

Quantitative assessment of inconsistency in meta-analysis using decision thresholds with two new indices

Bernardo Sousa-Pinto^{1,2}, Ignacio Neumann³, Rafael José Vieira^{1,2}, Antonio Bognanni^{4,5,6}, Manuel Marques-Cruz^{1,2}, Sara Gil-Mata^{1,2}, Simone Mordue⁷, Clareece Nevill⁸, Gianluca Baio⁹, Paul Whaley^{10,11}, Guido Schwarzer¹², James Steele¹³, Gavin Stewart¹⁴, Holger J Schünemann^{6,15*}, Luís Filipe Azevedo^{1,2}

1 – MEDCIDS - Department of Community Medicine, Information and Health Decision Sciences; Faculty of Medicine, University of Porto, Porto, Portugal

2 – CINTESIS@RISE - Health Research Network, MEDCIDS, Faculty of Medicine, University of Porto, Porto, Portugal

3 – School of Medicine, Universidad San Sebastián, Santiago, Chile

4 – Department of Health Research Methods, Evidence and Impact, McMaster University, Hamilton, Ontario, Canada

5 – Department of Medicine, Evidence in Allergy Group, McMaster University, Hamilton, Canada

6 – Clinical and Epidemiology and Research Center (CERC), IRCCS Humanitas Research Hospital, & Department of Biomedical Sciences, Humanitas University, Milan, Italy

7 – Centre for Conservation Science, Cambridge, UK

8 – Department of Population Health Sciences, University of Leicester, Leicester, UK

9 – Department of Statistical Science, University College London, London, UK

10 – Lancaster Environment Centre, Lancaster University, Lancaster, UK

11 – Evidence-Based Toxicology Collaboration (EBTC), Johns Hopkins Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, USA

12 – Institute of Medical Biometry and Statistics, Faculty of Medicine and Medical Center – University of Freiburg, Freiburg, Germany

13 – Department of Sport and Health, Solent University, Southampton, UK

14 – Evidence Synthesis Lab, School of Natural and Environmental Science, University of Newcastle, Newcastle-upon-Tyne, UK

15 – Fraunhofer Institute for Translational Medicine and Pharmacology ITMP, Allergology and Immunology, Berlin, Germany

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
61
62
63
64
65

31 Correspondence to: *Holger J Schünemann; Clinical and Epidemiology and Research Center (CERC),
32 Department of Biomedical Sciences, Humanitas University & IRCCS Humanitas Research Hospital
33 Milan, Italy; E-mail address: schuneh@mcmaster.ca

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
61
62
63
64
65

34 **Abstract**

35 **Objective:** In evidence synthesis, inconsistency is typically assessed visually and with the I^2 and the Q
36 statistics. However, these measures have important limitations (i) if there are few primary studies of
37 small sample sizes, or (ii) if there are multiple studies with precise estimates. In addition, with the
38 increasing use of decision thresholds (DT), for example in GRADE Evidence to Decision frameworks,
39 inconsistency judgments can be anchored around DTs. In this article, we developed quantitative
40 measures to assess inconsistency based on DTs.

41 **Study Design and Setting:** We developed two measures to quantify inconsistency based on DTs – the
42 Decision Inconsistency (*DI*) and the Across-Studies Inconsistency (*ASI*) indices. The *DI* and the *ASI*
43 are based on the distribution of the posterior samples studies' effect sizes across interpretation
44 categories defined by DTs. We developed these indices for the Bayesian context, followed by a
45 frequentist extension.

46 **Results:** The *DI* informs on the *overall inconsistency of effect sizes* across interpretation categories,
47 while the *ASI* quantifies how *different* studies are compared to each other (in relation to interpretation
48 categories) based on absolute effects. A $DI \geq 50\%$ and an $ASI \geq 25\%$ are suggestive of important
49 unexplained inconsistency. We provide an R package (*metainc*) and a web tool
50 (<https://metainc.med.up.pt/>) to support the computation of the *DI* and *ASI*, including in the context of
51 sensitivity analyses assessing the impact of potential uncertainty in inconsistency.

52 **Conclusion:** The *DI* and the *ASI* can contribute to quantitatively assess inconsistency, particularly as
53 DTs are gaining recognition in evidence synthesis and health decision-making.

54

55 **Key words:** GRADE; Heterogeneity; Inconsistency; Meta-analysis; Systematic review

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
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75

56 **What is new**

57 Key findings

- 58 • This study proposes two new quantitative measures to assess inconsistency in the evidence
59 synthesis context – the Decision Inconsistency (*DI*) and the Across-Studies Inconsistency (*ASI*)
60 indices. These indices differ from previously existing measures by considering effect size in the
61 context of decision thresholds (DTs).
- 62 • We have developed a R package (metainc) and a web tool (<https://metainc.med.up.pt/>) to easily
63 support the computation of the *DI* and the *ASI*.

64 What this adds to what was known

- 65 • The GRADE working group posits that inconsistency judgments can be made considering DTs.
66 Our methods allow such judgments to be supported by quantitative indices' results.

67 What is the implication, what should change now?

- 68 • Quantitative assessment of inconsistency based on DTs is now possible. Therefore, judging
69 inconsistency when assessing the certainty of evidence may now consider the quantification of
70 the DT-related *DI* and *ASI*, alongside other approaches for appraising inconsistency.

