



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Hemming, V;Hanea, AM;Walshe, T;Burgman, MA

Title:

Weighting and aggregating expert ecological judgments

Date:

2020-06-01

Citation:

Hemming, V., Hanea, A. M., Walshe, T. & Burgman, M. A. (2020). Weighting and aggregating expert ecological judgments. *Ecological Applications*, 30 (4), <https://doi.org/10.1002/eap.2075>.

Persistent Link:

<https://hdl.handle.net/11343/275290>

1

2 DR. VICTORIA HEMMING (Orcid ID : 0000-0003-3220-6161)

3

4

5 Article type : Articles

6

7

8 Running head: Weighting and aggregating experts

9 **Weighting and aggregating expert ecological judgments**10 Victoria Hemming^{1,2,3,*}, Anca M. Hanea^{1,2}, Terry Walshe², Mark A. Burgman^{2,4}11 ¹ The Centre of Excellence for Biosecurity Risk Analysis, The University of Melbourne,
12 Melbourne, Victoria, Australia13 ² School of Biosciences, The University of Melbourne, Melbourne, Victoria, Australia.14 ³ Department of Forest and Conservation Sciences, The University of British Columbia,
15 Vancouver, Canada.16 ⁴ The Centre for Environmental Policy, Imperial College London, London, United Kingdom17 *Corresponding author. Email victoria.hemming@ubc.ca, ORCID ID: 0000-0003-3220-
18 6161.

19

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/eap.2075](https://doi.org/10.1002/eap.2075)

This article is protected by copyright. All rights reserved

20

21 **Abstract**

22 Performance weighted aggregation of expert judgments, using calibration questions, has been
23 advocated to improve pooled quantitative judgments for ecological questions. However, there
24 is little discussion or practical advice in the ecological literature regarding the application,
25 advantages or challenges of performance weighting. In this paper we 1) illustrate how the
26 IDEA protocol with four-step question format can be extended to include performance
27 weighted aggregation from the Classical Model, and 2) explore the extent to which this
28 extension improves pooled judgments for a range of performance measures. Our case study
29 demonstrates that performance weights can improve judgments derived from the IDEA
30 protocol with four-step question format. However, there is no *a-priori* guarantee of
31 improvement. We conclude that the merits of the method lie in demonstrating that the final
32 aggregation of judgments provides the best representation of uncertainty (i.e. validation),
33 whether that be via equally weighted or performance weighted aggregation. Whether the time
34 and effort entailed in performance weights can be justified is a matter for decision-makers.
35 Our case study outlines the rationale, challenges, and benefits of performance weighted
36 aggregations. It will help to inform decisions about the deployment of performance weighting
37 and avoid common pitfalls in its application.

38

39 Key words: performance weights, equal weights, aggregation, expert judgment, calibration,
40 Classical Model

41

42

43 **1 Introduction**

44 Over the past 15 years a considerable body of research has emerged in the ecological
45 literature emphasizing the need for more rigorous and structured methods for collecting
46 quantitative expert judgments. The literature has summarised existing structured elicitation
47 protocols and key steps which could be adapted and applied to better suit the practical (e.g.
48 geographically dispersed experts) and financial (lack of funding) constraints of most

49 ecological contexts (Burgman 2004, Low Choy et al. 2009, Kuhnert et al. 2010, Burgman et
50 al. 2011a, Martin et al. 2012, McBride et al. 2012a, McBride et al. 2012b, Drescher et al.
51 2013).

52 A common approach that has been advocated is to recruit a diverse group of individuals and
53 take an equally weighted aggregation of their independent judgments (Burgman et al. 2011b,
54 Hemming et al. 2018b). This will often produce point estimates which are at least as accurate
55 (i.e. closer to the truth) and interval judgments which are better calibrated than the median-
56 ranked individual for these scores (Burgman et al. 2011b, Budescu and Chen 2014, Hemming
57 et al. 2018b). While one person can sometimes outperform the group aggregate, rarely can
58 that person be predicted by credentials conventionally associated with expertise such as age,
59 experience, or peer-identification (Aspinall and Cooke 2013, Burgman 2015, Mellers et al.
60 2015).

