

19 **Abstract**

20 Performance weighted aggregation of expert judgements, using calibration questions, has been
21 advocated to improve pooled quantitative judgements for ecological questions. However, there
22 is little discussion or practical advice in the ecological literature regarding the application,
23 advantages or challenges of performance weighting.

24 In this paper we 1) illustrate how the IDEA protocol with four-step question format can be
25 extended to include performance weighted aggregation from the Classical Model, and 2)
26 explore the extent to which this extension improves pooled judgements for a range of
27 performance measures.

28 Our case study demonstrates that performance weights can improve judgements derived from
29 the IDEA protocol with four-step question format. However, there is no *a-priori* guarantee of
30 improvement. We conclude that the merits of the method lie in demonstrating that the final
31 aggregation of judgements provides the best representation of uncertainty (i.e. validation),
32 whether that be via equally weighted or performance weighted aggregation.

33 Whether the time and effort entailed in performance weights can be justified is a matter for
34 decision-makers. Our case study outlines the rationale, challenges, and benefits of performance
35 weighted aggregations. It will help to inform decisions about the deployment of performance
36 weighting and avoid common pitfalls in its application.

37

38

39 **1 Introduction**

40 Over the past 15 years a considerable body of research has emerged in the ecological literature
41 emphasising the need for more rigorous and structured methods for collecting quantitative
42 expert judgements. The literature has summarised existing structured elicitation protocols and
43 key steps which could be adapted and applied to better suit the practical (e.g. geographically
44 dispersed experts) and financial (lack of funding) constraints of most ecological contexts
45 (Burgman 2004, Low Choy et al. 2009, Kuhnert et al. 2010, Burgman et al. 2011a, Martin et
46 al. 2012, McBride et al. 2012a, McBride et al. 2012b, Drescher et al. 2013).

47 A common approach that has been advocated is to recruit a diverse group of individuals and
48 take an equally weighted aggregation of their independent judgements (Burgman et al. 2011b,
49 Hemming et al. 2018b). This will often produce point estimates which are at least as accurate
50 (i.e. closer to the truth) and interval judgements which are better calibrated than the median-
51 ranked individual for these scores (Burgman et al. 2011b, Budescu and Chen 2014, Hemming
52 et al. 2018b). While one person can sometimes outperform the group aggregate, rarely can that
53 person be predicted by credentials conventionally associated with expertise such as age,
54 experience, or peer-identification (Aspinall and Cooke 2013, Burgman 2015, Mellers et al.
55 2015).

56 The performance of the equal weighted aggregation is largely explained as a statistical
57 phenomenon (Lorenz et al. 2011) in which the judgments of individuals represent random
58 independent samples. If those samples are diverse then not only should the information pool
59 related to the questions increase (Clemen and Winkler 1999), but the errors made by
60 individuals are more likely to cancel (Larrick and Soll 2006, Budescu and Chen 2014). This
61 phenomenon is often termed the ‘wisdom of the crowd’ (Surowiecki 2004), or the ‘staticised
62 group’ (Einhorn et al. 1977, Hogarth 1978). Interestingly, participants need not be experts and

63 can be biased, as long as they have some information related to the questions that can be
64 combined for prediction (Budescu and Chen 2014).

65 Equal weighting is advantageous as it's relatively simple to apply (Hogarth 1978, Hora 2004,
66 Hemming et al. 2018b). Typically, group sizes of 5-12 participants derive improved
67 judgements, with diminishing returns thereafter (Hogarth 1978, Hora 2004, Hemming et al.
68 2018b). It requires no additional work to develop questions or performance measures to score
69 and aggregate experts. It can be applied to any type of prediction including point estimates,
70 distributions and probabilities. The simplicity of equal weighting, and its ability to improve a
71 wide range of estimates make it suitable for aggregating judgements under the practical and
72 financial constraints typical of many ecological decisions.

73 However, despite substantial testing and real-world applications, many people find equal
74 weighting difficult to trust (Weiss and Shanteau 2004). This is partly because the method relies
75 on the recruitment of a diversity of individuals, often including individuals who may normally
76 be excluded from such elicitations because of their perceived limited knowledge (Shanteau et
77 al. 2002, Weiss and Shanteau 2004, Burgman et al. 2011a).

78 When uncertainty is elicited, the diversity of the group can also increase the uncertainty
79 associated with group judgements, sometimes leading to uninformative judgements
80 (MacDonald et al. 2008, Barons et al. 2018). Occasionally individuals will outperform the
81 group aggregation, and ideally decision-makers would like to restrict elicitation to these better
82 performing individuals, or at least have the judgments of those individuals weigh more than
83 those of lesser performers. Finally, there is no single method for generating an equally weighted
84 aggregation. For example, for point estimates, the arithmetic mean is commonly applied, but
85 one could also use the median, geometric mean or harmonic mean (Armstrong 2001, Colson
86 and Cooke 2017). Rarely is there any validation to support such choices made by the analyst,

87 which can lead to questions about the validity of the specific method chosen, and the influence
88 of analyst's subjective bias. The problems associated with equal weights can serve to
89 undermine the credibility of the final judgements derived and the subsequent decisions and
90 assessment based on such judgements.

91 Performance weighted aggregation is often suggested as a way to address these challenges and
92 perceived deficiencies (Cooke 1991, Budescu and Chen 2014, Mellers et al. 2015). It involves
93 developing sets of questions related to the main elicitation questions for which the answers can
94 be obtained but are not widely known to experts (Cooke 1991, Goossens et al. 2008, Tetlock
95 and Gardner 2015). These are referred to as test, seed or calibration questions (we use the term
96 calibration questions from hereon). Those who perform better on these questions are afforded
97 more weight in the final aggregation of the main elicitation questions. The method is
98 differentiated from other forms of weighted aggregation such as those based on self-rating,
99 peer-rating, trimming, or representativeness in that weights are obtained via validation of
100 judgements against an external truth (Armstrong 2001, Aspinall and Cooke 2013).

101 The main reason decision-makers seek to apply performance weights is to create aggregated
102 judgements which are more accurate (for point estimates), or well-calibrated and informative
103 (for interval judgements, probabilities and probability distributions) (Budescu and Chen 2014,
104 Mellers et al. 2015, Colson and Cooke 2017). However, the inclusion of calibration questions
105 is also seen to create a sense of legitimacy. It provides evidence that those who have been
106 included in the final aggregation have some knowledge in the relevant domain, and that they
107 can communicate their knowledge together with their uncertainty in the format required by the
108 analyst (Barons et al. 2018, Quigley et al. 2018). It can also be used to validate assumptions
109 made by the analyst in combining expert judgements.

110 Despite advocacy, there has been little progress in ecology towards understanding or applying
111 performance weighted aggregation, outside of a few applications (Metcalf and Wallace 2013,
112 Wittmann et al. 2015, Barons et al. 2018). We contend this has led to an under-appreciation of
113 the fundamental requirements of the method in ecology, of how the method can be practically
114 applied more widely in ecology, and the extent to which implementation may improve
115 outcomes.

116 In this paper we 1) illustrate how the IDEA protocol with four-step question format can be
117 extended to include performance weighted aggregation from the Classical Model, and 2)
118 explore the extent to which this extension improves pooled judgements for a range of
119 performance measures.

120 We choose the IDEA protocol (“Investigate”, “Discuss”, “Estimate”, and “Aggregate”) as it is
121 a structured elicitation protocol that has been tested and refined in the ecological literature and
122 (Hanea et al. 2016, Hemming et al. 2018a). The method involves first recruiting a diverse group
123 of individuals, and allowing each individual to “Investigate” the problem before making a
124 private individual estimate (often termed “Round 1”), following which experts see the
125 judgements of others and then enter into a “Discussion” phase. Experts then provide a final
126 private “Estimate” (“Round 2”). The judgements are “Aggregated”, typically using equal
127 weights (Figure 1).

128 Elicitation in the IDEA protocol can be undertaken remotely (i.e. via email), in a face-to-face
129 workshop, or by combining the two formats. This flexibility provides a practical advantage for
130 ecologists who usually have limited resources to convene experts face-to-face.

131 Most applications of the IDEA protocol in ecology aim to obtain quantitative judgements
132 together with uncertainty. When doing so, the four-step question format is often deployed

133 (Speirs-Bridge et al. 2010) (Figure 1). This method derives a credible interval with a ‘best’
134 point estimate based on the following questions:

- 135 1. Realistically what is the lowest plausible value for x?
- 136 2. Realistically what is the highest plausible value for x?
- 137 3. Realistically what is your best estimate for x?
- 138 4. Looking at your interval from lowest to highest, how confident are you that your
139 interval will capture the realised truth.

140 The four-step question format has been demonstrated to reduce overconfidence in interval
141 judgements relative to eliciting fixed quantiles (Speirs-Bridge et al. 2010). It has also helped
142 in obtaining quantitative judgements (with uncertainty) from experts who may eschew
143 quantification. Its development and application has improved the quality of information
144 derived from expert elicitation in ecology beyond that of categorical variables and point
145 estimates, which can be imbued with considerable ambiguity or fail to provide crucial
146 information about uncertainty (Wallsten et al. 1986, Gregory and Keeney 2017).

147 The practical advantages of the IDEA protocol with the four-step question format has seen the
148 adoption of the combined method spread rapidly in ecology (Adams-Hosking et al. 2016,
149 Hudson et al. 2017, Barons et al. 2018, Carwardine et al. 2019, Estévez et al. 2019). However,
150 it has been suggested that the aggregations derived could be further improved by incorporating
151 the performance weighted aggregation (Metcalf and Wallace 2013, Hemming et al. 2018a,
152 Hemming et al. 2018b).

153 The Classical Model (Cooke 1991) is a method for performance weighted aggregation often
154 cited in the ecological literature as a means to improve uncertain quantitative ecological
155 judgements (Burgman et al. 2011a, Martin et al. 2012, Drescher et al. 2013, Metcalf and
156 Wallace 2013, Hemming et al. 2018a). While it has been applied to a large number of

157 engineering case studies (Cooke and Goossens 2008, Colson and Cooke 2017) we are aware
158 of only two ecological examples, both in the Laurentian Great Lakes (Rothlisberger et al. 2009,
159 Wittmann et al. 2015).

160 In this this study we apply the Classical Model to a case study in which judgements were
161 elicited using the IDEA protocol (Hemming et al. 2018a) and four-step question format (Speirs-
162 Bridge et al. 2010). In doing so, we address the key aims of this study (outlined above), while
163 providing an insight into key considerations required for the deployment of performance
164 weighted aggregation more broadly.

165 **2 Methods**

166 **2.1 Fundamentals of performance weighting**

167 There is a considerable body of literature describing the application of performance weighting
168 with calibration questions, however, it is spread across a broad range of domains which can be
169 difficult to access and synthesise. We summarise key points to be considered prior to
170 application.

171 Generating performance weights with calibration questions entails (a) the development of
172 questions for which there are answers unknown to the participants, and (b) the selection of an
173 appropriate scoring rule to measure the performance of expert estimates.

