications for development of INSPECT-SR and future directions

516

These observations have informed the development of the INSPECT-SR tool and 517 accompanying guidance. The findings highlight the need for careful curation of the checks 518 included in INSPECT-SR, and suggest that any statistical checks included in the tool would have 519 to be accompanied by detailed guidance to enable their application, as well as to prevent 520 misuse and misinterpretation. As technological solutions become available to facilitate some 521 522 useful but difficult checks, they can become part of the tool implementation. As the role of 523 automation, including artificial intelligence, is likely to expand in evidence synthesis, it will be 524 important to examine how it might enable or hinder detection of problematic RCTs (39). For 525 example, some checks, such as statistical checks, may be more amenable to automation than 526 checks that require more content knowledge, such as the plausibility of participant recruitment or effect sizes. 527

Additional future directions informed by this study will be development of training for INSPECT-SR and tools that can be applied to individual patient data or observational study designs. Creating a searchable, open archive of trials that have been evaluated with INSPECT-SR will aid all systematic reviewers and users of trials. Lastly, although INSPECT-SR is being developed for use by systematic reviewers, adaptations of the tool could also be useful to journal editors or publishers who screen trials for research integrity problems.

534

535 Conclusion

The study appears to reinforce the need for routine trustworthiness assessment in RCTs, 536 537 suggesting that problematic studies in systematic reviews may not be infrequent, and are not 538 detected by Risk of Bias assessment. Only two of the studies judged to be concerning had 539 associated retraction or expression of concern notices at the point of assessment, highlighting the need to evaluate other features in order to identify these untrustworthy trials. The time 540 taken to complete the full barrage of checks for each RCT was long, and would likely not be 541 practicable in the context of a typical systematic review. The goal of subsequent stages of the 542 INSPECT-SR project will be to identify a subset of these checks that are both feasible and 543 544 useful, and to implement these in the form of a tool that can be implemented by systematic reviewers. The results from this study will be used to select checks for this purpose, alongside 545 a Delphi study of experts and potential users of the tool. Both sets of results will be presented 546 547 to experts at a series of consensus meetings, which will be used to determine the content of 548 a draft version of INSPECT-SR. The draft version of the tool will then be tested in the assessment of RCTs, and feedback will be used to finalise the tool in early 2025. 549

- 550
- 551

552 Declarations

JW, CH, GAA, LB, JJK declare funding from NIHR (NIHR203568) in relation to the current 553 project. JW additionally declares Stats or Methodological Editor roles for BJOG, Fertility and 554 555 Sterility, Reproduction and Fertility, Journal of Hypertension, and for Cochrane Gynaecology and Fertility. CH declares a Statistical Editor role for Cochrane Colorectal. GAA additionally 556 declares a Statistical Reviewer role for the European Journal of Vascular and Endovascular 557 Surgery. LB additionally declares a role as Academic Meta-Research Editor for PLoS Biology, 558 and that The University of Colorado receives remuneration for service as Senior Research 559 560 Integrity Editor, Cochrane. JJK additionally declares a Statistical Editor role for The BMJ. EF is 561 employed by the Cochrane Collaboration and on the Editorial Board of Cochrane Evidence 562 Synthesis and Methods. ES is a Senior (Sign-off) Editor for Cochrane and contributed to Cochrane's Policy for managing potentially problematic studies. SL is an editor for Cochrane 563 564 Gynaecology and Fertility, Human Reproduction, and Fertility and Sterility. TJL is the Deputy Editor in Chief of The Cochrane Library and is an employee of The Cochrane Collaboration. 565 DNB is an associate editor for Research Quarterly for Exercise and Sport and a section editor 566 for Communications in Kinesiology. NEO is a member of the Cochrane Editorial Board and 567 568 holds an ERA-NET Neuron Co-Fund grant for a separate project. RR declares acting as an author and editor on Cochrane reviews. KS is an editor for Cochrane Gynaecology and 569 570 Fertility, and Fertility and Sterility. MvW declares to be co-ordinating editor for Cochrane 571 Gynaecology and Fertility and Cochrane Sexually Transmitted Infections, methodological editor for Human Reproduction Update and Editorial Editor for Fertility & Sterility. HT is 572 573 Deputy Editor of The Lancet Gastroenterology & Hepatology and is an employee of Elsevier. 574 SL received funding from the French National Research Agency (ANR-23-CE36-0006-01). AK is an editorial board member for BJGP Open. TLi serves as the Principal Investigator on a grant 575 from the National Eye Institute, National Institutes of Health that funds the work of Cochrane 576 577 Eyes and Vision US Project. She also acts as a sign-off editor for The Cochrane Library. ZM is