61 The performance of the equal weighted aggregation is largely explained as a statistical
62 phenomenon (Lorenz et al. 2011) in which the judgments of individuals represent random
63 independent samples. If those samples are diverse then not only should the information pool
64 related to the questions increase (Clemen and Winkler 1999), but the errors made by
65 individuals are more likely to cancel (Larrick and Soll 2006, Budescu and Chen 2014). This
66 phenomenon is often termed the ‘wisdom of the crowd’ (Surowiecki 2004), or the ‘staticised
67 group’ (Einhorn et al. 1977, Hogarth 1978). Interestingly, participants need not be experts
68 and can be biased, as long as they have some information related to the questions that can be
69 combined for prediction (Budescu and Chen 2014).

70 Equal weighting is advantageous as it’s relatively simple to apply (Hogarth 1978, Hora 2004,
71 Hemming et al. 2018b). Typically, group sizes of 5-12 participants derive improved
72 judgments, with diminishing returns thereafter (Hogarth 1978, Hora 2004, Hemming et al.
73 2018b). It requires no additional work to develop questions or performance measures to score
74 and aggregate experts. It can be applied to any type of prediction including point estimates,
75 distributions and probabilities. The simplicity of equal weighting, and its ability to improve a
76 wide range of estimates make it suitable for aggregating judgments under the practical and
77 financial constraints typical of many ecological decisions.

78 However, despite substantial testing and real-world applications, many people find equal
79 weighting difficult to trust (Weiss and Shanteau 2004). This is partly because the method
80 relies on the recruitment of a diversity of individuals, often including individuals who may

81 normally be excluded from such elicitations because of their perceived limited knowledge
82 (Shanteau et al. 2002, Weiss and Shanteau 2004, Burgman et al. 2011a).

83 When uncertainty is elicited, the diversity of the group can also increase the uncertainty
84 associated with group judgments, sometimes leading to uninformative judgments
85 (MacDonald et al. 2008, Barons et al. 2018). Occasionally individuals will outperform the
86 group aggregation, and ideally decision-makers would like to restrict elicitation to these
87 better performing individuals, or at least have the judgments of those individuals weigh more
88 than those of lesser performers. Finally, there is no single method for generating an equally
89 weighted aggregation. For example, for point estimates, the arithmetic mean is commonly
90 applied, but one could also use the median, geometric mean or harmonic mean (Armstrong
91 2001, Colson and Cooke 2017). Rarely is there any validation to support such choices made
92 by the analyst, which can lead to questions about the validity of the specific method chosen,
93 and the influence of analyst's subjective bias. The problems associated with equal weights
94 can serve to undermine the credibility of the final judgments derived and the subsequent
95 decisions and assessment based on such judgments.

96 Performance weighted aggregation is often suggested as a way to address these challenges
97 and perceived deficiencies (Cooke 1991, Budescu and Chen 2014, Mellers et al. 2015). It
98 involves developing sets of questions related to the main elicitation questions for which the
99 answers can be obtained but are not widely known to experts (Cooke 1991, Goossens et al.
100 2008, Tetlock and Gardner 2015). These are referred to as test, seed or calibration questions
101 (we use the term calibration questions from hereon). Those who perform better on these
102 questions are afforded more weight in the final aggregation of the main elicitation questions.
103 The method is differentiated from other forms of weighted aggregation such as those based
104 on self-rating, peer-rating, trimming, or representativeness in that weights are obtained via
105 validation of judgments against an external truth (Armstrong 2001, Aspinall and Cooke
106 2013).

107 The main reason decision-makers seek to apply performance weights is to create aggregated
108 judgments which are more accurate (for point estimates), or well-calibrated and informative
109 (for interval judgments, probabilities and probability distributions) (Budescu and Chen 2014,
110 Mellers et al. 2015, Colson and Cooke 2017). However, the inclusion of calibration questions
111 is also seen to create a sense of legitimacy. It provides evidence that those who have been
112 included in the final aggregation have some knowledge in the relevant domain, and that they

113 can communicate their knowledge together with their uncertainty in the format required by
114 the analyst (Barons et al. 2018, Quigley et al. 2018). It can also be used to validate
115 assumptions made by the analyst in combining expert judgments.