174 There is little prescriptive guidance as to what makes a good calibration question, although
175 some features are self-evident (Cooke and Goossens 2000, Aspinall and Cooke 2013, Tetlock
176 and Gardner 2015, Quigley et al. 2018). They should relate to the knowledge needed to answer
177 the main elicitation questions (i.e. domain knowledge). They should ask questions about
178 uncertainty to capture an expert's ability to adapt and communicate their knowledge. They
179 should be in a similar format to the main elicitation questions. They should not be questions

180 which can be easily guessed, and not so hard that an expert could not reasonably form a
181 judgement. A substantial number of calibration questions may be required to differentiate luck
182 from good judgement, depending on the scoring rules. Ideally, questions should relate to
183 predictions of events or quantities rather than estimating the outcomes of past events
184 (retrodictions), although this is not always possible. The questions should be reviewed by at
185 least two people with domain knowledge to ensure they provide fair and reasonable
186 assessments of an expert's ability to make good judgements related to the main elicitation
187 questions.

188 One of the most important aspects of scoring rules is that they should not influence experts in
189 an undesirable way - termed proper scoring rules (Brier 1950). Strictly proper scoring rules are
190 those for which an expert maximises the expected score, if and only if they state their true
191 beliefs (Gneiting and Raftery 2007). There are many methods for scoring and assessing expert
192 judgements (Brier 1950, Cooke 1991, Flandoli et al. 2011, Budescu and Chen 2014, Satopää
193 et al. 2014, Hemming et al. 2018b), which vary depending on the types of judgements elicited
194 (probabilities, intervals, distributions etc). Not all scoring rules are proper scoring rules, and
195 few have been substantially tested and applied in real applications. The Brier Score is an
196 exception and has been used to assess performance of individuals and groups on single event
197 probabilities such as weather forecasts and geopolitical events, but has not been developed into
198 a formal weighting scheme (Brier 1950, Tetlock and Gardner 2015, Barons et al. 2018). The
199 other exception is the scoring rule of the Classical Model (discussed below) (Cooke 1991).

200 Scoring rules aim to optimise judgements and the way in which they do this depends on their
201 reward structure (Winkler and Murphy 1968, Tetlock 2005). For example, scoring rules for
202 interval judgements may penalise overconfidence (e.g. intervals that are too narrow, which
203 include the truth less often than the purported level of confidence provided by the expert) more
204 than under-confidence (e.g. intervals that are too broad and capture more realisations than the

205 purported level of confidence of the expert). It's therefore important to understand how such
206 transgressions of judgement are handled by a proposed scoring rule, to ensure that the reward
207 structure matches the preferences and needs of the decision-maker and the problem at hand.
208 This of course requires an awareness among decision-makers about what aspects of judgement
209 are most important to them.

210 Obtaining an understanding of the reward structure can be challenging as research papers
211 outlining the application of scoring rules rarely provide clear examples of how judgements are
212 incorporated and combined. Few adequately discuss their embedded reward structure. A
213 further complication arises in understanding scoring rules because the terms used to describe
214 judgement, such as 'calibration', 'accuracy' and 'overconfidence', are used interchangeably
215 and may refer to different concepts (Lichtenstein and Fischhoff 1977, Lichtendahl Jr et al.
216 2013, Cooke 2018b, Hemming et al. 2018b).

217 **2.2 The Classical Model**

218 In this paper we choose to investigate the application of the Classical Model (Cooke 1991).
219 The method was developed as a means for reaching rational consensus, which is defined by
220 Cooke and Goossens (2008) as an agreement as to how to derive a consensus distribution from
221 multiple, elicited distributions. Ultimately, it treats expert judgement as a form of empirical
222 data and promotes adherence to four critical elements of scientific inquiry: accountability,
223 empirical control, neutrality, and fairness (Cooke and Goossens 2008).

224 In elicitations employing the Classical Model, experts are asked a set of calibration questions
225 (usually 10-15), for which the answers can be obtained. As noted above, these questions should
226 relate to the main questions of the elicitation (termed target variables or questions of interest).
227 Unlike the four-step question format commonly used with the IDEA protocol, experts are asked
228 to specify their judgements as quantiles of a continuous non-parametric probability distribution

229 (usually 5th, 50th, and 95th) for both calibration questions and questions of interest. The
230 individual judgements of experts are typically elicited in a face-to-face elicitation with one or
231 more facilitators present (Wittmann et al. 2015). Experts are scored on their performance using
232 two performance measures (see section 2.4 for details): “statistical accuracy” (often termed
233 “calibration”), and “information” (sometimes termed “informativeness”, or “relative
234 information”). These are subsequently multiplied to provide an asymptotically proper scoring
235 rule (the CM Score) (refer to Appendix S1: Section 4.2.3), and to derive differential weights.

236 Experts who perform well on the calibration questions are afforded more weight in the final
237 aggregations for the questions of interest. Both equally weighted and performance weighted
238 linear pooled aggregations of distributions are then created and subsequently scored on their
239 performance on the calibration questions (i.e. via in-sample validation, where the same
240 questions used to develop the performance weighted aggregations are used to score the
241 aggregations). To achieve rational consensus, experts or decision makers usually agree prior to
242 the elicitation that the aggregation which achieves the highest combined score on the
243 calibration questions will be used to weight expert judgements of the target questions.

244 The primary purpose of performance weighting and calibration questions in the Classical
245 Model is to come to an unbiased and empirically validated decision on how to combine the
246 expert judgements. This step can help to overcome pre-judgements and exclusion of potentially
247 knowledgeable individuals, as well arbitrary choices by analysts and decision makers about
248 how to weight and aggregate experts. In analyses of 78 case studies using the Classical Model,
249 performance weighted aggregations have outperformed equal weights in 76 studies (in-sample
250 validation), suggesting the method can also be used to optimise aggregated judgements (Cooke
251 and Goossens 2008, Colson and Cooke 2017).

252 **2.3 Case study**

253 To demonstrate how the Classical Model could be applied in ecology, and to investigate
254 potential improvements from its application, we use estimates for ecological questions from a
255 previous case study by Hemming et al. (2018b). In brief, the case study used the IDEA protocol
256 with the four-step question format to elicit judgements for thirteen questions relating to future
257 abiotic and biotic events on the Great Barrier Reef. The elicitation was undertaken via email
258 and the experts volunteered their time. The questions related to the types of events experts may
259 be asked in assessing risk to the Great Barrier Reef (Ward 2014), for example, the percentage
260 cover of coral bleaching that may be detected in the next monitoring event at a specified reef
261 (see Appendix S1: Section 1). The questions related to future monitoring events, so that
262 judgements could be scored against outcomes once monitoring data were collected.

263 In total, 58 experts completed Round 2 of the elicitation exercise. These 58 individuals had
264 been randomly assigned to one of eight groups within which judgments were aggregated. In
265 Hemming et al. (2018b) the judgements were standardised to 80% credible intervals using
266 linear extrapolation (outlined in Appendix S1: Section 1) and subsequently aggregated using
267 an equal weighted quantile aggregation (taking the arithmetic mean) (refer to Appendix S1:
268 Section 5). The judgements were then scored using performance measures of the IDEA
269 protocol. The study found that 1) the equally weighted aggregate judgements were often more
270 accurate and better calibrated than the median individual, 2) individuals could outperform the
271 aggregation, however, they could not have been selected based on their credentials or
272 demographic data, and 3) discussion and feedback led to improved final judgements (Appendix
273 S1: Section 1). However, it was suggested further improvements may be made via performance
274 weighted aggregation.

275 **2.3.1 Four-step to quantiles**

276 To make responses of the four-step question format compatible with requirements of the
277 Classical Model (quantiles of a continuous non-parametric distribution), individual judgements
278 need to be standardised to 90% credible intervals. We then assume (a) that the best estimate is
279 the 50th percentile (i.e. a median), and, (b) upper and lower estimates represent a *central*
280 credible interval (i.e. whereby the probability mass beyond a judgment's interval is apportioned
281 equally above and below the upper and lower bounds, respectively). We interpret lower bounds
282 as 5th quantiles and upper bounds as 95th quantiles.

283 In zero-inflated settings it is possible for respondents to provide a judgment of zero for both
284 their 5th and 50th quantile (which occurred in our case study but is not consistent with a
285 continuous distribution - refer to Appendix S1: Section 2). In such cases, a small number may
286 be added or deducted to separate the quantiles. For example, zeros may be replaced by the
287 following numbers depending on where in the estimate the zeros occur (Cooke 2018a):

- 288 • Lower / 5th : 0.00001
- 289 • Best / 50th : 0.0001
- 290 • Upper / 95th: 0.001

291 In our case study, we also encountered circumstances where the lower estimate, or best estimate
292 reasonably coincided with the upper bounds which led to similar adjustments (see Appendix
293 S1: Sections 2-3).

294 **2.4 Scoring Judgements**

295 Assuming the judgements approximate quantiles of a continuous probability distribution, the
296 judgements can then be scored using the Classical Model's performance measures. There is
297 substantial ambiguity and confusion in the ecological literature as to what the performance
298 measures of the Classical Model actually reward. They have been cited as rewarding 'accuracy'

299 (Rothlisberger et al. 2009, Burgman et al. 2011a, Martin et al. 2012), which may give the
300 impression they reward the accuracy of point estimates. They have also been noted to score
301 ‘calibration’ and ‘precision’ (width) which may give the impression they are designed to assess
302 interval judgements according to definitions that arise in the psychological literature
303 (Lichtenstein and Fischhoff 1977, Yaniv and Foster 1997, Burgman et al. 2011a, Wittmann et
304 al. 2015).

305 Verbal clarifications contained within the Classical Model literature often fail to clarify the
306 reward structure, which may perpetuate misinterpretations. For example, statistical accuracy
307 has been described as a measure of the likelihood that “*at least 7 out of 10 realisations should*
308 *fall outside an expert's 90% confidence bands, if each value really had an independent 90%*
309 *chance of falling inside the bands?*” (Rothlisberger et al. 2009, Colson and Cooke 2017). This
310 may give the impression that it is designed primarily to score the calibration of the 90% credible
311 interval judgements, rather than the calibration of the expert’s interquantile ranges.

312 To better understand the reward structure of the Classical Model so that they are not misapplied
313 we will contrast the performance measures for the Classical Model with those commonly used
314 in the IDEA protocol (Hemming et al. 2018b). We outline these performance measures below.
315 Equations and a worked example are provided in Appendix S1: Section 4.

316 ***IDEA performance measures***

317 With the four-step question format in the IDEA protocol, individuals are scored by
318 performance measures of *accuracy*, *calibration* and *informativeness* (Hemming et al. 2018b).

319 *Accuracy* is designed to assesses the accuracy of point estimates. It is the difference between
320 *b*, the expert’s best estimate, and the observed value, *x*. It is measured using the average log
321 ratio error (ALRE) of expert responses. The measure is a relative measure, scale invariant, and
322 emphasizes order of magnitude errors rather than linear errors. Smaller ALRE scores indicate

323 more accurate responses. For any given question the log ratio score has a maximum possible
324 range of 0.31 ($=\log_{10}(2)$), which occurs when the true answer coincides with either the group
325 minimum or group maximum (Burgman et al. 2011b)

326 Calibration is the proportion of intervals provided by the experts containing the realised truth
327 relative to their assigned confidence (Lichtenstein and Fischhoff 1977, Lin and Bier 2008). For
328 example, if the expert's intervals are standardised to 90% credible intervals then we expect for
329 a well calibrated expert and 100 questions, that 90 of the realisations will fall between their 5th
330 and 95th quantiles. If they capture fewer realisations, they may be considered overconfident,
331 and if they capture more realisations they may be considered underconfident. The measure is
332 an absolute measure and is scale invariant. If the realisations are equal to the expert's 5th or 95th
333 quantiles, then they are usually assessed as being included within the expert's credible
334 intervals.