supported by an NHMRC Investigator Grant 1195676. ZM is an associate Editor for BMC 578 579 Medical Research Methodology and is on the Editorial Board for Clinical and Public Health Guidelines. RC is Editor-in-Chief at Meta-Psychology. CL is a work-package leader for the 580 581 doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101120360), funded 582 by the EU. CV received funding as part of the OSIRIS project (Open Science to Increase Reproducibility in Science); the OSIRIS (Open Science to Increase Reproducibility in Science) 583 project has received funding from the EU (grant agreement No. 101094725). FN received 584 funding from the French National Research Agency (ANR-23-CE36-0006-01), the French 585 586 ministry of health and the French ministry of research. He is a work-package leader in the OSIRIS project (Open Science to Increase Reproducibility in Science). The OSIRIS project has 587 588 received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101094725. He is a work-package leader for the 589 590 doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101120360), funded by the EU. DN declares having led/co-authored/co-authoring Cochrane Reviews. He also 591 592 declares having been part of the Cochrane Convenes initiative organised by Cochrane to consider the issue of misinformation, its impact on the health evidence ecosystem and 593 solutions to address it. LJ is the creator of the scrutiny package in R. WL is supported by an 594 595 NHMRC Investigator grant (GNT2016729). RW is supported by an NHMRC Investigator Grant (2009767) and acts as a Deputy Editor for Human Reproduction, and an editorial board 596 member for BJOG and Cochrane Gynaecology and Fertility. EF, SGTLa and RR declare 597 employment by Cochrane. TLa additionally declares authorship of a chapter in the Cochrane 598 Handbook for Systematic Reviews of Interventions and that he is a developer of standards for 599 600 Cochrane intervention reviews (MECIR). AL is on the editorial board of BMC Medical Ethics.

601

602 Ethical approval

603 The University of Manchester ethics decision tool was used on 30/09/22. Ethical approval

604 was not required for this study, since it involved appraisal of published research.

605

606 Funding

This study/project is funded by the NIHR Research for Patient Benefit programme (NIHR203568). The views expressed are those of the author(s) and not necessarily those of the NIHR or the Department of Health and Social Care.

610

611 References

Higgins JP, Altman DG, Gotzsche PC, Juni P, Moher D, Oxman AD, et al. The Cochrane
 Collaboration's tool for assessing risk of bias in randomised trials. BMJ. 2011;343:d5928.

Sterne JAC, Savovic J, Page MJ, Elbers RG, Blencowe NS, Boutron I, et al. RoB 2: a revised tool
 for assessing risk of bias in randomised trials. BMJ. 2019;366:14898.

616 3. Carlisle JB. False individual patient data and zombie randomised controlled trials submitted to 617 Anaesthesia. Anaesthesia. 2021;76(4):472-9.

4. Weeks J, Cuthbert A, Alfirevic Z. Trustworthiness assessment as an inclusion criterion for
systematic reviews—What is the impact on results? Cochrane Evidence Synthesis and Methods.
2023;1(10):e12037.

5. Hill A, Mirchandani M, Pilkington V. Ivermectin for COVID-19: Addressing Potential Bias and
Medical Fraud. Open Forum Infect Dis. 2022;9(2):ofab645.

6. Cochrane. Cochrane Policy for managing potentially problematic studies. Cochrane Database
 624 of Systematic Reviews: editorial policies Cochrane Library [Available from:
 625 <u>https://www.cochranelibrary.com/cdsr/editorial-policies</u>.

626 7. Boughton SL, Wilkinson J, Bero L. When beauty is but skin deep: dealing with problematic
627 studies in systematic reviews. Cochrane Database Syst Rev. 2021;6(6):ED000152.

628 8. O'Connell NE, Moore RA, Stewart G, Fisher E, Hearn L, Eccleston C, et al. Investigating the 629 veracity of a sample of divergent published trial data in spinal pain. Pain. 2023;164(1):72-83.