71 **Highlights**

- 72 • GRADE assessments of inconsistency are facilitated by decision thresholds (DT)
- 73 • Two inconsistency indices have been developed to measure inconsistency based on DTs
- 74 • The new indices allow to assess the impact of uncertainty in the evidence on inconsistency

75 1. Introduction

76 In the context of evidence synthesis and appraisal, there are several methods to assess inconsistency
77 [1]. One possibility is the visual inspection of the forest plot, which provides a simple but subjective
78 approach. The Cochran's Q test allows for the calculation of a p -value, based on which it is possible to
79 reject (or not) the null hypothesis of no heterogeneity. However, it has low power in meta-analyses if
80 there are few primary studies and/or studies with small sample sizes, while exhibiting over-inflated
81 power to detect small amounts of heterogeneity in meta-analyses with a large number of primary
82 studies[2]. The I^2 value is frequently used to assess the relative extent of inconsistency. Nevertheless, it
83 also has important limitations, as it is influenced by the sample size of the included primary studies
84 (e.g., it may overestimate inconsistency across studies with precise estimates) and may yield biased
85 results when used in the context of small-sample meta-analyses [1, 3, 4]. In addition, both the Q-
86 Cochran test and the I^2 value are based on frequentist methods [5], potentially limiting their application
87 to the Bayesian context [6].

88 These classical inconsistency measures exclusively rely on statistical criteria. However, there may be
89 scenarios where concerns for apparently large statistical heterogeneity might be mitigated. In fact, the
90 GRADE approach uses four items to judge inconsistency, namely the I^2 , Cochran's Q test p -value,
91 overlap in confidence intervals of primary studies by visual inspection, and the degree of difference in
92 the point estimates of relative effects. To bring value to considered judgment, GRADE states that
93 guideline and systematic review developers may abstain from rating down the certainty of evidence
94 (CoE) if point estimates of primary studies are on the same side of a prespecified threshold (i.e., fall
95 within the same target of the certainty range), despite the evidence of statistical heterogeneity [1, 7].
96 This assessment depends on providing context to outcomes' interpretation, by defining health outcome-
97 level decision thresholds (DTs) – effect size measures suggesting whether an intervention translates to
98 trivial or no, small, moderate or large effects [8, 9]. In a Bayesian context, it is possible to directly assess
99 the proportion of effect sizes, sampled from the posterior distribution of the different primary studies,
100 falling into the different ranges (interpretations) defined by DTs. Based on that, and on the concept of
101 incorporating outcome-level DTs in inconsistency assessment, we developed two measures to support
102 the assessment of inconsistency in meta-analysis. While the concept for this approach has been
103 developed considering a Bayesian framework, it is also fully applicable to the frequentist context.

104 Given the limitations in existing approaches and their interpretation, our objective was to develop
105 measures to support assessments of inconsistency. In this article we describe the development and
106 application of two new measures : (i) one assessing overall outcome-level-related inconsistency (the
107 Decision Inconsistency index), and (ii) one assessing across-studies inconsistency (the Across-Studies
108 Inconsistency index). These measures are not intended to replace but rather to complement existing
109 approaches to appraise inconsistency. We will start by reviewing the concept of the Dissimilarity Index

110 as the foundation from which our approach is derived. We will provide the concepts and formulae for
111 the Decision Inconsistency Index and for the Across-Studies Inconsistency Index. Subsequently, we
112 will apply our proposed approach using two practical examples. We will then present a web app and the
113 metainc R package to implement the Decision Inconsistency Index and the Across-Studies
114 Inconsistency Index. Finally, we will discuss potential limitations of our approach and how it may
115 contribute to interpreting inconsistency in the GRADE CoE framework [10].

116 2. The Decision Inconsistency Index

117 2.1. The Dissimilarity Index

118 The Dissimilarity Index is one of the most commonly used demographics measures of segregation [11],
119 reflecting the relative distributions of two groups over a set of geographic units [12]. It ranges between
120 0 and 1, with 0 indicating perfect integration (i.e., each geographic unit has the same percentage of
121 members of each group as the total population) and 1 indicating maximum segregation (i.e., each
122 geographic unit exclusively includes members of one of the two groups [11]) (Supplementary Figure
123 1). The formula for the computation of the Dissimilarity Index is the following [12, 13]:

$$124 \text{Dissimilarity ind} = \frac{1}{2} \sum_i^n \left| \frac{N_{1i}}{N_1} - \frac{N_{2i}}{N_2} \right|$$

125 with n corresponding to the number of geographic units, N_{1i} corresponding to the population of group
126 1 in the geographical unit i , N_1 corresponding to the total population of group 1 in all considered
127 geographical units, N_{2i} corresponding to the population of group 2 in the geographical unit i , and N_2
128 corresponding to the total population of group 2 in all considered geographical units being.

129 2.2. The Decision Inconsistency Index

130 Consider a meta-analysis including k primary studies comparing an intervention I versus a comparator
131 C on a certain outcome whose reduction would correspond to a benefit and whose increase would
132 correspond to a harm. For that outcome, an outcome-level DT has been established so that:

- 133 • If the effect size (ES) $> DT$, I would be associated with at least small harms.
- 134 • If $ES < -DT$, I would be associated with at least small benefits;
- 135 • If $-DT \leq ES \leq DT$, I would be associated with a trivial or no effect (henceforth referred to as
136 “trivial effect”);

137 Therefore, in this scenario, we consider three interpretation categories (at least small benefits, trivial or
138 no effects, and at least small harms) for ES, with two DT ($-DT$ and DT). We consider ES to be presented
139 as absolute effects, as recommended by the GRADE working group for contextualizing ES in relation
140 to DTs.