335 Informativeness is used to denote a measure of the width (or precision) of the intervals provided
336 by experts (Yaniv and Foster 1997). It is a relative measure and scale invariant. For each
337 question, the expert's intervals are divided by a background range for the question, where the
338 range is based on all estimates provided by the pool of experts for that question. Answers close
339 to 0 indicate that an expert was highly informative, while a 1 would indicate the expert's
340 uncertainty spanned the entire range of responses for that question. The final score for
341 informativeness for an expert is their average across all questions.

342 ***Performance measures of the Classical Model***

343 The Classical Model has two main performance measures that assess the ability of an expert to
344 provide useful probability distributions, statistical accuracy and information.

345 Statistical accuracy (often referred to as calibration and often denoted by 'C') assesses the
346 ability of experts to answer according to a theoretically optimal multinomial distribution. It

347 assesses the interquantile calibration of experts. For example, over a set of questions for which
348 realisations could be obtained, we would expect for any high performing expert that:

- 349 • For 5% of their judgements, the realisations would fall below their 5th quantile. We express
350 the observed proportion as Q_1 .
- 351 • For 45% of their judgements, the realisations would fall between their 5th and their 50th
352 quantile. We express the observed proportion as Q_2 .
- 353 • For 45% of their judgements, the realisations would fall between their 50th and their 95th
354 quantile. We express the observed proportion as Q_3 .
- 355 • For 5% of their judgements, the realisations would fall above their 95th quantile. We express
356 the observed proportion as Q_4 .

357 The expectation of where the realisations fall in relation to an expert's interquantile ranges can
358 be expressed as a theoretical multinomial distribution $p=(0.05, 0.45, 0.45, 0.05)$ (Bedford and
359 Cooke 2001). Under the Classical Model, the actual proportion of realisations within each
360 inter-quantile range for each expert (or aggregation) e , is tallied to create a multinomial
361 distribution for each expert: $s(e) = (Q_1, Q_2, Q_3, Q_4)$.

362 The realised distribution is then compared to the theoretical distribution using the Kullback-
363 Leibler (KL) divergence measure and a chi-square test with three degrees of freedom.
364 Statistical Accuracy is the p -value of this test. Higher values indicate an expert's distribution
365 more closely matches the theoretical distribution. A statistical accuracy below 0.05 is often
366 used as a cut-off point at which an expert is considered statistically inaccurate (i.e. Bamber et
367 al. (2016), Colson and Cooke (2017)). The 0.05 level is often used in meta-analyses comparing
368 the weighting and aggregation schemes in the Classical Model literature, but can also be used
369 by the analyst as a cut-off point at which zero weight may be assigned to the expert's
370 judgement.

371 In scoring expert judgements, if the realisations are equal to the values provided by the experts
372 for the 5th, 50th, and 95th quantiles, then the following rules are used to decide which probability
373 bin the realisation should be placed into:

- 374 • If the realisation equals the 5th quantile, it is placed in the first probability bin Q₁.
- 375 • If the realisation equals the 50th quantile, it is placed in the second probability bin Q₂.
- 376 • If the realisation equals the 95th, it is placed in the third probability bin Q₃.

377 We highlight this assumption as (on rare occasions) it can affect the score participants receive.
378 For example, in the unlikely case that a participant was to estimate the median perfectly for 9
379 of 10 questions, they could obtain a multinomial distribution of $S(e) = (1, 9, 0, 0)$, which when
380 compared to the theoretically optimal multinomial distribution means they would be
381 considered statistically inaccurate at the 0.05 level, despite having perfect calibration and
382 exceptional accuracy under the IDEA protocol scoring rules.

383 Information (often referred to as relative information, or informativeness) under the Classical
384 Model measures the degree to which the expert's distribution is concentrated and to which it
385 differs from a uniform or log-uniform distribution (which are considered the least informative
386 distributions). It uses the KL divergence measure, which is scale invariant (Quigley et al. 2018).
387 Information is calculated per question and does not depend on the realisation. The final
388 information score of an expert is an average taken across all calibration questions. Larger
389 numbers indicate better performance because they represent distributions which show greater
390 departure from a uniform or log-uniform distribution.

391 A simple example contrasting the performance measures is provided in Box 1, and Figure 2.
392 In the results section, we plot outcomes for these measures against each other to gain a better
393 understanding of the underlying reward structures.

394 **2.5 Weighting and aggregating**

395 There are notable trade-offs between statistical accuracy and information in the Classical
396 Model. By providing very wide intervals, an expert may achieve near perfect statistical
397 accuracy, but will have low information (Quigley et al. 2018). Likewise, by providing very
398 narrow intervals, they will have a high level of information, but usually at the cost of poor
399 statistical accuracy. Ideally an expert should have both high statistical accuracy and
400 information (Quigley et al. 2018). Therefore, the performance measures of the Classical Model
401 are only proper if they are combined.

402 Under the Classical Model, the scores for statistical accuracy and information are combined to
403 provide weights for each expert. There are five basic ways in which experts may be weighted
404 and combined (equations provided in the Appendix S1):

405 Equal Weights (EW): is a linear pool of all expert distributions using the arithmetic mean of
406 their distributions. It affords all experts the same weight regardless of how well they performed
407 on calibration questions. It can be calculated without calibration questions.

408 Global Weights (GW): is calculated based on the combined statistical accuracy and information
409 scores (CM Score) averaged across all calibration questions. Experts who performed better on
410 the calibration questions are afforded more weight than those who performed poorly.

411 Itemised Weights (IW): uses the same statistical accuracy scores as Global Weights, however,
412 the weight each expert is awarded will change per question because it considers the information
413 of the expert for each question of interest rather than the average calculated based on all of the
414 calibration questions. This often leads to aggregations with higher information (and
415 informativeness) on average than Global Weights.

416 Global Weights Optimised (GWO) and Itemised Weights Optimised (IWO): are similar to their
417 un-optimised variants described above (i.e. Global Weights (GW) and Itemised Weights (IW)).
418 However, they optimise the statistical accuracy score by successively raising the level at which
419 an expert is considered statistically inaccurate from an alpha level equal to the lowest
420 calibration score. The weights are calculated and used to generate weighted aggregations that
421 are scored on the calibration questions. The weighted aggregation with the highest performance
422 on the calibration questions is chosen (Quigley et al. 2018). In decisions with one or two well
423 calibrated experts, most or all of the weight may be assigned to those experts with no weight
424 given to the other experts.

425 For a set of calibration questions, an analyst may create a set of pooled judgements for each
426 question under each weighting scheme. These pooled judgements can then be scored for their
427 statistical accuracy and information (i.e. in-sample validation). These scores are then multiplied
428 to create an overall score, which we term the Classical Model (CM) Score. The aggregation
429 method which produces the highest CM Score on the calibration questions is usually taken as
430 the preferred weighting scheme when combining expert judgements on the questions of interest
431 (for which answers are not known). If two aggregations result in the same statistical accuracy,
432 that with a higher information score is preferred (Bedford and Cooke 2001).

433 **2.5.1 Linear pooling versus quantile aggregation**

434 The Classical Model uses linear pooling of distributions for both equal weighted and
435 performance weighted aggregations, which differs from quantile aggregation commonly used
436 by the IDEA protocol when the four-step question format is used (Hemming et al. 2018a) (refer
437 to Appendix S1: Section 5 for discussion and a worked example).

438 Quantile aggregation is simple to apply, and entails no additional assumptions about what the
439 estimates represent beyond a best estimate with a credible interval. In general, it provides more
440 accurate and better calibrated judgements compared to the best-regarded experts (Burgman et
441 al. 2011b, Hemming et al. 2018b). However, Bamber et al. (2016) and Colson and Cooke
442 (2017) found that quantile aggregation is much more overconfident than linear pooling (when
443 assessed using the Classical Model's Statistical Accuracy measure). To investigate these
444 findings, we extend our analysis to compare how the two methods of equally weighted
445 aggregation can affect judgements. Henceforth we use the term 'equal weights' (abbreviated
446 to EW) to refer to linear pooling of distributions, and 'quantile aggregation' (abbreviated to
447 QuA) to refer to quantile aggregation.

448 **2.6 Analysis**

449 For the eight groups of experts in our case study, we assessed the six alternative approaches to
450 aggregation (two forms of equal weighted aggregations (EW (Classical Model), QuA (IDEA)),
451 and four forms of performance weighted aggregation from the Classical Model (IW, GW, IWO,
452 GWO) (described in Section 2.5). Individual and group performance was assessed using the
453 five performance measures (statistical accuracy, information, calibration, informativeness, and
454 accuracy) (described in Section 2.4), and the Classical Model scoring rule (CM
455 Score)(described in Section 2.5).

456 To obtain the performance measures and aggregations associated with the Classical Model, the
457 analyst must enter judgements in software called *Excalibur* (Lightwist 2013, Cooke

458 2018a)(Appendix S1: Section 3). For measures associated with the IDEA protocol we
459 developed *R*-code (available on the Open Science Framework (Hemming 2019)). More details
460 are available in Appendix S1: Section 3.

461 To contrast the differences of the aggregations, we use boxplots, constructed in *R* (version 3.4.1
462 (2017-06-30) -- "Single Candle"), using the *ggplot2* package. The boxes represent the 25th,
463 50th and 75th percentiles. The whiskers represent the spread of the data referenced on the inter-
464 quartile range, $(Q1-1.5*IQR, Q3+1.5*IQR)$. For normally distributed data this is
465 approximately 2.7 standard deviations, or 99.3% of the data (Krzywinski and Altman 2014).

466 **3 Results**

467 **3.1 Comparison of performance measures**

468 In Figure 3, we plot the two performance measures underpinning weights obtained under the
469 Classical Model for the 58 participants. When scored on statistical accuracy, less than half (23)
470 of participants were statistically accurate (obtaining scores higher than 0.05). For 13 questions
471 the highest possible statistically accuracy score would have been 0.93. No individuals achieved
472 this score (highest statistical accuracy score was 0.53).

473 High statistical accuracy usually came at the expense of lower information. Participants who
474 were statistically accurate were more likely to have lower information scores. Such
475 observations reflect the trade-offs between statistical accuracy and information discussed by
476 Quigley et al. (2018) and re-enforce the need to combine the two measures to derive a proper
477 scoring rule.

478 Figure 4 shows the scatter between statistical accuracy (the Classical Model) and calibration
479 (IDEA Protocol) for the 58 participants over 13 questions. While the difference in the highest
480 statistical accuracy score possible and that obtained by experts appears large (i.e. a change from

481 0.93 to 0.53 implies a 43% reduction in statistical accuracy), we can see that this change was
482 due to just one additional realisation falling outside of the experts' credible intervals. Thus,
483 statistical accuracy can be highly sensitive to seemingly small variations in performance.

484 Figure 4 also shows that while there is a positive correlation between the two measures
485 (Spearman rank correlation= 0.84, 95% CI: 0.74, 0.90) there are also some notable differences.
486 Importantly, an expert may have near perfect calibration under the scoring rules employed by
487 the IDEA protocol, but be statistically inaccurate at the 0.05 level according to the Classical
488 Model. These results further clarify that statistical accuracy does not reward calibration
489 primarily between the expert's 90% credible intervals, on which IDEA's calibration depends.