9. Wilkinson J, Heal C, Antoniou GA, Flemyng E, Alfirevic Z, Avenell A, et al. Protocol for the
development of a tool (INSPECT-SR) to identify problematic randomised controlled trials in systematic
reviews of health interventions. BMJ Open. 2024;14(3):e084164.

Mol BW, Lai S, Rahim A, Bordewijk EM, Wang R, van Eekelen R, et al. Checklist to assess
Trustworthiness in RAndomised Controlled Trials (TRACT checklist): concept proposal and pilot. Res
Integr Peer Rev. 2023;8(1):6.

Weibel S, Popp M, Reis S, Skoetz N, Garner P, Sydenham E. Identifying and managing
problematic trials: A research integrity assessment tool for randomized controlled trials in evidence
synthesis. Res Synth Methods. 2023;14(3):357-69.

Hunter KE, Aberoumand M, Libesman S, Sotiropoulos JX, Williams JG, Aagerup J, et al. The
Individual Participant Data Integrity Tool for assessing the integrity of randomised trials. Res Synth
Methods. 2024.

642 13. Abbott J, Acharya G, Aviram A, Barnhart K, Berghella V, Bradley CS, et al. Trustworthiness
643 criteria for meta-analyses of randomized controlled studies: OBGYN journal guidelines. Elsevier; 2024.
644 p. 101481.

64514.Grey A, Bolland MJ, Avenell A, Klein AA, Gunsalus CK. Check for publication integrity before646misconduct. Nature. 2020;577(7789):167-9.

Wilkinson J, Heal C, Antoniou GA, Flemyng E, Avenell A, Barbour V, et al. A survey of experts
to identify methods to detect problematic studies: stage 1 of the INveStigating ProblEmatic Clinical
Trials in Systematic Reviews project. J Clin Epidemiol. 2024;175:111512.

Bordewijk EM, Li W, van Eekelen R, Wang R, Showell M, Mol BW, et al. Methods to assess
research misconduct in health-related research: A scoping review. J Clin Epidemiol. 2021;136:189-202.

65217.Parker L, Boughton S, Lawrence R, Bero L. Experts identified warning signs of fraudulent653research: a qualitative study to inform a screening tool. J Clin Epidemiol. 2022;151:1-17.

In Jung L. scrutiny: Error Detection in Science. R package version 0.3.0. 2023 [Available from:
 <u>https://CRAN.R-project.org/package=scrutiny.</u>

Anaya J. The GRIMMER test: A method for testing the validity of reported measures ofvariability. PeerJ Preprints. 2016;4:e2400v1.

Heathers JA, Anaya J, van der Zee T, Brown NJ. Recovering data from summary statistics:
Sample parameter reconstruction via iterative techniques (SPRITE). PeerJ Preprints; 2018. Report No.:

660 2167-9843.

661 21. Microsoft Corporation. Microsoft Excel. 2018.

662 22. Guyatt GH, Oxman AD, Vist GE, Kunz R, Falck-Ytter Y, Alonso-Coello P, et al. GRADE: an 663 emerging consensus on rating quality of evidence and strength of recommendations. BMJ. 664 2008;336(7650):924-6.

665 23. Viechtbauer W. Conducting meta-analyses in R with the metafor package. 2010.

666 24. DerSimonian R, Laird N. Meta-analysis in clinical trials. Control Clin Trials. 1986;7(3):177-88.

667 25. DerSimonian R, Laird N. Meta-analysis in clinical trials revisited. Contemp Clin Trials.
668 2015;45(Pt A):139-45.

669 26. Mheissen S, Khan H, Normando D, Vaiid N, Flores-Mir C. Do statistical heterogeneity methods
670 impact the results of meta- analyses? A meta epidemiological study. PLoS One. 2024;19(3):e0298526.

671 27. Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K. cluster: Cluster Analysis Basics and
672 Extensions. R package version 2.1.0. [Available from: <u>https://cran.r-</u>
673 project.org/web/packages/cluster/index.html]. 2019.

Pearson ES. The choice of statistical tests illustrated on the interpretation of data classed in a
2× 2 table. Biometrika. 1947;34(1/2):139-67.

676 29. Campbell I. Chi-squared and Fisher-Irwin tests of two-by-two tables with small sample 677 recommendations. Stat Med. 2007;26(19):3661-75.

67830.Simonsohn U. Just post it: the lesson from two cases of fabricated data detected by statistics679alone. Psychol Sci. 2013;24(10):1875-88.