141 In a Bayesian meta-analytical context with no overall decision-related inconsistency, all samples from
 142 the posterior distributions of the ES of included primary studies will have their values associated with
 143 the same interpretation (either at least small benefits, at least small harms, or trivial effect). On the other
 144 extreme, if there is complete inconsistency, there will be a perfectly even distribution of the posterior
 145 ES samples across the three interpretation categories. That is, one third of the samples will indicate at
 146 least small benefits, another third will indicate at least small harms and the final third will indicate a
 147 trivial effect. Therefore, a situation of no clinical inconsistency would be analogous to one with full
 148 segregation in the geographical context (Dissimilarity Index=1), while a situation with maximum
 149 inconsistency would be analogous to one with full integration (Dissimilarity Index=0). However, in
 150 contrast with the Dissimilarity Index, we do not compare two groups but rather one distribution of
 151 posterior samples for the ES (across interpretation categories) with the expected distribution that would
 152 have been observed if there was maximum inconsistency. This concept forms the basis of a novel
 153 measure of inconsistency we propose – the Decision Inconsistency Index (*DI*). The *DI* quantifies overall
 154 inconsistency from a decision point of view, and may be calculated by:

$$DI = 1 - \left(\frac{\frac{1}{2} \sum_j \left| \frac{N_j - 1}{N - J} \right|}{\frac{J-1}{J}} \right)$$

155 with N_j corresponding to the number of ES posterior samples per interpretation category, N
 156 corresponding to the total number of ES samples (i.e., the sum, for all primary studies, of all study-level
 157 posterior samples), and J corresponding to the number of interpretation categories.

158 Dividing by $\left(\frac{J-1}{J}\right)$ ensures that the *DI* lies between 0 and 1, while subtracting the ratio from 1 ensures
 159 that higher values are associated with higher inconsistency (as with the I^2 value). That is, the *DI* ranges
 160 between 0 and 1 (or 0-100%, if multiplied by 100), with 0 indicating no DT-related inconsistency and
 161 1 indicating maximum DT-related inconsistency.

162 The *DI* can be calculated for as many interpretation categories as desired. If only two DTs are being
 163 considered (i.e., DT distinguishing trivial effects from at least small benefits and DT distinguishing
 164 trivial effects from at least small harms), three interpretation categories are possible (at least small
 165 benefits, at least small harms, and trivial effect), the formula of the *DI* can be stated as:

$$DI_{[2 \text{ DTs}]} = 1 - \left(\frac{\frac{1}{2} \sum_j \left| \frac{N_j - 1}{N - 3} \right|}{\frac{2}{3}} \right)$$

166 However, when using the GRADE Evidence to Decision (EtD) framework, decision-makers are usually
 167 interested in knowing not only whether an intervention is associated with non-trivial effects but also
 168 their magnitude (i.e., whether the interventions' desirable and undesirable health effects are trivial or
 169 none, small, moderate, or large) [14, 15]. The *DI* can be applied to these situations with three DTs on

172 each side of the no-effect (i.e., three DTs for benefits and three for harms) and, therefore, seven
173 interpretation categories as recommended by GRADE [14, 15]. In this case, the formula of the *DI* can
174 be stated as:

$$175 \quad DI_{[6 \text{ DTs}]} = 1 - \left(\frac{\frac{1}{2} \sum_j \left| \frac{N_j - 1}{N} - \frac{1}{7} \right|}{\frac{6}{7}} \right)$$

176 **3. The Across-Studies Inconsistency index**

177 Although the *DI* provides a measure of the overall inconsistency of ES across interpretation categories,
178 it does not quantify how inconsistent the ES of primary studies are when compared with each other. Let
179 us consider the examples depicted in Figure 1. Figure 1A provides an example in which all ES samples
180 of primary studies point to a trivial effect. The *DI* would be of 0%. Figure 1B and Figure 1C are two
181 examples for which one third of ES samples suggest important benefits, one third points to important
182 harms, and one third to a trivial effect. Therefore, in both examples, the *DI* would be expected to be
183 large. However, while in Figure 1B all primary studies display a similar proportion of ES samples in
184 each decision category, in Figure 1C the different primary studies display a different proportion of ES
185 samples in each decision category (e.g., the first primary study would be expected to have most samples
186 suggesting at least small harms while the third would be expected to have most samples suggesting at
187 least small benefits). Therefore, the examples depicted in Figure 1B and 1C would differ on across-
188 studies inconsistency, which would be larger in the latter.

189 Therefore, as a complement to the *DI*, we suggest that across-studies inconsistency should also be
190 measured. For this purpose, Dissimilarity Index-based measures would not be suitable, as (i) the
191 Dissimilarity Index has been devised to consider two groups, and (ii) it displays important limitations
192 when dealing with small unit sizes (the Dissimilarity Index is extremely sensitive to small differences
193 in cases when a small number of observations falls within a certain category)[16]. The adjustments
194 proposed to address this limitation cannot be applied to measures generalizing the Dissimilarity Index
195 to more than two groups[16, 17].