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

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

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

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

137 Most applications of the IDEA protocol in ecology aim to obtain quantitative judgments
138 together with uncertainty. When doing so, the four-step question format is often deployed
139 (Speirs-Bridge et al. 2010) (Figure 1). This method derives a credible interval with a ‘best’
140 point estimate based on the following questions:

- 141 1. Realistically what is the lowest plausible value for x?
- 142 2. Realistically what is the highest plausible value for x?
- 143 3. Realistically what is your best estimate for x?

144 4. Looking at your interval from lowest to highest, how confident are you that your
145 interval will capture the realised truth.

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

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

159 The Classical Model (Cooke 1991) is a method for performance weighted aggregation often
160 cited in the ecological literature as a means to improve uncertain quantitative ecological
161 judgments (Burgman et al. 2011a, Martin et al. 2012, Drescher et al. 2013, Metcalf and
162 Wallace 2013, Hemming et al. 2018a). While it has been applied to a large number of
163 engineering case studies (Cooke and Goossens 2008, Colson and Cooke 2017) we are aware
164 of only two ecological examples, both in the Laurentian Great Lakes (Rothlisberger et al.
165 2009, Wittmann et al. 2015).

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

171 **2 Methods**

172 **2.1 Fundamentals of performance weighting**

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

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

180 There is little prescriptive guidance as to what makes a good calibration question, although
181 some features are self-evident (Cooke and Goossens 2000, Aspinall and Cooke 2013, Tetlock
182 and Gardner 2015, Quigley et al. 2018). They should relate to the knowledge needed to
183 answer the main elicitation questions (i.e. domain knowledge). They should ask questions
184 about uncertainty to capture an expert's ability to adapt and communicate their knowledge.
185 They should be in a similar format to the main elicitation questions. They should not be
186 questions which can be easily guessed, and not so hard that an expert could not reasonably
187 form a judgment. A substantial number of calibration questions may be required to
188 differentiate luck from good judgment, depending on the scoring rules. Ideally, questions
189 should relate to predictions of events or quantities rather than estimating the outcomes of past
190 events (retrodictions), although this is not always possible. The questions should be reviewed
191 by at least two people with domain knowledge to ensure they provide fair and reasonable
192 assessments of an expert's ability to make good judgments related to the main elicitation
193 questions.

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

205 al. 2018). The other exception is the scoring rule of the Classical Model (discussed below)
206 (Cooke 1991).

207 Scoring rules aim to optimize judgments and the way in which they do this depends on their
208 reward structure (Winkler and Murphy 1968, Tetlock 2005). For example, scoring rules for
209 interval judgments may penalise overconfidence (e.g. intervals that are too narrow, which
210 include the truth less often than the purported level of confidence provided by the expert)
211 more than under-confidence (e.g. intervals that are too broad and capture more realisations
212 than the purported level of confidence of the expert). It's therefore important to understand
213 how such transgressions of judgment are handled by a proposed scoring rule, to ensure that
214 the reward structure matches the preferences and needs of the decision-maker and the
215 problem at hand. This of course requires an awareness among decision-makers about what
216 aspects of judgment are most important to them.

217 Obtaining an understanding of the reward structure can be challenging as research papers
218 outlining the application of scoring rules rarely provide clear examples of how judgments are
219 incorporated and combined. Few adequately discuss their embedded reward structure. A
220 further complication arises in understanding scoring rules because the terms used to describe
221 judgment, such as 'calibration', 'accuracy' and 'overconfidence', are used interchangeably
222 and may refer to different concepts (Lichtenstein and Fischhoff 1977, Lichtendahl Jr et al.
223 2013, Cooke 2018b, Hemming et al. 2018b).