490 In Appendix S1: Section 7, we demonstrate that the differences occur because the Classical
491 Model's rules score a multinomial distribution with three degrees of freedom $p = (0.05, 0.45,$
492 $0.45, 0.05)$. As such, beyond very low levels of calibration (i.e. <50% calibration for 13
493 questions), the statistical accuracy measure cannot be used to assess the calibration of 90%
494 credible intervals (i.e. a multinomial distribution with one degree of freedom, or a binomial
495 distribution).

496 Figure 5 shows the correlation between information (Classical Model) and informativeness
497 (IDEA Protocol). The two scores are negatively correlated (Spearman rank correlation = -0.69,
498 95%CI: -0.80, -0.52), an artefact of the scoring rules, whereby under the IDEA protocol
499 participants who receive a low score are more informative (narrower intervals), whereas for
500 the Classical Model a higher score indicates that they provide more information relative to a
501 uniform or log-uniform distribution. Figure 5 demonstrates that information and
502 informativeness are slightly different measures of an expert's judgement. The Classical Model
503 does not only assess the width of intervals, it also accounts for their departure from a uniform
504 distribution. This can mean that higher information score may be obtained in some cases simply

505 by reducing the symmetry of the ranges between an expert's 2nd and 3rd quantiles (i.e. if the
506 median does not fall squarely in the centre of the range then information can be increased).

507 **3.2 Performance of aggregations**

508 Figure 6 shows the CM Score for each of the aggregations. In the Classical Model, this
509 combined score would be used to select the final aggregation for uncertainty by a decision-
510 maker. For this case study, if we were to use the median values of these scores, we would not
511 choose quantile aggregation (QuA) because it has a low CM Score (median value of 0.14).
512 Equally weighted linear pooling employed by the Classical Model does better (EW) (median
513 value of 0.41), and there is some indication that performance weighted aggregation by the
514 optimised variants (IWO, and GWO) may lead to further improvements (median values of 0.60
515 and 0.50 respectively).

516 Figures 7a and 7b decompose the CM Score provided in Figure 6, into statistical accuracy and
517 information scores of the Classical Model. While quantile aggregation (QuA) performs well
518 on information (median value of 1.51), it performs poorly in terms of statistical accuracy
519 (median value of 0.10) compared to equal weights (EW) (median value of 0.36) and
520 performance weighted aggregations (IW, IWO, GW, GWO) (all achieving a median value of
521 0.36, except for IW which achieves 0.27), with two groups considered statistically inaccurate
522 at the 0.05 level. This supports the finding by Bamber et al. (2016) and Colson and Cooke
523 (2017) that quantile aggregation used in the IDEA protocol with four-step question format can
524 be overconfident relative to linear-pooling of distributions when assessed by statistical
525 accuracy.

526 There is little or no difference in the median performance of equally weighted (EW) and the
527 performance weighted aggregations (IW, IWO, GW, GWO) in terms of statistical accuracy.
528 However, both the optimised aggregations (GWO, and IWO) and itemised weights (IW) have

529 higher information (1.56 and 1.68) than equal weights (EW, median of 1.14) or global weights
530 (GW, median of 1.18), and are equivalent to quantile aggregation (QuA, 1.51) suggesting
531 performance weighting improves estimates in this case study by being more informative than
532 equal weights (EW).

533 Figures 7c-e assess each of the aggregation methods according to measures commonly used in
534 the IDEA protocol. Even when scored according to calibration between the expert's 90%
535 credible intervals, the study finds that quantile aggregation (QuA) generates more
536 overconfident estimates (median calibration of 0.77), having a lower calibration than all other
537 aggregations (0.85, or on average by one question). It does, however, have a higher level of
538 informativeness (0.25) than all other aggregations (medians ranging between 0.33 and 0.42),
539 including optimised aggregations. The median accuracy of the best estimate is better for all
540 aggregations than the median ranked individual for this measure. However, the optimised
541 aggregations have some groups which perform worse than the median individual. This may not
542 be surprising because (as discussed) the Classical Model was not designed to optimise point
543 estimates.

544 Quantile aggregation (QuA) performed relatively poorly on statistical accuracy and calibration
545 (Figure 7). Recall that some questions related to count data, and the upper and lower bounds
546 were adjusted so that they did not contain zero. In our case study, the lowest estimate which
547 could be provided by an expert was 0.00001. This adjustment may have led to overconfident
548 judgements for two questions which contained zeros.

549 To check this, we replaced the answers for these two questions with 0.000011 and re-calculated
550 the calibration and statistical accuracy of judgements of each of the groups (Figure 8, see also
551 Appendix S1: Section 8). The adjustment improved the statistical accuracy of many groups
552 across all aggregations. All but one aggregation (quantile aggregation, QuA) had a median

553 statistical accuracy above 0.53 (Figure 8a). Only one group was considered statistically
554 inaccurate when their judgements were combined via quantile aggregation (QuA).

555 Quantile aggregation (QuA) was overconfident, even when assessed according to calibration
556 of interval judgements but many groups were less so than prior to accounting for the two
557 questions with zeros (Figure 8b). Group judgements for quantile aggregation achieved good
558 but not perfect median calibration of 0.76, although no group reached perfect calibration when
559 quantile aggregation (QuA) was used. In contrast, each of the linear pooled distributions except
560 for the itemised optimised weights achieved a median group calibration of 0.90 (i.e. perfect
561 calibration). The data adjustments improved calibration and statistical accuracy, but they did
562 not substantially alter the information or informativeness scores which meant that quantile
563 aggregation (QuA) was still substantially more informative than the equal weights (EW) (a
564 median informativeness score of 0.24 compared to 0.42).

565 **4 Discussion**

566 Performance weights have been proposed to improve expert judgements in ecology. However,
567 there have been few applications and little discussion of their strengths and weaknesses in the
568 ecological literature. Here, we outlined the key rationales and theories of performance weights,
569 then described one of the most well-known methods, the Classical Model (Cooke 1991), and
570 examined how it might be applied to improve judgements derived from the IDEA protocol with
571 four-step question format (Hemming et al. 2018a, Hemming et al. 2018b).

572 This study highlighted how the Classical Model and the IDEA protocol may be integrated, but
573 clarified important differences between them that should be considered before applying
574 performance weights.

575 The four-step question format needs to first be converted into quantiles of a continuous
576 probability distribution. It may be better to remove these assumptions by eliciting these

577 quantiles directly. However, the four-step question format is often used because it helps to
578 overcome overconfidence relative to eliciting fixed intervals (Speirs-Bridge et al. 2010), and
579 because experts who are unfamiliar with the language of statistical distributions are
580 comfortable in providing quantitative judgements of uncertainty (a problem not only
581 encountered in ecological domains (Walls and Quigley 2001, Hirsch et al. 2004)). These trade-
582 offs need to be considered when deciding how best to elicit estimates. If the four-step question
583 format is to be used with the Classical Model, then we suggest that the assumptions about how
584 the estimates will be interpreted are communicated to experts in introductory material and
585 through the feedback and discussion stages of the IDEA protocol.

586 Once judgements were converted into quantiles of a continuous probability distribution, we
587 described key steps required to incorporate the judgements into *Excalibur* and to generate
588 scores and aggregations for the Classical Model (outlined in more detail in the Appendix S1:
589 Sections 2-3). These steps have not been substantially documented in the literature, inhibiting
590 use of performance weights. The advice outlined here will make implementation of the method
591 more accessible to those unfamiliar with the Classical Model and improve efficiencies when
592 analysing data.

593 We then described the performance measures underpinning the Classical Model, noting that
594 there was considerable ambiguity in the literature as to how the Classical Model rewards
595 judgements, with terms such as “calibration”, “accuracy”, “information”, and “overconfidence”
596 being differently interpreted (Rothlisberger et al. 2009, Burgman et al. 2011b, Metcalf and
597 Wallace 2013, Wittmann et al. 2015, Colson and Cooke 2017).

598 Insights from our results emphasise that the Classical Model was designed to assess probability
599 distributions rather than point estimates or interval judgements (as some interpretations
600 suggest). Specifically, ‘statistical accuracy’ measures the degree to which an expert’s

601 multinomial distribution matches a theoretically optimal multinomial distribution, and
602 ‘information’ measures the departure from a uniform or log-uniform background measure. As
603 such the Classical Model is not focused primarily on avoiding surprises outside of the 90%
604 confidence intervals, or the precision of the intervals (as assessed in the IDEA protocol) and
605 may lead to counterintuitive outcomes in settings where this is a primary concern.

606 The question therefore arises as to when each performance measure may be more appropriate?
607 Calibration, informativeness and accuracy (as scored in the IDEA protocol) tend to be
608 important in the contexts of risk assessments and structured decision-making in which
609 decision-makers are deciding to take action, and are using the best estimate to understand the
610 most likely scenario, or the uncertainty bounds to investigate how sensitive their decisions are
611 to different risk attitudes (Gregory et al. 2012, Addison et al. 2015). In other words, the
612 measures normally associated with IDEA may be most useful when assessing the outputs of a
613 model or risk analysis (Morgan and Henrion (1990), page 78).

614 On the other hand, it may be more important to understand the calibration within the expert’s
615 interquartile ranges (i.e. the 2nd and 3rd quantiles) (as scored by the Classical Model) when they
616 estimate probability distributions as inputs to a model, for example sampling in Monte Carlo
617 simulations, especially where tail risks are a key concern (Morgan and Henrion (1990), page
618 78).

619 While calibration and informativeness of interval judgements may be of interest they have not
620 yet been combined into a proper scoring rule (although telling experts they will be scored on
621 both should minimise gaming behaviour). Our results demonstrate that the Classical Model
622 does not by itself provide this information, which may be disappointing to those who seek to
623 apply the Classical Model to optimise or assess such judgements. However, if this information
624 was of interest the performance measures of the IDEA protocol may be used to provide this

625 information. Agreement as to which performance measures will be used should be made prior
626 to application.

627 Equal weighted aggregations are often used in ecology when combining expert judgments.
628 However, there are numerous methods by which an equal weighted aggregation can be derived,
629 and not all will perform equally well or have been validated. We contrasted two forms of equal
630 weighted aggregation, quantile aggregation (QuA, used in Hemming et al. 2018), and equal
631 weighting via linear pooling of distributions (EW, used by the Classical Model). We found that
632 both forms of equal weighted aggregation were better than the median ranked individual for
633 each measure of statistical accuracy, calibration, and accuracy. Furthermore, as was
634 demonstrated in Hemming et al. (2018b), while some individuals could outperform the group
635 aggregation they could not be predicted by standard metrics of expertise (years of experience,
636 peer-recommendation, or self-rating). This suggests that taking the equal weighted aggregation
637 is a more robust method than trying to select a single expert with good judgement based on
638 their credentials and status.

639 Our results corroborate those of the Bamber et al. (2016) and Colson and Cooke (2017), that
640 while quantile aggregation is simpler to apply, and was more informative, it led to
641 overconfident estimates compared to linear pooling of equally weighted distributions, and
642 performance weighted distributions (Figure 7). This was true regardless of whether we assessed
643 the judgements based on calibration or statistical accuracy.

644 We found that the degree of overconfidence was reduced when we accounted for questions
645 with zeros, and the way in which the Classical Model accounts for realisations which equal a
646 participant's estimates (i.e. if the realisation coincides with the lower bound it will be
647 considered as falling outside of the expert's 90% credible intervals). As these adjustments are
648 not made when the four-step question format is used in the IDEA protocol, the degree of

649 overconfidence from quantile aggregation may not typically be as severe for many applications
650 of the IDEA protocol. Nonetheless, we would suggest these findings warrant further
651 investigation on more case studies with the four-step question-format.