680 31. Bolland MJ, Avenell A, Gamble GD, Grey A. Systematic review and statistical analysis of the 681 integrity of 33 randomized controlled trials. Neurology. 2016;87(23):2391-402.

32. Hunter KE, Seidler AL, Askie LM. Prospective registration trends, reasons for retrospective
registration and mechanisms to increase prospective registration compliance: descriptive analysis and
survey. BMJ Open. 2018;8(3):e019983.

685 33. Harriman SL, Patel J. When are clinical trials registered? An analysis of prospective versus
686 retrospective registration. Trials. 2016;17:187.

Azar M, Riehm KE, Saadat N, Sanchez T, Chiovitti M, Qi L, et al. Evaluation of Journal
Registration Policies and Prospective Registration of Randomized Clinical Trials of Nonregulated
Health Care Interventions. JAMA Intern Med. 2019;179(5):624-32.

35. Hamilton DG, Hong K, Fraser H, Rowhani-Farid A, Fidler F, Page MJ. Prevalence and predictors
of data and code sharing in the medical and health sciences: systematic review with meta-analysis of
individual participant data. BMJ. 2023;382:e075767.

693 36. Brown NJ, Heathers J. Rounded Input Variables, Exact Test Statistics (RIVETS). 2019.

694 37. BOX GEP. NON-NORMALITY AND TESTS ON VARIANCES. Biometrika. 1953;40(3-4):318-35.

38. Davey J, Turner RM, Clarke MJ, Higgins JP. Characteristics of meta-analyses and their
component studies in the Cochrane Database of Systematic Reviews: a cross-sectional, descriptive
analysis. BMC Med Res Methodol. 2011;11:160.

Synthesis (RAISE): guidance and recommendations (Draft for consultation and revision) Open Science
Framework2024 [Available from: <u>https://osf.io/cn7x4</u>.

- 701
- 702
- 703
- 704
- 705 Figure legends
- Figure 1: Responses to trustworthiness checks in four domains
- 707 S Figure 1: Study-level responses to trustworthiness checks

S Figure 2: Dendrogram displaying hierarchical clustering of checks using completeagglomeration based on Gower dissimilarity

710

Inspecting results in the paper (28 checks)	Inspecting the research team and their work (19 checks)
1. Are any baseline data implausible with respect to magnitude, frequency, or variance?	29. Check whether withdrawal and loss to follow-up in multiple trials by the same author are
Is the number of participant withdrawals compatible with the disease, age and timeline?	consistent with the expected (random) binomial distribution
3. Are subgroup means incompatible with those for the whole cohort?	30. Are contributorship statements present?
4. Are the reported summary data compatible with the reported range?	31. Are contributorship statements complete?
5. Are correct units reported?	32. Have the data been published elsewhere by the research team in an illegitimate fashion?
6. Are calculations of proportions and percentages correct?	33. Are duplicate-reported data consistent between publications?
Are numbers of participants correct and consistent throughout the publication?	34. Are relevant methods consistent between publications?
8. Are there any discrepancies between data reported in figures, tables and text?	35. Is any duplicate reporting acknowledged or explained?
9. Are any outcome data, including estimated treatment effects, implausible?	36. Is there evidence of duplication of figures?
10. Are baseline statistical tests correct?	37. Does the statistics methods section use generic language, suggesting lack of expert
11. Are statistical tests of outcomes correct?	statistical input?
12. Are differences in variances in baseline variables between randomised groups plausible?	38. Is the distribution of non-first digits in manuscripts from one author compatible with a
13. Are any of the baseline data excessively similar between randomized groups?	genuine measurement process?
14. Are any of the baseline data excessively different between randomised groups?	39. Does consideration of other studies from members of the research team highlight causes
15. Are the summary outcome data identical or nearly identical across study groups?	for concern (including expressions of concern, relevant post-publication amendment, or
16. Are there any discrepancies between the values for percentage and absolute change?	critical retraction)?
17. Are there any discrepancies between reported data and participant inclusion criteria?	40. Is the standard deviation of summary statistics in multiple studies by same authors
18. Are the variances in biological variables surprisingly consistent over time?	plausible (when compared to simulated or bootstrapped data?)
19. Are results internally consistent?	41. Do all authors meet criteria for authorship?
20. Are coefficients of variation unusually similar when calculated across variables reported in the paper?	42. Is authorship of related papers consistent?
21. Is the amount of missing data plausible?	43. Are the authors on staff of institutions they list?
22. Are the results substantially divergent from the results of multiple other studies in meta-analysis?	44. Do any authors have a professorial title but no other publications on PubMed?
23. Are non-first digits compatible with a genuine measurement process?	45. Can co-authors attest to the reliability of the paper?
24. Are the variances of integer data possible?	46. Given the nature of the study, does the author list make sense? - e.g.does a simple study
25. Are the means of integer data possible?	have dozens of authors from different institutions and with diverse expertise?
26. Is there heterogeneity across studies in degree of imbalance in baseline characteristics (in meta-analysis) – only once per review	47. In which country was the study conducted?
27. Are integer data simulated from reported summary statistics plausible?	
28. Are important features missing from the paper?	
	I