196 Therefore, to assess across-studies inconsistency, we propose a measure comparing (i) the observed
197 number of samples per interpretation category for each study with (ii) the expected number of samples
198 (per interpretation category for each study) if the proportion of samples per interpretation category had
199 been the same for all studies and equal to the overall proportion. Potentially adequate candidates for
200 such measures would be, for example, those based on the chi-squared statistic (χ^2), particularly relative
201 to the maximum value (in order to allow for obtaining a measure ranging between 0 and 1). Given that,
202 in this context, χ^2 would be given by:

$$\chi^2 = \sum \frac{(N_{ij} - E_{ij})^2}{E_{ij}}$$

(with N_{ij} corresponding to the number of ES samples per interpretation category and primary study, E_{ij} corresponding to the expected number of ES samples per interpretation category and primary study), and that the maximum of the chi-square statistic ($\max \chi^2$) would be given by:

$$\max \chi^2 = N(\min(J, k) - 1),$$

(with J corresponding to the number of interpretation categories and k to the number of primary studies), then the Across-Studies Inconsistency Index (*ASI*) would be given by:

$$ASI = \sqrt{\frac{\sum \frac{(N_{ij} - E_{ij})^2}{E_{ij}}}{N(\min(J, k) - 1)}}$$

The *ASI* measures across-study inconsistency considering decision interpretation categories. Its values can range between 0 and 1 (or 0-100%, if multiplied by 100), with 0 indicating no across-studies inconsistency and 1 indicating complete across-studies inconsistency.

4. Extension to the frequentist context

The *DI* and the *ASI* have been developed as measures using Bayesian meta-analysis. Indeed, the fact that Bayesian models yield a posterior distribution for the ES measures of primary studies renders that context particularly suited for the computation of the *DI* and *ASI*. However, the *DI* and *ASI* can also be calculated in the frequentist context. To accomplish that, one first needs to obtain the best linear unbiased predictions of the ES of each primary study, which can be calculated for frequentist random effects models empirically utilizing the obtained estimate for the between-study variability [τ^2] (which is analogous to the posterior ES estimates obtained in the Bayesian context). Subsequently, using these values and the corresponding standard-errors, it is possible to fit a probability distribution for each primary study, based on which samples can be drawn. The *DI* and *ASI* can then be calculated – in a similar way as in the Bayesian context – based on the sampled values for the best linear unbiased predictions of the ES of each primary study. Given that this is not straightforward, we developed an application for meta-analysts (see below).

5. Practical examples and application in sensitivity analyses

The Supplement displays two examples in which the *DI* and *ASI* are calculated. Supplementary Example 1 involves the computation of these indices in what was Cochrane's first living systematic review and meta-analysis of 18 primary studies comparing heparin with placebo on mortality at 12 months in a population of ambulatory patients with cancer [18] (Supplementary Tables 1-2). In brief,

232 this meta-analysis indicated that heparin was associated with decreased odds of mortality (5 fewer
233 deaths per 1000 individuals; 95% credible interval: between 49 fewer deaths and 28 more deaths per
234 1000 individuals), with a I^2 value suggesting moderate inconsistency ($I^2=35.2\%$). Computing the DT-
235 based inconsistency indices, we would obtain a DI of 80.1% and a ASI of 19.3%, pointing to the
236 possibility of relevant inconsistency (see below). The DI reflects the wide spread of posterior samples
237 of ES across interpretation categories defined by information sizes, pointing to the potential difficulty
238 in judging the size of the effect associated with the use of heparin.

239 We explored approaches for conducting sensitivity analyses, including (i) leave-one-out sensitivity
240 analysis (removing each primary study at once and recalculating the DI and the ASI), (ii) sensitivity
241 analysis based on risk of bias, (iii) sensitivity analysis based on uncertainty in baseline risk, and (iv)
242 sensitivity analysis based on DTs. Overall, leave-one-out meta-analysis did not allow us to identify any
243 individual primary study largely responsible for the observed inconsistency (Supplementary Table 3).
244 The sensitivity analysis based on the risk of bias suggested that, in this example, inconsistency may be
245 higher for studies displaying a lower risk of bias ($DI=83.8\%$; $ASI=23.2\%$) than for those with a high
246 risk of bias ($DI=77.9\%$; $ASI=18.4\%$), but the difference was small. Finally, an increase in the baseline
247 risk was associated with a decreasing trend for the DI (Supplementary Figure 2); the results of the
248 sensitivity analysis based on DTs are displayed in Supplementary Figure 3.

249 Supplementary Example 2 computes the DI and ASI in a sample of 100 published meta-analyses of
250 outcomes with available DTs [19, 20] (Supplementary Table 4). We observed that, in our sample, the
251 median DI and ASI values were of 32% and 19%, respectively. On the other hand, the second tertile
252 values are close to $DI \geq 50\%$ and $ASI \geq 25\%$ (Table 1). Considering only the DT of going from trivial
253 or no to a small effect (sometimes consistent with the “minimal important difference”) instead of three
254 decision thresholds on each side of the null effect (trivial or no to small, small to moderate and moderate
255 to large effects) did not produce a predictable effect on the DI (mean difference of 0.4 percent points)
256 but was associated with an increase in the ASI (mean difference of 7.5 percent points) (Supplementary
257 Figure 4).

258 **6. Implementation in practice**

259 We have developed an online app allowing for the computation of the DI and the ASI . The app, which
260 is available at <https://metainc.med.up.pt/>, takes as input a dataset containing the ES and the variance
261 for each primary study. Based on the provided input, it can perform either frequentist meta-analysis
262 using meta, or Bayesian meta-analysis using brms. Based on the posterior samples of the ES measures
263 of the primary studies, the app provides information on the DI and on the ASI . Users of this online app
264 are not required to (i) conduct a Bayesian meta-analysis or have knowledge on how to do it, or (ii) have
265 knowledge on how to obtain the best linear unbiased predictions of the ES of primary studies in the

266 context of frequentist meta-analysis. Therefore, the app allows to overcome potential barriers for the
267 computation of the *DI* and the *ASI*.