652 We then examined how performance weighting could be used to improve aggregated
653 judgements. We found that there was little difference in the calibration or statistical accuracy
654 of performance weighting and equal weighted linear pooled distributions. However,
655 performance weighting produced more informative bounds than equal weighted linear pooling
656 (by 10% of the background range when measured according to informativeness). These results
657 suggest that if the aim is to reduce arbitrary uncertainty while achieving well-calibrated
658 intervals, then performance weights can better achieve this.

659 In our study, we demonstrated a modest improvement by performance weighted aggregation.
660 However, we note that there is no guarantee that performance weighted aggregation will lead
661 to improvements in all cases. However, a clear advantage of the Classical Model, and other
662 methods which utilise calibration questions is that they provide empirical evidence for the
663 legitimacy of final aggregations (often lacking in studies that use expert judgement). This is
664 especially important because decisions regarding who should be included in an elicitation and
665 how to aggregate these judgements may exclude potentially knowledgeable individuals, and
666 often lack validation.

667 Whether or not the decision context justifies (or can afford) the additional time and expense
668 ultimately depends on the context of the case study, the decision-maker and the value of
669 additional information. Wittmann et al. (2015) and Rothlisberger et al. (2009) justify their
670 application based on the immense value of fisheries to the Great Lakes and the possibility of
671 litigation following mismanagement. This suggests that there are contexts in ecology in which
672 this additional time and expense can be justified. If resources are not available to deploy

673 calibration questions and performance weighted aggregation, then our study shows that an
674 equal weighted aggregation (i.e. quantile aggregation or linear pooling of distributions)
675 provides an effective means to improve judgements relative to selecting a single seemingly
676 well-credentialed expert.

677 However, there are obstacles to wider uptake of performance weighting and lines for further
678 research. We found it difficult to develop questions about future events on the Great Barrier
679 Reef for which we could obtain data in a reasonable time (3-6 months). Despite the substantial
680 amount of monitoring which takes place there (GBRMPA 2014). Others have noted problems
681 in obtaining access to ecological datasets (Meek et al. 2015). It may be possible to use existing
682 datasets to generate calibration questions. However, especially with remote elicitation, there
683 will always be a risk that experts discover the sources of the data when forming their
684 judgements (as occurred in (Hemming et al. 2019a)).

685 We found that questions relating to count data (particularly where the realisations are often
686 zero inflated) should be avoided when using the Classical Model. In ecology, zero inflated
687 count data are common (Martin et al. 2005).

688 Calibration questions should be related to target variables, for which the answer is known or
689 will become known (Cooke and Goossens 2000). However, ascertaining whether or not a
690 question is *relevant* in many domains may be difficult because domains are often ill-defined,
691 making the selection of *relevant* questions a subjective decision (Colyvan and Ginzburg 2003).
692 If datasets are difficult to obtain, then the analyst may need to rely on past questions for which
693 the data are available, or questions which are less relevant to the questions of interest. It would
694 be useful to understand at what point calibration questions become so distantly related to target
695 questions that in-sample validation is not a good predictor of performance.

696 We used *Excalibur* to generate aggregations and score experts, however, the program was
697 challenging to use. The analysis was time consuming and it was difficult to provide a
698 reproducible workflow for our analysis. The methods of aggregation and the scoring rules
699 should be simple enough to re-code in *R* and other freely available software (we note that
700 recently they have been re-coded in *MATLAB* (Leontaris and Morales-Nápoles 2018)). A
701 revision of *Excalibur* could help to increase adoption of the method.

702 Our study explored the effect of performance weights using in-sample validation (i.e. on the
703 same questions used to score experts and generate aggregations) for one case study. However,
704 the ideal test is how well it performs out-of-sample (i.e. on questions not used in the training
705 set) (Clemen 2008). This has not been addressed by this study. When Colson and Cooke (2017)
706 addressed this question they found some differences in out-of-sample performance that were
707 not revealed by in-sample validation and suggested this would be the focus of further research.

708 The scoring rules and aggregation methods of the Classical Model may not always be well-
709 understood. To avoid confusion, we suggest that in future, statistical accuracy scores should be
710 accompanied by their corresponding multinomial distributions. We provide *R* and *MATLAB*
711 code for this (Hemming et al. 2019b). While, it's less easy to convey the reward structure of
712 the information score, we believe it would be useful to display the intervals of the aggregations
713 so that the relative improvements can be compared (this is already often presented in
714 applications of the Classical Model).

715 **5 Conclusions**

716 Performance weighted aggregations with calibration questions has been proposed as a means
717 to improve expert judgements in ecology, however, applications have been scarce. We
718 explored how the Classical Model could be applied to the IDEA protocol with four-step
719 question format.

720 Our study found that the Classical Model could be applied to the IDEA protocol with four-step
721 question format provided the values of the four-step elicitation can be assumed to represent
722 quantiles of a continuous distribution. A key finding of this paper is that the reward structures
723 embedded in the performance measures of the two approaches to elicitation are often confused
724 and differ in important ways. This should be understood prior to application to ensure that the
725 methods for optimisation match the decision-maker's preferences and problem setting.

726 We demonstrated that equal weighted aggregations can achieve relatively well-calibrated
727 aggregated judgements. However, linear pooling of distributions may produce better calibrated
728 but less informative distributions than quantile aggregation as found by Bamber et al. (2016)
729 and Colson and Cooke (2017). We found that performance weighted aggregations can
730 outperform equal weighted aggregations, in our case by providing more informative
731 judgments, however, we emphasise that there is no guarantee they will do so in every case. The
732 main reason that the candidate alternatives for aggregation should be explored is to ensure the
733 final representation of uncertainty is the best possible (whether that be via equal weights or
734 performance weights).

735 Whether the time and investment in applying performance weights is worth the benefits is
736 ultimately a matter of context. Our example illustrates that there are contexts in which this
737 additional time and effort may be justified.

738 Our paper will help ecologists to better understand the fundamental steps, challenges, and
739 advantages involved in deploying performance weighted aggregation, and to avoid common
740 pitfalls which may arise. We welcome more research to understand how these methods could
741 be adapted to better suit the practical and financial constraints of a wider range of ecological
742 applications and estimates (i.e. point estimates, interval judgements, and single event
743 probabilities).

744 **6 Acknowledgements**

745 The authors would like to thank the experts who volunteered their time for the case study
746 presented. VH received funding to draft this publication by the Australian Research Training
747 Program, and the David Hay Memorial Fund. VH and AH were funded by the Australian
748 Centre of Excellence for Biosecurity Risk Analysis, VH, AH, TW and MB were funded by the
749 School of BioSciences at the University of Melbourne, MB was also funded by Centre for
750 Environmental Policy, Imperial College London.

751

752 **7 Literature citations**

- 753 Adams-Hosking, C., M. F. McBride, G. Baxter, M. Burgman, D. de Villiers, R. Kavanagh,
754 I. Lawler, D. Lunney, A. Melzer, P. Menkhorst, R. Molsher, B. D. Moore, D.
755 Phalen, J. R. Rhodes, C. Todd, D. Whisson, and C. A. McAlpine. 2016. Use of
756 expert knowledge to elicit population trends for the koala (*Phascolarctos cinereus*).
757 *Diversity and Distributions* **22**:249-262.
- 758 Addison, P. F. E., K. de Bie, and L. Rumpff. 2015. Setting conservation management
759 thresholds using a novel participatory modeling approach. *Conservation Biology*
760 **29**:1411-1422.
- 761 Armstrong, J. S. 2001. Combining forecasts. Pages 417-439 *in* J. S. Armstrong, editor.
762 *Principles of forecasting: A handbook for researchers and practitioners*. Springer
763 US, Boston, United States of America.
- 764 Aspinall, W. P., and R. M. Cooke. 2013. Quantifying scientific uncertainty from expert
765 judgement elicitation. Pages 64-99 *in* J. Rougier, S. Sparks, and L. Hill, editors.
766 *Risk and Uncertainty Assessment for Natural Hazards*. Cambridge University
767 Press, Cambridge, United Kingdom.
- 768 Bamber, J., W. Aspinall, and R. Cooke. 2016. A commentary on “how to interpret expert
769 judgment assessments of twenty-first century sea-level rise” by Hylke de Vries and
770 Roderik SW van de Wal. *Climatic Change* **137**:321-328.
- 771 Barons, M. J., A. M. Hanea, S. K. Wright, K. C. Baldock, L. Wilfert, D. Chandler, S. Datta,
772 J. Fannon, C. Hartfield, and A. Lucas. 2018. Assessment of the response of

773 pollinator abundance to environmental pressures using structured expert elicitation.
774 *Journal of Apicultural Research*:1-12.

775 Bedford, T., and R. M. Cooke. 2001. *Mathematical tools for probabilistic risk analysis*.
776 Cambridge University Press, Cambridge, United Kingdom.

777 Brier, G. W. 1950. Verification of forecasts expressed in terms of probability. *Monthly*
778 *weather review* **78**:1-3.

779 Budescu, D. V., and E. Chen. 2014. Identifying expertise to extract the wisdom of crowds.
780 *Management Science* **61**:267-280.

781 Burgman, M., A. Carr, L. Godden, R. Gregory, M. McBride, L. Flander, and L. Maguire.
782 2011a. Redefining expertise and improving ecological judgment. *Conservation*
783 *Letters* **4**:81-87.

784 Burgman, M. A. 2004. Expert frailties in conservation risk assessment and listing
785 decisions. Pages 20-29 *in* P. Hutchings, D. Lunney, and C. Dickman, editors.
786 *Threatened species legislation: is it just an Act?* Royal Zoological Society,
787 Mosman, NSW, Australia.

788 Burgman, M. A. 2015. *Trusting Judgements: How to get the best out of experts*.
789 Cambridge University Press, Cambridge, United Kingdom.

790 Burgman, M. A., M. McBride, R. Ashton, A. Speirs-Bridge, L. Flander, B. Wintle, F.
791 Fidler, L. Rumpff, and C. Twardy. 2011b. Expert status and performance. *PLoS*
792 *One* **6**:1-7.

793 Carwardine, J., T. G. Martin, J. Firm, R. P. Reyes, S. Nicol, A. Reeson, H. S. Grantham,
794 D. Stratford, L. Kehoe, and I. Chadès. 2019. Priority Threat Management for
795 biodiversity conservation: A handbook. *Journal of Applied Ecology* **56**:481-490.

796 Clemen, R. T. 2008. Comment on Cooke's classical method. *Reliability Engineering &*
797 *System Safety* **93**:760-765.

798 Clemen, R. T., and R. L. Winkler. 1999. Combining probability distributions from experts
799 in risk analysis. *Risk Analysis* **19**:187-203.

800 Colson, A. R., and R. M. Cooke. 2017. Cross validation for the classical model of
801 structured expert judgment. *Reliability Engineering & System Safety* **163**:109-120.

802 Colyvan, M., and L. R. Ginzburg. 2003. Laws of nature and laws of ecology. *Oikos*
803 **101**:649-653.

804 Cooke, R. 2018a. Macro converting XL file to EXCALIBUR dtt file. *in* R. Cooke, editor.
805 <http://rogermcooke.net/>

806 Cooke, R., and L. Goossens. 2000. Procedures guide for structural expert judgement in
807 accident consequence modelling. *Radiation Protection Dosimetry* **90**:303-309.