712 Table 1: Trustworthiness checks in the first and second domains of the assessed list

Inspecting text and publication details (7 checks)
70. Are there typographical errors?
71. Has the study been retracted or does it have an expression of concern, a relevant post-publication
amendment, a critical Retraction Watch or PubPeer comment or has been previously excluded from a
systematic review?
72. Is there evidence of copied work, such as duplicated or partially duplicated tables?
73. Is there evidence of text reuse (cutting and pasting text between papers), including text that is
inconsistent with the study?
74. Is there evidence of automatically-generated text?
75. Was the study published in journal from a list of predatory/ low quality journals?
76. Is there evidence of manipulation or duplication of images?

720 Table 2: Trustworthiness checks in the third and fourth domains of the assessed list

Review-level summary (n = 50)

Number of RCTs in assessed meta-analysis	1	27 (54%)
	2	10 (20%)
	3	8 (16%)
	4	2 (4%)
	5	2 (4%)
	6*	1 (2%)
Number of participants in assessed meta-analysis	Median (IQR)	147 (53 to 341)
Outcome Type	Binary	26 (52%)
	Continuous	24 (48%)
Year of publication	2023	27 (54%)
	2022	10 (20%)
	2021	3 (6%)
	2020	5 (10%)
	2019	4 (8%)
	2014	1 (2%)
GRADE assessment	High	3 (6%)
	Moderate	7 (14%)
	Low	25 (50%)
	Very low	15 (30%)

10000		

727	Table 3: Characteristics for 50 Cochrane Reviews assessed in the study. Frequency (%) or median (1 st quartile to 3 rd quartile)
728	*Assessed in error, included in analysis.
729	
730	

Inspecting results in the paper (28 checks)	Inspecting the research team and their work (19 checks)
1. Are any baseline data implausible with respect to magnitude, frequency, or variance?	29. Check whether withdrawal and loss to follow-up in multiple trials by the same author are
2. Is the number of participant withdrawals compatible with the disease, age and timeline?	consistent with the expected (random) binomial distribution
3. Are subgroup means incompatible with those for the whole cohort?	30. Are contributorship statements present?
4. Are the reported summary data compatible with the reported range?	31. Are contributorship statements complete?
5. Are correct units reported?	32. Have the data been published elsewhere by the research team in an illegitimate fashion?
6. Are calculations of proportions and percentages correct?	33. Are duplicate-reported data consistent between publications?
7. Are numbers of participants correct and consistent throughout the publication?	34. Are relevant methods consistent between publications?
8. Are there any discrepancies between data reported in figures, tables and text?	35. Is any duplicate reporting acknowledged or explained?
9. Are any outcome data, including estimated treatment effects, implausible?	36. Is there evidence of duplication of figures?
10. Are baseline statistical tests correct?	37. Does the statistics methods section use generic language, suggesting lack of expert
11. Are statistical tests of outcomes correct?	statistical input?
12. Are differences in variances in baseline variables between randomised groups plausible?	38. Is the distribution of non-first digits in manuscripts from one author compatible with a
13. Are any of the baseline data excessively similar between randomized groups?	genuine measurement process?
14. Are any of the baseline data excessively different between randomised groups?	39. Does consideration of other studies from members of the research team highlight causes
15. Are the summary outcome data identical or nearly identical across study groups?	for concern (including expressions of concern, relevant post-publication amendment, or
16. Are there any discrepancies between the values for percentage and absolute change?	critical retraction)?
17. Are there any discrepancies between reported data and participant inclusion criteria?	40. Is the standard deviation of summary statistics in multiple studies by same authors
18. Are the variances in biological variables surprisingly consistent over time?	plausible (when compared to simulated or bootstrapped data?)
19. Are results internally consistent?	41. Do all authors meet criteria for authorship?
20. Are coefficients of variation unusually similar when calculated across variables reported in the paper?	42. Is authorship of related papers consistent?
21. Is the amount of missing data plausible?	43. Are the authors on staff of institutions they list?
22. Are the results substantially divergent from the results of multiple other studies in meta-analysis?	44. Do any authors have a professorial title but no other publications on PubMed?
23. Are non-first digits compatible with a genuine measurement process?	45. Can co-authors attest to the reliability of the paper?
24. Are the variances of integer data possible?	46. Given the nature of the study, does the author list make sense? - e.g.does a simple study
25. Are the means of integer data possible?	have dozens of authors from different institutions and with diverse expertise?
26. Is there heterogeneity across studies in degree of imbalance in baseline characteristics (in meta-analysis) – only once	47. In which country was the study conducted?
per review	
27. Are integer data simulated from reported summary statistics plausible?	
28. Are important features missing from the paper?	