224 **2.2 The Classical Model**

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

231 In elicitations employing the Classical Model, experts are asked a set of calibration questions
232 (usually 10-15), for which the answers can be obtained. As noted above, these questions
233 should relate to the main questions of the elicitation (termed target variables or questions of
234 interest). Unlike the four-step question format commonly used with the IDEA protocol,

235 experts are asked to specify their judgments as quantiles of a continuous non-parametric
236 probability distribution (usually 5th, 50th, and 95th) for both calibration questions and
237 questions of interest. The individual judgments of experts are typically elicited in a face-to-
238 face elicitation with one or more facilitators present (Wittmann et al. 2015). Experts are
239 scored on their performance using two performance measures (see section 2.4 for details):
240 “statistical accuracy” (often termed “calibration”), and “information” (sometimes termed
241 “informativeness”, or “relative information”). These are subsequently multiplied to provide
242 an asymptotically proper scoring rule (the CM Score) (refer to Appendix S1: Section 4.2.3),
243 and to derive differential weights.

244 Experts who perform well on the calibration questions are afforded more weight in the final
245 aggregations for the questions of interest. Both equally weighted and performance weighted
246 linear pooled aggregations of distributions are then created and subsequently scored on their
247 performance on the calibration questions (i.e. via in-sample validation, where the same
248 questions used to develop the performance weighted aggregations are used to score the
249 aggregations). To achieve rational consensus, experts or decision makers usually agree prior
250 to the elicitation that the aggregation which achieves the highest combined score on the
251 calibration questions will be used to weight expert judgments of the target questions.

252 The primary purpose of performance weighting and calibration questions in the Classical
253 Model is to come to an unbiased and empirically validated decision on how to combine the
254 expert judgments. This step can help to overcome pre-judgments and exclusion of potentially
255 knowledgeable individuals, as well arbitrary choices by analysts and decision makers about
256 how to weight and aggregate experts. In analyses of 78 case studies using the Classical
257 Model, performance weighted aggregations have outperformed equal weights in 76 studies
258 (in-sample validation), suggesting the method can also be used to optimize aggregated
259 judgments (Cooke and Goossens 2008, Colson and Cooke 2017).

260 **2.3 Case study**

261 To demonstrate how the Classical Model could be applied in ecology, and to investigate
262 potential improvements from its application, we use estimates for ecological questions from a
263 previous case study by Hemming et al. (2018b). In brief, the case study used the IDEA
264 protocol with the four-step question format to elicit judgments for thirteen questions relating
265 to future abiotic and biotic events on the Great Barrier Reef. The elicitation was undertaken

266 via email and the experts volunteered their time. The questions related to the types of events
267 experts may be asked in assessing risk to the Great Barrier Reef (Ward 2014), for example,
268 the percentage cover of coral bleaching that may be detected in the next monitoring event at a
269 specified reef (see Appendix S1: Section 1). The questions related to future monitoring
270 events, so that judgments could be scored against outcomes once monitoring data were
271 collected.

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

284 **2.3.1 Four-step to quantiles**

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

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

297 • Lower / 5th : 0.00001

298 • Best / 50th : 0.0001

299 • Upper / 95th: 0.001

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

303 2.4 Scoring Judgments

304 Assuming the judgments approximate quantiles of a continuous probability distribution, the
305 judgments can then be scored using the Classical Model's performance measures. There is
306 substantial ambiguity and confusion in the ecological literature as to what the performance
307 measures of the Classical Model actually reward. They have been cited as rewarding
308 'accuracy' (Rothlisberger et al. 2009, Burgman et al. 2011a, Martin et al. 2012), which may
309 give the impression they reward the accuracy of point estimates. They have also been noted
310 to score 'calibration' and 'precision' (width) which may give the impression they are
311 designed to assess interval judgments according to definitions that arise in the psychological
312 literature (Lichtenstein and Fischhoff 1977, Yaniv and Foster 1997, Burgman et al. 2011a,
313 Wittmann et al. 2015).