808 Cooke, R. M. 1991. *Experts in uncertainty: Opinion and subjective probability in science.*
809 Oxford University Press, New York.

810 Cooke, R. M. 2018b. Validation in the Classical Model. Pages 37-59 *in* L. C. Dias, A.
811 Morton, and J. Quigley, editors. *Elicitation: The science and art of structuring*
812 *judgement.* Springer International Publishing, Cham, Switzerland.

813 Cooke, R. M., and L. L. Goossens. 2008. TU Delft expert judgment data base. *Reliability*
814 *Engineering & System Safety* **93**:657-674.

815 Drescher, M., A. Perera, C. Johnson, L. Buse, C. Drew, and M. Burgman. 2013. Toward
816 rigorous use of expert knowledge in ecological research. *Ecosphere* **4**:1-26.

817 Einhorn, H. J., R. M. Hogarth, and E. Klempler. 1977. Quality of group judgment.
818 *Psychological Bulletin* **84**:158.

819 Estévez, R. A., F. O. Mardones, F. Álamos, G. Arriagada, J. Carey, C. Correa, J. Escobar-
820 Doderó, A. Gaete, A. Gallardo, and R. Ibarra. 2019. Eliciting expert judgements to
821 estimate risk and protective factors for Piscirickettsiosis in Chilean salmon
822 farming. *Aquaculture*.

823 Flandoli, F., E. Giorgi, W. P. Aspinall, and A. Neri. 2011. Comparison of a new expert
824 elicitation model with the Classical Model, equal weights and single experts, using
825 a cross-validation technique. *Reliability Engineering & System Safety* **96**:1292-
826 1310.

827 GBRMPA. 2014. Great Barrier Reef Outlook Report. Townsville, Australia.

828 Gneiting, T., and A. E. Raftery. 2007. Strictly proper scoring rules, prediction, and
829 estimation. *Journal of the American Statistical Association* **102**:359-378.

830 Goossens, L. H. J., R. M. Cooke, A. R. Hale, and L. Rodić-Wiersma. 2008. Fifteen years
831 of expert judgement at TUDelft. *Safety Science* **46**:234-244.

832 Gregory, R., L. Failing, M. Harstone, G. Long, T. McDaniels, and D. Ohlson. 2012.
833 *Structured Decision Making: a practical guide to environmental management*
834 *choices*, Chichester, West Sussex.

835 Gregory, R., and R. L. Keeney. 2017. A Practical Approach to Address Uncertainty in
836 Stakeholder Deliberations. *Risk Analysis* **37**:487-501.

837 Hanea, A., M. McBride, M. Burgman, B. Wintle, F. Fidler, L. Flander, B. Manning, and
838 S. Mascaro 2016. InvestigateDiscussEstimateAggregate for structured expert judgement.
839 International journal of forecasting **33**:267-269.

840 Hemming, V. 2019. Code: Weighting and Aggregating Expert Ecological Judgements.
841 The Open Science Framework. DOI 10.17605/OSF.IO/FXQVK.
842 <http://osf.io/fxqvk>

843 Hemming, V., N. Armstrong, M. A. Burgman, and A. M. Hanea. 2019a. Improving expert
844 forecasts in reliability: Application and evidence for structured elicitation
845 protocols. Quality and Reliability Engineering International **n/a**.

846 Hemming, V., M. A. Burgman, A. M. Hanea, M. F. McBride, and B. C. Wintle. 2018a. A
847 practical guide to structured expert elicitation using the IDEA protocol. Methods
848 in Ecology and Evolution **9**:169-181.

849 Hemming, V., S. Lane, and A. Hanea. 2019b. Classical Model Calculator. The Open
850 Science Framework. DOI 10.17605/OSF.IO/BGYZU. <http://osf.io/bgyzu>

851 Hemming, V., T. V. Walshe, A. M. Hanea, F. Fidler, and M. A. Burgman. 2018b. Eliciting
852 improved quantitative judgements using the IDEA protocol: A case study in natural
853 resource management. PLoS One **13**:e0198468.

854 Hirsch, K. G., J. J. Podur, R. F. Janser, R. S. McAlpine, and D. L. Martell. 2004.
855 Productivity of Ontario initial-attack fire crews: results of an expert-judgement
856 elicitation study. Canadian Journal of Forest Research **34**:705-715.

857 Hogarth, R. M. 1978. A note on aggregating opinions. Organizational Behavior and
858 Human Performance **21**:40-46.

859 Hora, S. C. 2004. Probability judgments for continuous quantities: Linear combinations
860 and calibration. *Management Science* **50**:597-604.

861 Hudson, E. G., V. J. Brookes, and M. P. Ward. 2017. Assessing the risk of a canine rabies
862 incursion in Northern Australia. *Frontiers in Veterinary Science* **4**:141.

863 Krzywinski, M., and N. Altman. 2014. Points of Significance: Visualizing samples with
864 box plots. *Nat Meth* **11**:119-120.

865 Kuhnert, P. M., T. G. Martin, and S. P. Griffiths. 2010. A guide to eliciting and using
866 expert knowledge in Bayesian ecological models. *Ecology Letters* **13**:900-914.

867 Larrick, R. P., and J. B. Soll. 2006. Intuitions about combining opinions: Misappreciation
868 of the averaging principle. *Management Science* **52**:111-127.

869 Leontaris, G., and O. Morales-Nápoles. 2018. ANDURIL — A MATLAB toolbox for
870 ANalysis and Decisions with UnceRtaInty: Learning from expert judgments.
871 *SoftwareX* **7**:313-317.

872 Lichtendahl Jr, K. C., Y. Grushka-Cockayne, and R. L. Winkler. 2013. Is it better to
873 average probabilities or quantiles? *Management Science* **59**:1594-1611.

874 Lichtenstein, S., and B. Fischhoff. 1977. Do those who know more also know more about
875 how much they know? *Organizational Behavior and Human Performance* **20**:159-
876 183.

877 Lightwist. 2013. Excalibur. <http://www.lighttwist.net/wp/excalibur>

878 Lin, S.-W., and V. M. Bier. 2008. A study of expert overconfidence. *Reliability*
879 *Engineering & System Safety* **93**:711-721.

880 Lorenz, J., H. Rauhut, F. Schweitzer, and D. Helbing. 2011. How social influence can
881 undermine the wisdom of crowd effect. *Proceedings of the National Academy of*
882 *Sciences* **108**:9020-9025.

883 Low Choy, S., R. O'Leary, and K. Mengersen. 2009. Elicitation by design in ecology:
884 using expert opinion to inform priors for Bayesian statistical models. *Ecology*
885 **90**:265-277.

886 MacDonald, J. A., M. J. Small, and M. Morgan. 2008. Explosion probability of
887 unexploded ordnance: expert beliefs. *Risk Analysis* **28**:825-841.

888 Martin, T. G., M. A. Burgman, F. Fidler, P. M. Kuhnert, S. Low-Choy, M. McBride, and
889 K. Mengersen. 2012. Eliciting expert knowledge in conservation science.
890 *Conservation Biology* **26**:29-38.

891 Martin, T. G., B. A. Wintle, J. R. Rhodes, P. M. Kuhnert, S. A. Field, S. J. Low-Choy, A.
892 J. Tyre, and H. P. Possingham. 2005. Zero tolerance ecology: improving ecological
893 inference by modelling the source of zero observations. *Ecology Letters* **8**:1235-
894 1246.

895 McBride, M. F., F. Fidler, and M. A. Burgman. 2012a. Evaluating the accuracy and
896 calibration of expert predictions under uncertainty: predicting the outcomes of
897 ecological research. *Diversity and Distributions* **18**:782-794.

898 McBride, M. F., S. T. Garnett, J. K. Szabo, A. H. Burbidge, S. H. Butchart, L. Christidis,
899 G. Dutson, H. A. Ford, R. H. Loyn, and D. M. Watson. 2012b. Structured elicitation
900 of expert judgments for threatened species assessment: a case study on a continental
901 scale using email. *Methods in Ecology and Evolution* **3**:906-920.

902 Meek, M. H., C. Wells, K. M. Tomalty, J. Ashander, E. M. Cole, D. A. Gille, B. J. Putman,
903 J. P. Rose, M. S. Savoca, and L. Yamane. 2015. Fear of failure in conservation: the
904 problem and potential solutions to aid conservation of extremely small populations.
905 *Biological Conservation* **184**:209-217.

906 Mellers, B., E. Stone, T. Murray, A. Minster, N. Rohrbaugh, M. Bishop, E. Chen, J. Baker,
907 Y. Hou, and M. Horowitz. 2015. Identifying and cultivating superforecasters as a
908 method of improving probabilistic predictions. *Perspectives on Psychological*
909 *Science* **10**:267-281.

910 Metcalf, S. J., and K. J. Wallace. 2013. Ranking biodiversity risk factors using expert
911 groups – Treating linguistic uncertainty and documenting epistemic uncertainty.
912 *Biological Conservation* **162**:1-8.

913 Morgan, M. G., and M. Henrion. 1990. *Uncertainty: A guide to dealing with uncertainty*
914 *in quantitative risk and policy analysis* Cambridge University Press. New York,
915 NY, United States of America.

916 Quigley, J., A. Colson, W. Aspinall, and R. M. Cooke. 2018. Elicitation in the Classical
917 Model. Pages 15-36 *in* L. C. Dias, A. Morton, and J. Quigley, editors. *Elicitation:*
918 *The science and art of structuring judgement*. Springer International Publishing,
919 Cham, Switzerland.

920 Rothlisberger, J. D., D. M. Lodge, R. M. Cooke, and D. C. Finnoff. 2009. Future declines
921 of the binational Laurentian Great Lakes fisheries: the importance of environmental
922 and cultural change. *Frontiers in Ecology and the Environment* **8**:239-244.

923 Satopää, V. A., J. Baron, D. P. Foster, B. A. Mellers, P. E. Tetlock, and L. H. Ungar. 2014.
924 Combining multiple probability predictions using a simple logit model.
925 International journal of forecasting **30**:344-356.

926 Shanteau, J., D. J. Weiss, R. P. Thomas, and J. C. Pounds. 2002. Performance-based
927 assessment of expertise: How to decide if someone is an expert or not. European
928 Journal of Operational Research **136**:253-263.

929 Speirs-Bridge, A., F. Fidler, M. McBride, L. Flander, G. Cumming, and M. Burgman.
930 2010. Reducing overconfidence in the interval judgments of experts. Risk Analysis
931 **30**:512-523.

932 Surowiecki, J. 2004. The wisdom of crowds: Why the many are smarter than the few and
933 how collective wisdom shapes business, economies, societies, and nations. Little,
934 Brown, London, United Kingdom.

935 Tetlock, P. 2005. Expert political judgement: How good is it? How can we know?
936 Princetown University Press, Princetown, United States of America.

937 Tetlock, P., and D. Gardner. 2015. Superforecasting: The art and science of prediction.
938 Random House, New York.

939 Walls, L., and J. Quigley. 2001. Building prior distributions to support Bayesian reliability
940 growth modelling using expert judgement. Reliability Engineering & System
941 Safety **74**:117-128.

942 Wallsten, T. S., D. V. Budescu, A. Rapoport, R. Zwick, and B. Forsyth. 1986. Measuring
943 the vague meanings of probability terms. Journal of Experimental Psychology:
944 General **115**:348.

945 Ward, T. 2014. The rapid assessment workshop to elicit expert consensus to inform the
946 development of the Great Barrier Reef Outlook Report 2014. Townsville.