268 We have also developed a R package – *metainc* – to assess overall decision-related inconsistency and
269 across-studies inconsistency (computing the *DI* and the *ASI*, respectively) after performing Bayesian
270 or frequentist meta-analysis. It is available on CRAN ([https://cran.r-](https://cran.r-project.org/web/packages/metainc/index.html)
271 [project.org/web/packages/metainc/index.html](https://cran.r-project.org/web/packages/metainc/index.html)). Supplementary Boxes 1-2 provide information and
272 guided examples on how to use the *metainc* package.

273

274

275 7. Discussion

276 In this paper, we propose a quantitative approach to assess inconsistency in the meta-analytical context
277 using DTs. This approach involves the computation of the *DI* and the *ASI*, which provide
278 complementary information – the *DI* informs about the overall inconsistency of ES across interpretation
279 categories, while the *ASI* quantifies across-studies inconsistency. The proposed measures have been
280 developed in the Bayesian context, but they can also be computed for frequentist meta-analysis.

281 The GRADE guidance states that the assessment of inconsistency should not solely rely on classical
282 measures of heterogeneity, as they have statistical limitations [1]. However, the current guidance for
283 assessing inconsistency beyond those measures is centred on the visual inspection of the forest plot and
284 plausibility of subgroup analyses. While our proposed methods are not intended to replace other
285 approaches, they could provide valuable complementary information. In particular, by making use of
286 DTs, the proposed methods can help interpreting the importance of observed inconsistency, something
287 which is in line with recent statements to move away from interpreting results solely based on statistical
288 criteria [21-23]. As an example, the inconsistency indices and their use of DTs can help identify
289 situations in which (i) across-study inconsistency would impact the importance of findings (i.e., by
290 providing a formal degree of contextualization), or (ii) statistical measures of heterogeneity may be
291 overestimating inconsistency (e.g., due to primary studies with high precision estimates). Therefore,
292 our proposed approach can be applied not only within GRADE, but also in the context of any meta-
293 analysis, in order to help interpreting and framing the observed inconsistency.

294 Importantly, the proposed methods allow for sensitivity analyses based on the uncertainty in baseline
295 risk or DTs. This accounts for potential uncertainty in the baseline risk or DTs when assessing
296 inconsistency, which is not otherwise possible using only classical measures of heterogeneity or the
297 visual inspection of the forest plot. Accounting for uncertainty may be particularly relevant since
298 GRADE has proposed an approach to empirically obtain DTs that not only allows for the computation

299 of DT point estimates but also of best- and worst-case scenario DTs (Wiercioch et al, accepted pending
300 revision, [9]). Performing sensitivity analyses based on DTs allows assessing whether inconsistency
301 results are similar across a set of plausible DTs that can range between the best- and worst-case scenario
302 DT values. In addition, based on sensitivity analyses, it is also possible to appraise unexplained
303 inconsistency (i.e., inconsistency not explained by any pre-specified or convincing effect modifier)
304 using the *DI* and the *ASI*. This is particularly relevant, since GRADE recommends that judgments on
305 inconsistency are based on unexplained inconsistency.

306 Box 1 provides demonstrative examples (and a suggested reporting language) on how the *DI* and the
307 *ASI* can be used to support inconsistency judgements in the assessment of the CoE in the GRADE
308 approach. Although both the meta-analyses of the examples A and B display severe inconsistency as
309 assessed by the I^2 value (example A: $I^2=96\%$; example B: $I^2=69\%$), the DT-related indices suggest that
310 inconsistency may be at least a serious concern in example A ($DI=73\%$; $ASI=60\%$), but not in example
311 B ($DI=2\%$; $ASI=9\%$; despite quantitative differences, the effect sizes are pointing to large or moderate
312 benefits for all primary studies). While these examples illustrate a possible use of the *DI* and the *ASI*,
313 this paper does intend to provide guidance on how to judge inconsistency in the GRADE approach. This
314 is dealt with elsewhere including in the updated GRADE handbook (now called the GRADE Book
315 <https://book.gradepro.org/guideline/inconsistency> [10]). It will require a broad agreement on a
316 framework on how to consider different scenarios based on possible agreements and disagreements
317 between the different items related to inconsistency (i.e., visual inspection of the forest plot, statistical
318 measures of heterogeneity, and decision threshold-based inconsistency indices). This will also require
319 definite guidance when it is adequate to downgrade inconsistency by two or even three levels for
320 inconsistency. In addition, it will serve to adequately and jointly assess inconsistency and imprecision
321 in order to avoid double penalisation. The *DI* and *ASI* could support answering such questions by
322 providing information based on DTs in line with GRADE guidance (Wiercioch et al., accepted for
323 publication pending revision).