Table 1: Trustworthiness checks in the first and second domains of the assessed list

Inspecting conduct, governance, and transparency (22 checks)	Inspecting text and publication details (7 checks)
48. Is the grant funding number identical to the number in unrelated studies?	70. Are there typographical errors?
49. Is a funding source reported?	71. Has the study been retracted or does it have an expression of concern, a relevant post-publication
50. Is the volume of work reported by research group plausible, including that indicated by concurrent studies from the same group?	amendment, a critical Retraction Watch or PubPeer comment or has been previously excluded from a
51. Is the reported staffing adequate for the study conduct as reported?	systematic review?
52. Is the recruitment of participants plausible within the stated time frame for the research?	72. Is there evidence of copied work, such as duplicated or partially duplicated tables?
53. Is the recruitment of participants plausible considering the epidemiology of the disease in the area of the study location?	73. Is there evidence of text reuse (cutting and pasting text between papers), including text that is
54. Is the interval between study completion and manuscript submission plausible?	inconsistent with the study?
55. Is there evidence that the work has been approved by a specific, recognized committee? (ethics)	74. Is there evidence of automatically-generated text?
56. Are there any concerns about unethical practice?	75. Was the study published in journal from a list of predatory/ low quality journals?
57. Could the study plausibly be completed as described?	76. Is there evidence of manipulation or duplication of images?
58. Are the study methods plausible, at the location specified?	
59. Are the locations where the research took place specified, and is this information plausible?	
60. Do the authors agree to share individual participant data?	X
61. Are the data publically available?	
62. Are additional patient data recorded in patient case records beyond what is reported in the paper?	
63. Does the trial registration number refer to other studies?	
64. Has the study been prospectively registered?	
65. Are details such as dates and study methods in the publication consistent with those in the registration documents?	
66. Do authors cooperate with requests for information?	K
67. Do authors provide satisfactory responses to requests?	
68. Was the time between submission to acceptance reasonable?	
69. Is the procedure of the study aligned with local legislations?	

Table 2: Trustworthiness checks in the third and fourth domains of the assessed list

Review-level summary (n = 50)			
	_	/ /	
Number of RCTs in assessed meta-analysis	1	27 (54%)	
	2	10 (20%)	
	3	8 (16%)	
	4	2 (4%)	
	5	2 (4%)	
	6*	1 (2%)	
Number of participants in assessed meta-analysis	Median (IQR)	147 (53 to 341)	

ournal Pre-proof	
Binary	26 (52%)
Continuous	24 (48%)
2023	27 (54%)
2022	10 (20%)
2021	3 (6%)
2020	5 (10%)
2019	4 (8%)
2014	1 (2%)
High	3 (6%)
Moderate	7 (14%)
Low	25 (50%)
Very low	15 (30%)
	BinaryContinuous202320222021202020192014HighModerateLowVery low

Table 3: Characteristics for 50 Cochrane Reviews assessed in the study. Frequency (%) or median (1st quartile to 3rd quartile)

*Assessed in error, included in analysis.