947 Weiss, D. J., and J. Shanteau. 2004. The vice of consensus and the virtue of consistency.
948 Psychological investigations of competent decision making:226-240.

949 Winkler, R. L., and A. H. Murphy. 1968. “Good” probability assessors. Journal of applied
950 Meteorology 7:751-758.

951 Wittmann, M. E., R. M. Cooke, J. D. Rothlisberger, E. S. Rutherford, H. Zhang, D. M.
952 Mason, and D. M. Lodge. 2015. Use of structured expert judgment to forecast
953 invasions by bighead and silver carp in Lake Erie. Conservation Biology 29:187-
954 197.

955 Yaniv, I., and D. P. Foster. 1997. Precision and accuracy of judgmental estimation. Journal
956 of Behavioral Decision Making 10:21-32.

957

958 **8 Data availability statement**

959 Open Access: All data and code for the analyses presented in this paper are available on the
960 Open Science Framework (Hemming 2019).

961

962

963

964

965

Box 1 Scoring rules IDEA protocol vs. The Classical Model

In Figure 2, two experts have been asked to provide their estimates for 10 calibration questions. They have then been scored on their performance using the scoring rules outlined Section 2.4 from the Classical Model and the IDEA protocol.

Statistical accuracy (Classical Model) vs Calibration (IDEA)

Expert A, has an inter-quantile distribution of $s(A) = (0.10, 0.40, 0.40, 0.10)$, that is, over 10 questions one realisation fell below their 5th interval, four between their 5th and their 50th, four between their 50th and their 95th, and one above their 95th. When compared to the theoretically optimal inter-quantile distribution of $p = (0.05, 0.45, 0.45, 0.05)$, using a chi-squared test with three-degrees of freedom they receive a statistical accuracy (SA) of 0.83, which is the highest statistical accuracy that can be achieved on 10 questions.

Expert B, provides a theoretical distribution $s(B) = (0.10, 0.90, 0.0, 0.0)$, which is quite different to the theoretically optimal inter-quantile distribution p . Their statistical accuracy is low, 0.003. Having a statistical accuracy below 0.05 they would be deemed statistically inaccurate under the Classical Model.

In contrast, when scored using calibration (CA) from the IDEA protocol, Expert B would be perfectly calibrated having nine of their ten 90% credible intervals capturing the realised truth. Expert A would also be considered well-calibrated, but less so than Expert B, only capturing eight out of 10 realisations in their 90% credible intervals.

Information (Classical Model) vs Informativeness (Four-step question format)

Expert A and B provide intervals which are exactly the same width for each question. However, Expert B consistently provides a median close to the tails. This means the mass of their intervals departs from a uniform distribution whereby we would expect 5% of the total width of their interval to fall below their 5th quantile, 45% between their 5th and 50th, and again between their 50th and 95th, and 5% above their 95th quantile. Assuming this is the only difference in their intervals, Expert B would achieve a higher information score under the Classical Model than Expert A. However, as experts have intervals that are the same width, both experts would receive the same score for informativeness under the IDEA protocol.

967

968

969 **Figure 1** Key steps of the IDEA protocol (figure from Hemming et al. (2018b)). The
970 four-step question format (Speirs-Bridge et al. 2010) (depicted in Step 2) is commonly
971 used to derive a best estimate and credible interval in Round 1 and Round 2.

972 **Figure 2** Judgements provided by two hypothetical experts over 10 questions. The
973 blue lines represent their 90% credible intervals, the blue dots their ‘best estimate’ or
974 their ‘median’. The crosses represent where the realisation fell in relation to their
975 estimates. To calculate statistical accuracy (SA) according to the Classical Model, the
976 proportion of questions answered where realisations fell, 1) below their lowest interval
977 (i.e. 5th quantile), 2) between their lowest estimate and their best estimate / median, 3)
978 between their best estimate / median and their upper estimate / 95th quantile, and 4) above
979 their upper / 95th quantile is calculated and compared to a theoretically optimal
980 distribution $p=(0.05, 0.45, 0.45, 0.05)$. CA refers to calibration as calculated according to
981 the IDEA protocol, which is defined as the proportion of credible intervals capturing the
982 realisation.

983 **Figure 3** The statistical accuracy and information of $n = 58$ participants. A trade-off
984 exists between the two measures used by the Classical Model. Those who are statistically
985 accurate (above 0.05, red horizontal line) often have a lower information score than the
986 median score for individuals (grey vertical line). The blue dashed line shows the highest
987 statistical accuracy score possible for 13 questions (0.93), and the black line shows the
988 highest score obtained by individuals in the elicitation (0.53).

989 **Figure 4** Statistical accuracy of the Classical Model (CM) compared to IDEA
990 calibration for $n = 58$ participants. The graph shows that participants with perfect
991 calibration when assessed by the IDEA protocol, can have poor statistical accuracy for
992 the Classical Model. On the righthand side, we show where the realisations fell in each of

993 the expert's multinomial distributions (used to calculate statistical accuracy), and
994 contrast this with how many realisations fell within the participant's 90% credible
995 intervals (calibration). Bold numbers indicate the highest scores possible for statistical
996 accuracy and calibration.

997 **Figure 5** The spearman correlation between information calculated for the Classical
998 Model, and informativeness calculated for the IDEA protocol for $n = 58$ participants. The
999 shaded area represents a 95% confidence interval.

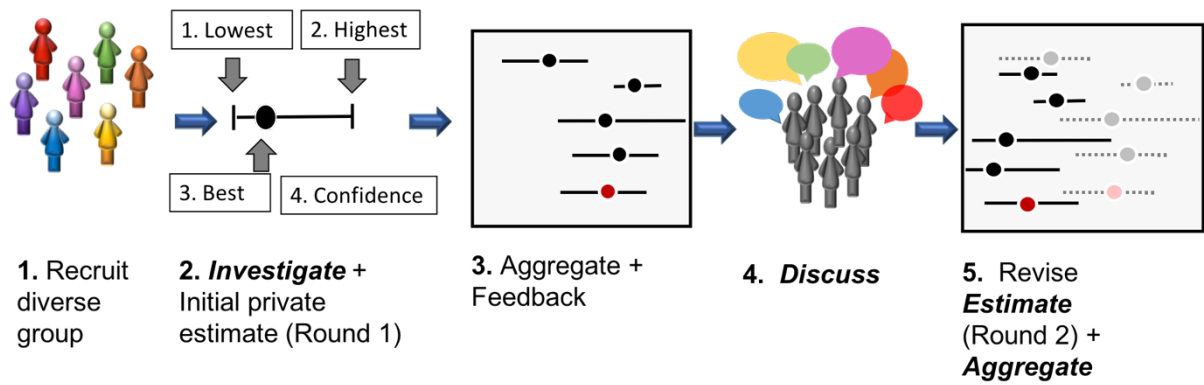
1000 **Figure 6** CM Scores derived for each aggregation.

1001 **Figure 7** Component performance measures of the Classical Model (CM) and IDEA
1002 protocol for $n = 8$ groups under six alternative procedures for aggregation. a) Statistical
1003 accuracy, the red-dashed line represents the 0.05 threshold for statistically inaccurate
1004 scores, (Classical Model), the blue dashed line represents a perfect statistical accuracy
1005 score for 13 questions, and the black dashed line represents the highest score obtained by
1006 any individual, b) information score (Classical Model), the red line represents the median
1007 information of an individual c) calibration, (IDEA) the red line represents perfect
1008 calibration (0.90), d) informativeness (IDEA), the red line represents the informativeness
1009 of the median individual, e) accuracy (IDEA) of the best estimate, the red line represents
1010 the accuracy of the median individual.

1011 **Figure 8** The scores of aggregations under the Classical Model and the IDEA
1012 protocol when adjustments are made to correct for questions for which the realised truth
1013 had been zero. a) Statistical accuracy, the red-dashed line represents the 0.05 threshold
1014 for statistically inaccurate scores, (Classical Model), the blue dashed line represents a
1015 perfect statistical accuracy score for 13 questions, and the black dashed line represents

1016 **the highest score obtained by any individual prior to the adjustment; b) calibration,**
1017 **(IDEA) the red line represents perfect calibration (0.90).**