324 The proposed approach has some limitations. Firstly, the assessment of inconsistency based on DTs
325 should not be based solely on the calculation of the *DI* and the *ASI*. These measures should not be
326 understood as definite indicators of inconsistency but as additional tools to use when assessing this
327 domain. Another limitation concerns the absence of cut-off points defining low, moderate and severe
328 inconsistency in the context of the *DI* and the *ASI*. We have evaluated the distribution of these indices
329 in a sample of 100 meta-analyses, and this may provide suggestions to users on how to interpret their
330 *DI* and *ASI* results. That is, users may hint at the magnitude of their inconsistency by comparison with
331 the percentiles of *DI* and *ASI* for other systematic reviews. Nevertheless, setting of cut-off points for
332 claiming inconsistency may require a more comprehensive approach and may depend on the number of
333 considered decision thresholds. However, it should be noted that universally-agreed or even sensible
334 cut-off points do not exist for the I^2 either [2]. Exploration of other approaches to explore inconsistency

335 ratings may be useful and we are beginning work in a GRADE project group to provide this guidance.
336 Presenting only the proportions point estimates falling into within different certainty ranges may be one
337 alternative. However, it brings complexity of having to interpret many data points simultaneously,
338 requiring judgments with unknown reliability (particularly in meta-analyses with a small amount of
339 primary studies) and not considering the confidence intervals of the studies' estimates. Finally, not all
340 functions are currently available for non-R users. However, efforts are being made to increase the
341 number of functions accessible through different software or platforms.

342 **8. Conclusion**

343 In meta-analysis, the assessment of inconsistency based solely on classical measures of heterogeneity
344 has important limitations. Considering DTs may allow for that assessment in the respective health
345 decision context. However, no quantitative approaches had been proposed so far. In this paper, we
346 describe two measures – the *DI* and the *ASI* – that can be used to quantitatively assess inconsistency
347 using DTs. While their computation does not replace other methods for assessing inconsistency, used
348 together they can be particularly helpful for (i) interpreting the health importance of observed
349 inconsistency, (ii) giving the evaluator additional information about the impact of potential uncertainty
350 in baseline risk or DTs, and (iii) supporting the rating of inconsistency in the GRADE appraisal of the
351 CoE. Based on the developed R package and web tool, this approach can be easily implemented both
352 in the Bayesian and frequentist contexts.

353
354 **Acknowledgements:** The authors would like to thank the Evidence Synthesis Hackathon, whose 2023
355 edition allowed for relevant advances on this study and on the respective R package.

356 **Author contribution statement:**

- 357 • BSP has participated in conceptualization, data curation, formal analysis, methodology, and
358 writing - original draft;
- 359 • MMC and SGM have participated in formal analysis, and writing – review & editing;
- 360 • SM, CN, GB and PW have participated in methodology, and writing – review & editing;
- 361 • GuS, JMS and GaS have participated in data curation, formal analysis, methodology, and
362 writing – review & editing;
- 363 • IN, RJV, AB, HJS and LFA have participated in conceptualization, methodology, and writing -
364 review & editing.

365 **Declarations of interest:** HJS is co-chair of the GRADE Working Group, but this is not an official
366 GRADE Working Group article (although the concepts herein may be used by GRADE in the future
367 but this will require formal approval). All other authors declare no conflict.

368 **Funding sources:** This research did not receive any specific grant from funding agencies in the public,
1
2 369 commercial, or not-for-profit sectors.

3
4 370 **Generative AI and AI-assisted technologies in the writing process:** No generative AI or AI-assisted
5
6 371 technologies were used in the writing process or in any other task for this manuscript.