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

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

326 ***IDEA performance measures***

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

329 *Accuracy* is designed to assesses the accuracy of point estimates. It is the difference between
330 b , the expert's best estimate, and the observed value, x . It is measured using the average log
331 ratio error (ALRE) of expert responses. The measure is a relative measure, scale invariant,
332 and emphasizes order of magnitude errors rather than linear errors. Smaller ALRE scores
333 indicate more accurate responses. For any given question the log ratio score has a maximum
334 possible range of 0.31 ($=\log_{10}(2)$), which occurs when the true answer coincides with either
335 the group minimum or group maximum (Burgman et al. 2011b)

336 *Calibration* is the proportion of intervals provided by the experts containing the realised truth
337 relative to their assigned confidence (Lichtenstein and Fischhoff 1977, Lin and Bier 2008).
338 For example, if the expert's intervals are standardized to 90% credible intervals then we
339 expect for a well calibrated expert and 100 questions, that 90 of the realisations will fall
340 between their 5th and 95th quantiles. If they capture fewer realisations, they may be
341 considered overconfident, and if they capture more realisations they may be considered
342 underconfident. The measure is an absolute measure and is scale invariant. If the realisations
343 are equal to the expert's 5th or 95th quantiles, then they are usually assessed as being included
344 within the expert's credible intervals.

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

352 ***Performance measures of the Classical Model***

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

355 *Statistical accuracy* (often referred to as calibration and often denoted by ‘C’) assesses the
356 ability of experts to answer according to a theoretically optimal multinomial distribution. It
357 assesses the interquantile calibration of experts. For example, over a set of questions for
358 which realisations could be obtained, we would expect for any high performing expert that:

- 359 • For 5% of their judgments, the realisations would fall below their 5th quantile. We express
360 the observed proportion as Q_1 .
- 361 • For 45% of their judgments, the realisations would fall between their 5th and their 50th
362 quantile. We express the observed proportion as Q_2 .
- 363 • For 45% of their judgments, the realisations would fall between their 50th and their 95th
364 quantile. We express the observed proportion as Q_3 .
- 365 • For 5% of their judgments, the realisations would fall above their 95th quantile. We
366 express the observed proportion as Q_4 .

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

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

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

- 384 • If the realisation equals the 5th quantile, it is placed in the first probability bin Q_1 .
- 385 • If the realisation equals the 50th quantile, it is placed in the second probability bin Q_2 .

386 • If the realisation equals the 95th, it is placed in the third probability bin Q_3 .
387 We highlight this assumption as (on rare occasions) it can affect the score participants
388 receive. For example, in the unlikely case that a participant was to estimate the median
389 perfectly for 9 of 10 questions, they could obtain a multinomial distribution of $S(e) = (1, 9, 0,$
390 $0)$, which when compared to the theoretically optimal multinomial distribution means they
391 would be considered statistically inaccurate at the 0.05 level, despite having perfect
392 calibration and exceptional accuracy under the IDEA protocol scoring rules.

393 *Information* (often referred to as relative information, or informativeness) under the Classical
394 Model measures the degree to which the expert's distribution is concentrated and to which it
395 differs from a uniform or log-uniform distribution (which are considered the least informative
396 distributions). It uses the KL divergence measure, which is scale invariant (Quigley et al.
397 2018). Information is calculated per question and does not depend on the realisation. The
398 final information score of an expert is an average taken across all calibration questions.
399 Larger numbers indicate better performance because they represent distributions which show
400 greater departure from a uniform or log-uniform distribution.

401 A simple example contrasting the performance measures is provided in

402 Data availability statement

403 All data and code for the analyses presented in this paper are available on the Open Science
404 Framework (Hemming 2019): <https://doi.org/10.17605/OSF.IO/FXQVK>

405

406

407

408

409

410 **Box 1**, and Figure 2. In the results section, we plot outcomes for these measures against each
411 other to gain a better understanding of the underlying reward structures.

412 **2.5 Weighting and aggregating**

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

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

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

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

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

435 Global Weights Optimized (GWO) and Itemized Weights Optimized (IWO): are similar to
436 their un-optimized variants described above (i.e. Global Weights (GW) and Itemized Weights
437 (IW)). However, they optimize the statistical accuracy score by successively raising the level
438 at which an expert is considered statistically inaccurate from an alpha level equal to the
439 lowest calibration score. The weights are calculated and used to generate weighted
440 aggregations that are scored on the calibration questions. The weighted aggregation with the
441 highest performance on the calibration questions is chosen (Quigley et al. 2018). In decisions
442 with one or two well calibrated experts, most or all of the weight may be assigned to those
443 experts with no weight given to the other experts.