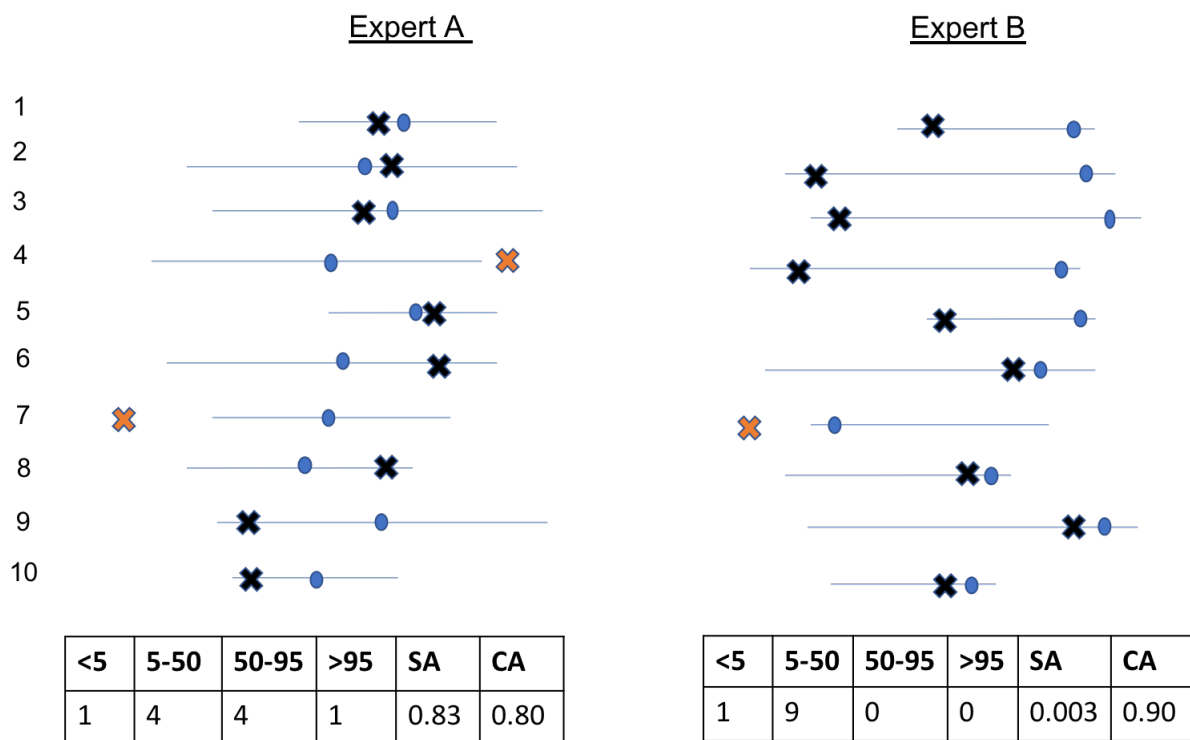
1018 **Figure 1**



1019

1020

1021 **Figure 2**



1022

1023

1024

1025

1026

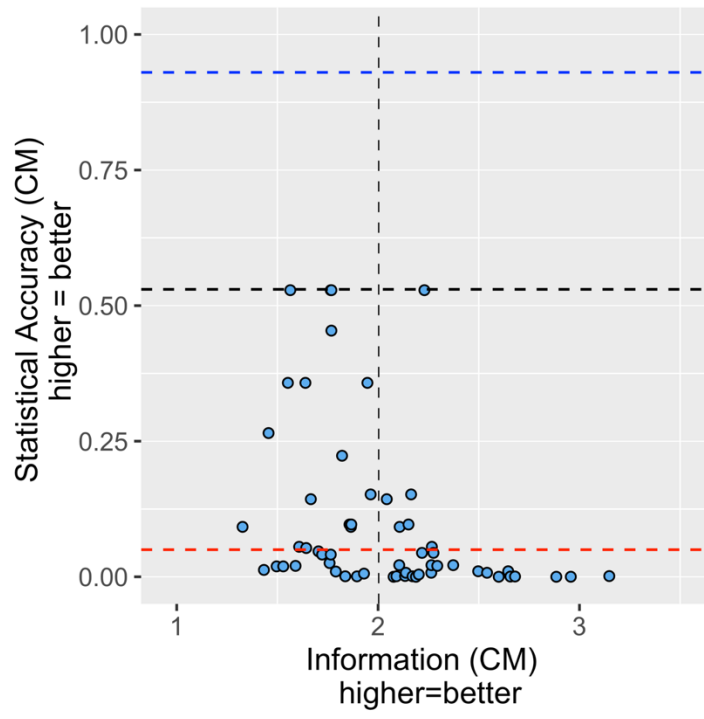
1027

1028

1029

1030

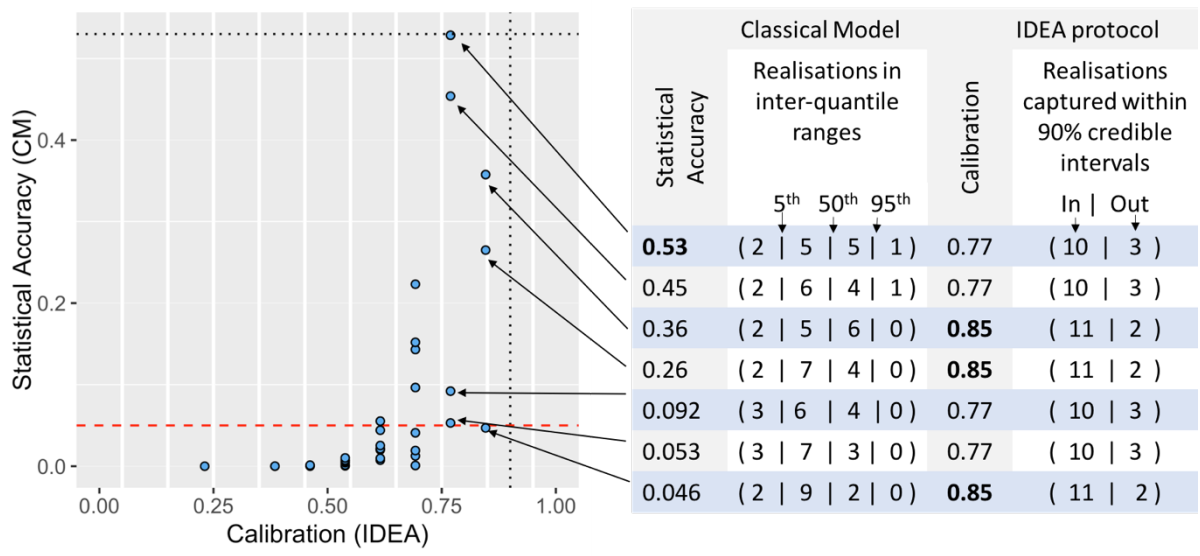
1031 **Figure 3**



1032

1033

1034 **Figure 4**



1035

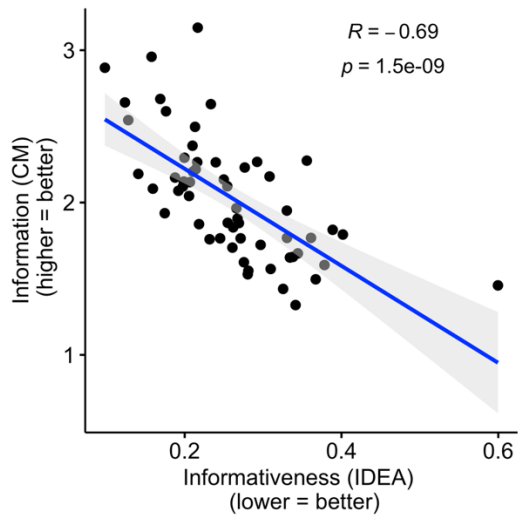
1036

1037

1038

1039

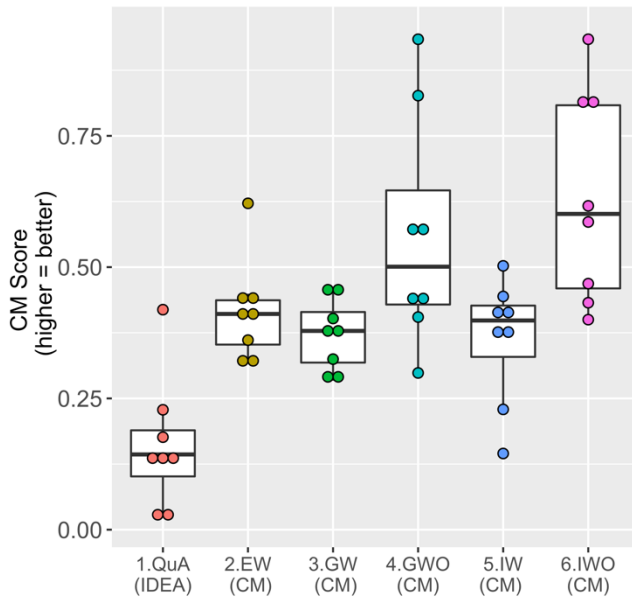
1040 **Figure 5**



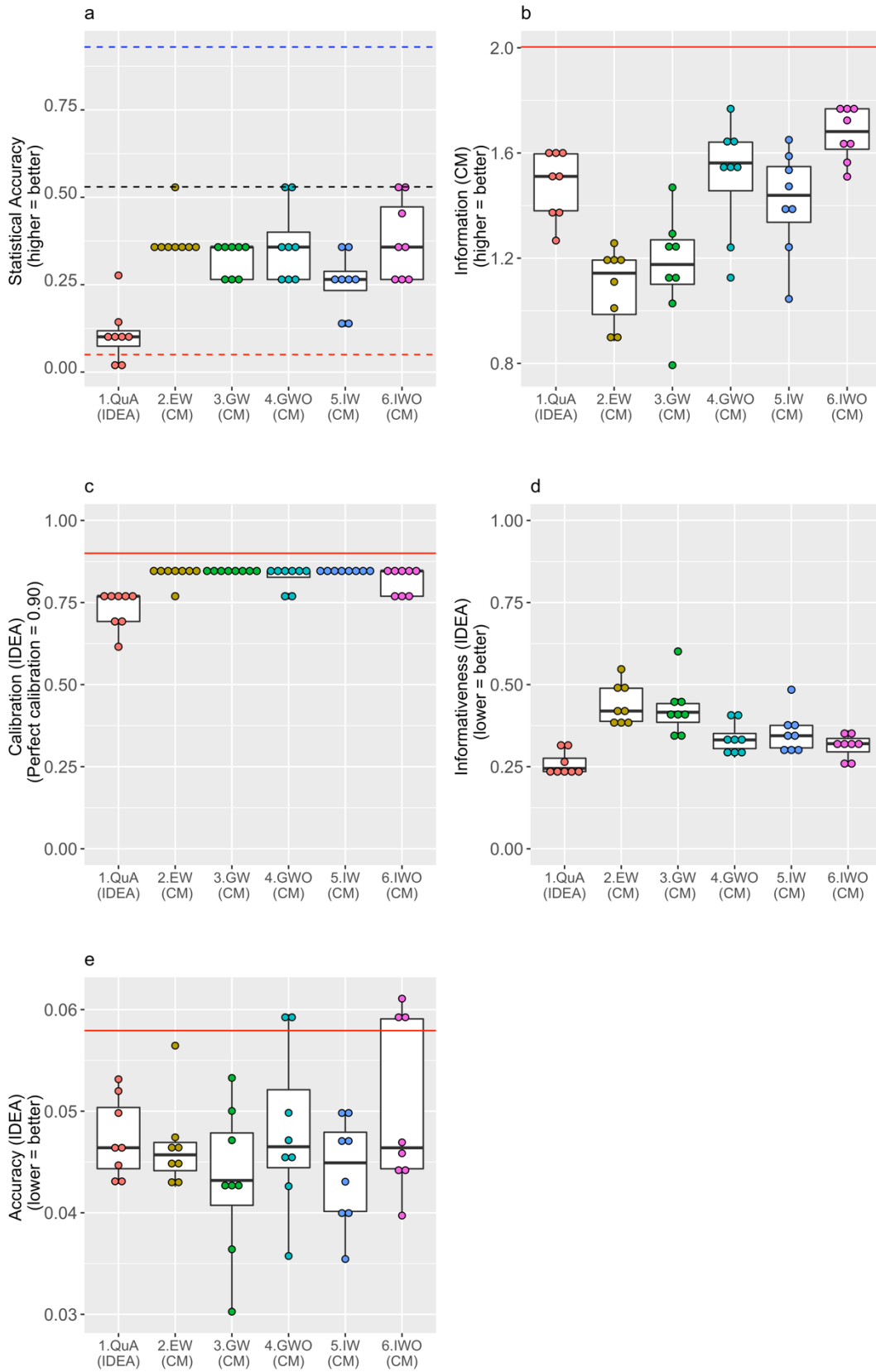
1041

1042

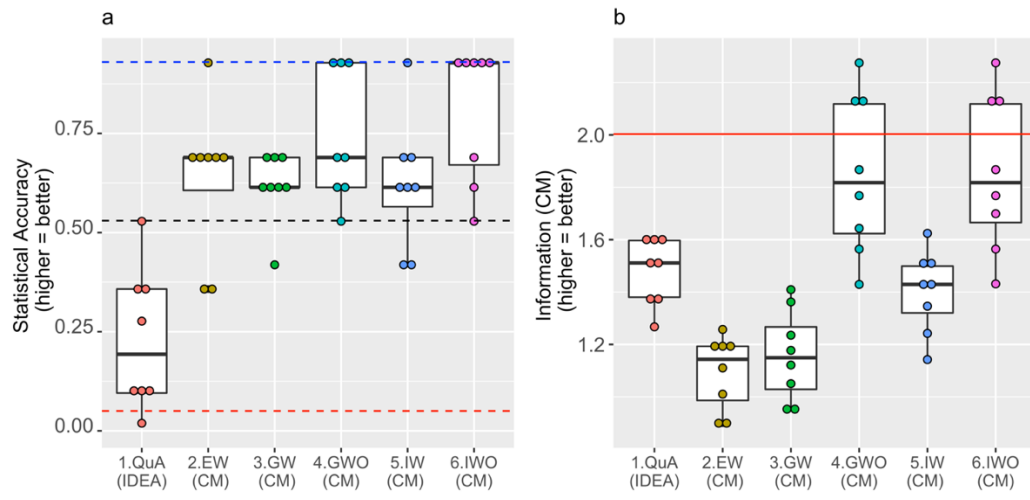
1043 **Figure 6**



1044



1047 **Figure 8**



1048

1049