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
61
62
63
64
65

372 **References**

- 1
2
3 373 [1] Guyatt GH, Oxman AD, Kunz R, Woodcock J, Brozek J, Helfand M, et al. GRADE guidelines: 7. Rating
4 374 the quality of evidence--inconsistency. *J Clin Epidemiol*. 2011;64:1294-302.
- 5 375 [2] Deeks JJ, Higgins JP, Altman DG, Group CSM. Analysing data and undertaking meta-analyses.
6 376 *Cochrane handbook for systematic reviews of interventions*. 2019:241-84.
- 7 377 [3] von Hippel PT. The heterogeneity statistic I(2) can be biased in small meta-analyses. *BMC Med Res*
8 378 *Methodol*. 2015;15:35.
- 9
10 379 [4] Schunemann HJ, Neumann I, Hultcrantz M, Brignardello-Petersen R, Zeng L, Murad MH, et al.
11 380 GRADE guidance 35: update on rating imprecision for assessing contextualized certainty of evidence
12 381 and making decisions. *J Clin Epidemiol*. 2022;150:225-42.
- 13 382 [5] Higgins JP, Thompson SG. Quantifying heterogeneity in a meta-analysis. *Stat Med*. 2002;21:1539-
14 383 58.
- 15 384 [6] Welton NJ, Sutton AJ, Cooper N, Abrams KR, Ades A. Evidence synthesis for decision making in
16 385 healthcare: John Wiley & Sons; 2012.
- 17 386 [7] Guyatt G, Zhao Y, Mayer M, Briel M, Mustafa R, Izcovich A, et al. GRADE guidance 36: updates to
18 387 GRADE's approach to addressing inconsistency. *J Clin Epidemiol*. 2023;158:70-83.
- 19 388 [8] Morgano GP, Mbuagbaw L, Santesso N, Xie F, Brozek JL, Siebert U, et al. Defining decision thresholds
20 389 for judgments on health benefits and harms using the Grading of Recommendations Assessment,
21 390 Development and Evaluation (GRADE) Evidence to Decision (EtD) frameworks: a protocol for a
22 391 randomised methodological study (GRADE-THRESHOLD). *BMJ Open*. 2022;12:e053246.
- 23 392 [9] Morgano GP, Wiercioch W, Piovani D, Neumann I, Nieuwlaat R, Piggott T, et al. Defining decision
24 393 thresholds for judgments on health benefits and harms using the Grading of Recommendations
25 394 Assessment, Development and Evaluation (GRADE) Evidence to Decision (EtD) frameworks: a
26 395 randomised methodological study (GRADE-THRESHOLD). *J Clin Epidemiol*. 2024:111639.
- 27 396 [10] Neumann I, Sousa-Pinto B, Meerpohl J, Dahm P, Brennan S, Alonso-Coello P, et al. Inconsistency.
28 397 In: Neumann I, Schünemann H, editors. *GRADE Book*. <https://book.gradepro.org>: GRADE Working
29 398 Group; 2024.
- 30 399 [11] Massey DS, Rothwell J, Domina T. The Changing Bases of Segregation in the United States. *Ann Am*
31 400 *Acad Pol Soc Sci*. 2009;626.
- 32 401 [12] Massey DS, Denton NA. The dimensions of residential segregation. *Social forces*. 1988;67:281-
33 402 315.
- 34 403 [13] The Social Science Data Analysis Network. *CensusScope: ABOUT DISSIMILARITY INDICES*.
- 35 404 [14] Alonso-Coello P, Schunemann HJ, Moberg J, Brignardello-Petersen R, Akl EA, Davoli M, et al.
36 405 GRADE Evidence to Decision (EtD) frameworks: a systematic and transparent approach to making well
37 406 informed healthcare choices. 1: Introduction. *BMJ*. 2016;353:i2016.
- 38 407 [15] Alonso-Coello P, Oxman AD, Moberg J, Brignardello-Petersen R, Akl EA, Davoli M, et al. GRADE
39 408 Evidence to Decision (EtD) frameworks: a systematic and transparent approach to making well
40 409 informed healthcare choices. 2: Clinical practice guidelines. *BMJ*. 2016;353:i2089.
- 41 410 [16] Allen R, Burgess S, Davidson R, Windmeijer F. More reliable inference for the dissimilarity index of
42 411 segregation. *Econom J*. 2015;18:40-66.
- 43 412 [17] Sakoda JM. A generalized index of dissimilarity. *Demography*. 1981;18:245-50.
- 44 413 [18] Akl EA, Kahale LA, Hakoum MB, Matar CF, Sperati F, Barba M, et al. Parenteral anticoagulation in
45 414 ambulatory patients with cancer. *Cochrane Database Syst Rev*. 2017;9:CD006652.
- 46 415 [19] Cuker A, Tseng EK, Nieuwlaat R, Angchaisuksiri P, Blair C, Dane K, et al. American Society of
47 416 Hematology living guidelines on the use of anticoagulation for thromboprophylaxis in patients with
48 417 COVID-19: January 2022 update on the use of therapeutic-intensity anticoagulation in acutely ill
49 418 patients. *Blood Adv*. 2022;6:4915-23.
- 50 419 [20] Piggott T, Nonino F, Baldin E, Filippini G, Rijke N, Schunemann H, et al. Multiple Sclerosis
51 420 International Federation guideline methodology for off-label treatments for multiple sclerosis. *Mult*
52 421 *Scler J Exp Transl Clin*. 2021;7:20552173211051855.

422 [21] Hernán MA, Greenland S. Why Stating Hypotheses in Grant Applications Is Unnecessary. JAMA.
1 423 2024.
2 424 [22] Wasserstein RL, Lazar NA. The ASA statement on p-values: context, process, and purpose. Taylor
3 425 & Francis; 2016. p. 129-33.
4 426 [23] Wasserstein RL, Schirm AL, Lazar NA. Moving to a world beyond “ $p < 0.05$ ”. Taylor & Francis; 2019.
5 427 p. 1-19.
6 428
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
61
62
63
64
65

Tables

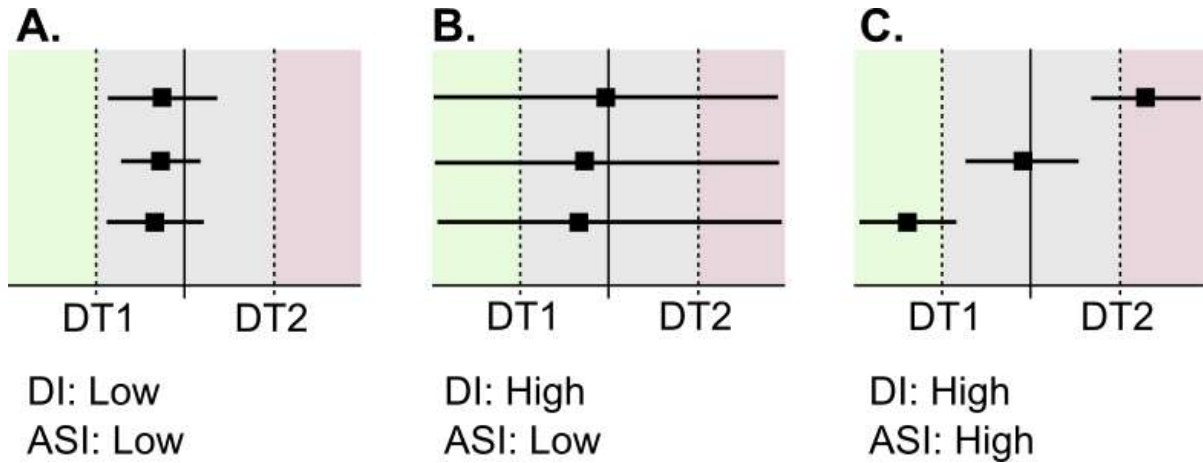
Table 1. Summary of the distributions of the Decision Inconsistency Index (DI) and Across-Studies Inconsistency Index (ASI) in a sample of 100 published meta-analysis