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

453 **2.5.1 Linear pooling versus quantile aggregation**

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

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

468 **2.6 Analysis**

469 For the eight groups of experts in our case study, we assessed the six alternative approaches
470 to aggregation (two forms of equal weighted aggregations (EW (Classical Model), QuA
471 (IDEA)), and four forms of performance weighted aggregation from the Classical Model
472 (IW, GW, IWO, GWO) (described in Section 2.5). Individual and group performance was
473 assessed using the five performance measures (statistical accuracy, information, calibration,

474 informativeness, and accuracy) (described in Section 2.4), and the Classical Model scoring
475 rule (CM Score)(described in Section 2.5).

476 To obtain the performance measures and aggregations associated with the Classical Model,
477 the analyst must enter judgments in software called *Excalibur* (Lightwist 2013, Cooke
478 2018a)(Appendix S1: Section 3). For measures associated with the IDEA protocol we
479 developed *R*-code (available on the Open Science Framework (Hemming 2019)). More
480 details are available in Appendix S1: Section 3.

481 To contrast the differences of the aggregations, we use boxplots, constructed in *R* (version
482 3.4.1 (2017-06-30) -- "Single Candle"), using the *ggplot2* package. The boxes represent the
483 25th, 50th and 75th percentiles. The whiskers represent the spread of the data referenced on
484 the inter-quartile range, $(Q1-1.5*IQR, Q3+1.5*IQR)$. For normally distributed data this is
485 approximately 2.7 standard deviations, or 99.3% of the data (Krzywinski and Altman 2014).

486 **3 Results**

487 **3.1 Comparison of performance measures**

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

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

498 Figure 4 shows the scatter between statistical accuracy (the Classical Model) and calibration
499 (IDEA Protocol) for the 58 participants over 13 questions. While the difference in the highest
500 statistical accuracy score possible and that obtained by experts appears large (i.e. a change
501 from 0.93 to 0.53 implies a 43% reduction in statistical accuracy), we can see that this change
502 was due to just one additional realisation falling outside of the experts' credible intervals.

503 Thus, statistical accuracy can be highly sensitive to seemingly small variations in
504 performance.

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

512 In Appendix S1: Section 7, we demonstrate that the differences occur because the Classical
513 Model's rules score a multinomial distribution with three degrees of freedom $p = (0.05, 0.45,$
514 $0.45, 0.05)$. As such, beyond very low levels of calibration (i.e. <50% calibration for 13
515 questions), the statistical accuracy measure cannot be used to assess the calibration of 90%
516 credible intervals (i.e. a multinomial distribution with one degree of freedom, or a binomial
517 distribution).

518 Figure 5 shows the correlation between information (Classical Model) and informativeness
519 (IDEA Protocol). The two scores are negatively correlated (Spearman rank correlation = -
520 0.69 , 95%CI: $-0.80, -0.52$), an artefact of the scoring rules, whereby under the IDEA protocol
521 participants who receive a low score are more informative (narrower intervals), whereas for
522 the Classical Model a higher score indicates that they provide more information relative to a
523 uniform or log-uniform distribution. Figure 5 demonstrates that information and
524 informativeness are slightly different measures of an expert's judgment. The Classical Model
525 does not only assess the width of intervals, it also accounts for their departure from a uniform
526 distribution. This can mean that higher information score may be obtained in some cases
527 simply by reducing the symmetry of the ranges between an expert's 2nd and 3rd quantiles (i.e.
528 if the median does not fall squarely in the centre of the range then information can be
529 increased).