100

Inspecting conduct, governance & transparency



Cluster Dendrogram

Declarations

JW, CH, GAA, LB, JJK declare funding from NIHR (NIHR203568) in relation to the current project. JW additionally declares Stats or Methodological Editor roles for BJOG, Fertility and Sterility, Reproduction and Fertility, Journal of Hypertension, and for Cochrane Gynaecology and Fertility. CH declares a Statistical Editor role for Cochrane Colorectal. GAA additionally declares a Statistical Reviewer role for the European Journal of Vascular and Endovascular Surgery. LB additionally declares a role as Academic Meta-Research Editor for PLoS Biology, and that The University of Colorado receives remuneration for service as Senior Research Integrity Editor, Cochrane. JJK additionally declares a Statistical Editor role for The BMJ. EF is employed by the Cochrane Collaboration and on the Editorial Board of Cochrane Evidence Synthesis and Methods. ES is a Senior (Sign-off) Editor for Cochrane and contributed to Cochrane's Policy for managing potentially problematic studies. SL is an editor for Cochrane Gynaecology and Fertility, Human Reproduction, and Fertility and Sterility. TJL is the Deputy Editor in Chief of The Cochrane Library and is an employee of The Cochrane Collaboration. DNB is an associate editor for Research Quarterly for Exercise and Sport and a section editor for Communications in Kinesiology. NEO is a member of the Cochrane Editorial Board and holds an ERA-NET Neuron Co-Fund grant for a separate project. RR declares acting as an author and editor on Cochrane reviews. KS is an editor for Cochrane Gynaecology and Fertility, and Fertility and Sterility.MvW declares to be coordinating editor for Cochrane Gynaecology and Fertility and Cochrane Sexually Transmitted Infections, methodological editor for Human Reproduction Update and Editorial Editor for Fertility & Sterility. HT is Deputy Editor of The Lancet Gastroenterology & Hepatology and is an employee of Elsevier. SL received funding from the French National Research Agency (ANR-23-CE36-0006-01). AK is an editorial board member for BJGP Open. TLi serves as the Principal Investigator on a grant from the National Eye Institute, National Institutes of Health that funds the work of Cochrane Eyes and Vision US Project. She also acts as a sign-off editor for The Cochrane Library. ZM is supported by an NHMRC Investigator Grant 1195676. ZM is an associate Editor for BMC Medical Research Methodology and is on the Editorial Board for Clinical and Public Health Guidelines. RC is Editor-in-Chief at Meta-Psychology. CL is a work-package leader for the doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101120360), funded by the EU. CV received funding as part of the OSIRIS project (Open Science to Increase Reproducibility in Science); the OSIRIS (Open Science to Increase Reproducibility in Science) project has received funding from the EU (grant agreement No. 101094725). FN received funding from the French National Research Agency (ANR-23-CE36-0006-01), the French ministry of health and the French ministry of research. He is a work-package leader in the OSIRIS project (Open Science to Increase Reproducibility in Science). The OSIRIS project has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101094725. He is a work-package leader for the doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101120360), funded by the EU. DN declares having led/co-authored/co-authoring Cochrane Reviews. He also declares having been part of the Cochrane Convenes initiative organised by Cochrane to consider the issue of misinformation, its impact on the health evidence ecosystem and solutions to address it. LJ is the creator of the scrutiny package in R. WL is supported by an NHMRC Investigator grant (GNT2016729). RW is supported by an NHMRC Investigator Grant (2009767) and acts as a Deputy Editor for Human Reproduction, and an editorial board member for BJOG and Cochrane Gynaecology and Fertility. EF, SGTLa and RR declare employment by Cochrane. TLa additionally declares authorship of a chapter in the Cochrane Handbook for Systematic Reviews of Interventions and that he is a developer of standards for Cochrane intervention reviews (MECIR). AL is on the editorial board of BMC Medical Ethics.

What is new

- An extensive list of potential checks for assessing study trustworthiness was assessed via an application to 95 randomised controlled trials (RCTs) in 50 Cochrane Reviews.
- Following application of the checks, assessors had concerns about the authenticity of 32% of the RCTs.
- If these RCTs were excluded, 22% of meta-analyses would have no remaining RCTs.
- However, the study showed that some checks were frequently infeasible, and others could be easily misunderstood or misinterpreted.
- The study restricted assessment to meta-analyses including five or fewer RCTs, which might distort the impact of applying the checks.

bulleton