Inconsistency measure	Percentile 50 (median)	Percentile 67	Percentile 75	Maximum
DI (%)	32.2	49.1	58.3	86.2
ASI (%)	19.2	24.8	28.6	60.5

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
61
62
63
64
65

Figures

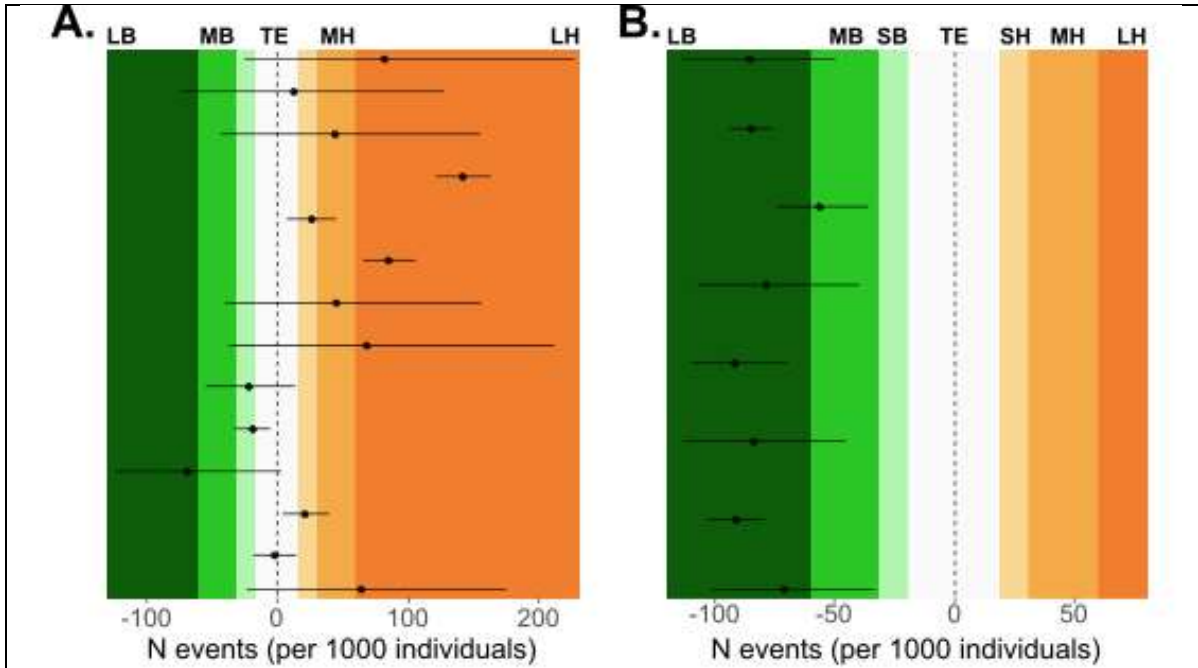
Figure 1. Hypothetical meta-analytical examples illustrating the differences in the concepts of Decision Inconsistency (DI) and Across-Studies Inconsistency (ASI). The green zone (left) indicates at least small benefits, the grey zone (centre) indicates a trivial or no effect, and the red zone (right) indicates at least small harms.



In Figure 1B, even though the point estimates of primary studies would all indicate a trivial effect, the DI would be high as there would be a large proportion of posterior samples also indicating at least small benefits and at least small harms (the effects of the primary studies are all compatible with at least small benefits, a trivial effect and at least small harms); however, all studies would be similar among themselves (hence, the ASI would be low). In Figure 1C, both the DI and the ASI would be high as (i) there would be a large proportion of posterior samples indicating at least small benefits, a trivial effect or at least small harms, and (ii) all studies would be very different among themselves. DT=Decision thresholds

Boxes

Box 1. Examples on how the Decision Inconsistency Index (*DI*) and the Across-Studies Inconsistency (*ASI*) indices can be used to support judgements on inconsistency in GRADE



Example A: We had at least serious concerns about inconsistency of the evidence. High inconsistency was suggested both by statistical measures of heterogeneity ($I^2 = 96\%$) and by threshold-based inconsistency indices ($DI = 73\%$; $ASI = 60\%$). No variable was identified that would potentially explain the inconsistency. We therefore rated down the certainty of evidence for inconsistency by at least 1 level.

Example B: We had no serious concerns about inconsistency of the evidence. Although high inconsistency was suggested by statistical measures of heterogeneity ($I^2 = 69\%$), decision threshold-based inconsistency indices suggested low inconsistency ($DI = 2\%$; $ASI = 9\%$). This shows that inconsistency may not be important. We therefore rated down the certainty of evidence for inconsistency by 0 levels.

Declaration of Interest Statement

Declarations of interest: HJS is co-chair of the GRADE Working Group, but this is not an official GRADE Working Group article (although the concepts herein may be used by GRADE in the future but this will require formal approval). All other authors declare no conflict.

Author contribution statement:

- BSP has participated in conceptualization, data curation, formal analysis, methodology, and writing - original draft;
- MMC and SGM have participated in formal analysis, and writing – review & editing;
- SM, CN, GB and PW have participated in methodology, and writing – review & editing;
- GuS, JMS and GaS have participated in data curation, formal analysis, methodology, and writing – review & editing;
- IN, RJV, AB, HJS and LFA have participated in conceptualization, methodology, and writing - review & editing.