530 **3.2 Performance of aggregations**

531 Figure 6 shows the CM Score for each of the aggregations. In the Classical Model, this
532 combined score would be used to select the final aggregation for uncertainty by a decision-

533 maker. For this case study, if we were to use the median values of these scores, we would not
534 choose quantile aggregation (QuA) because it has a low CM Score (median value of 0.14).
535 Equally weighted linear pooling employed by the Classical Model does better (EW) (median
536 value of 0.41), and there is some indication that performance weighted aggregation by the
537 optimized variants (IWO, and GWO) may lead to further improvements (median values of
538 0.60 and 0.50 respectively).

539 Figures 7a and 7b decompose the CM Score provided in Figure 6, into statistical accuracy
540 and information scores of the Classical Model. While quantile aggregation (QuA) performs
541 well on information (median value of 1.51), it performs poorly in terms of statistical accuracy
542 (median value of 0.10) compared to equal weights (EW) (median value of 0.36) and
543 performance weighted aggregations (IW, IWO, GW, GWO) (all achieving a median value of
544 0.36, except for IW which achieves 0.27), with two groups considered statistically inaccurate
545 at the 0.05 level. This supports the finding by Bamber et al. (2016) and Colson and Cooke
546 (2017) that quantile aggregation used in the IDEA protocol with four-step question format
547 can be overconfident relative to linear-pooling of distributions when assessed by statistical
548 accuracy.

549 There is little or no difference in the median performance of equally weighted (EW) and the
550 performance weighted aggregations (IW, IWO, GW, GWO) in terms of statistical accuracy.
551 However, both the optimized aggregations (GWO, and IWO) and itemized weights (IW)
552 have higher information (1.56 and 1.68) than equal weights (EW, median of 1.14) or global
553 weights (GW, median of 1.18), and are equivalent to quantile aggregation (QuA, 1.51)
554 suggesting performance weighting improves estimates in this case study by being more
555 informative than equal weights (EW).

556 Figures 7c-e assess each of the aggregation methods according to measures commonly used
557 in the IDEA protocol. Even when scored according to calibration between the expert's 90%
558 credible intervals, the study finds that quantile aggregation (QuA) generates more
559 overconfident estimates (median calibration of 0.77), having a lower calibration than all other
560 aggregations (0.85, or on average by one question). It does, however, have a higher level of
561 informativeness (0.25) than all other aggregations (medians ranging between 0.33 and 0.42),
562 including optimized aggregations. The median accuracy of the best estimate is better for all
563 aggregations than the median ranked individual for this measure. However, the optimized
564 aggregations have some groups which perform worse than the median individual. This may

565 not be surprising because (as discussed) the Classical Model was not designed to optimize
566 point estimates.

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

572 To check this, we replaced the answers for these two questions with 0.000011 and re-
573 calculated the calibration and statistical accuracy of judgments of each of the groups (Figure
574 8, see also Appendix S1: Section 8). The adjustment improved the statistical accuracy of
575 many groups across all aggregations. All but one aggregation (quantile aggregation, QuA)
576 had a median statistical accuracy above 0.53 (Figure 8a). Only one group was considered
577 statistically inaccurate when their judgments were combined via quantile aggregation (QuA).

578 Quantile aggregation (QuA) was overconfident, even when assessed according to calibration
579 of interval judgments but many groups were less so than prior to accounting for the two
580 questions with zeros (Figure 8b). Group judgments for quantile aggregation achieved good
581 but not perfect median calibration of 0.76, although no group reached perfect calibration
582 when quantile aggregation (QuA) was used. In contrast, each of the linear pooled
583 distributions except for the itemized optimized weights achieved a median group calibration
584 of 0.90 (i.e. perfect calibration). The data adjustments improved calibration and statistical
585 accuracy, but they did not substantially alter the information or informativeness scores which
586 meant that quantile aggregation (QuA) was still substantially more informative than the equal
587 weights (EW) (a median informativeness score of 0.24 compared to 0.42).

588 **4 Discussion**

589 Performance weights have been proposed to improve expert judgments in ecology. However,
590 there have been few applications and little discussion of their strengths and weaknesses in the
591 ecological literature. Here, we outlined the key rationales and theories of performance
592 weights, then described one of the most well-known methods, the Classical Model (Cooke
593 1991), and examined how it might be applied to improve judgments derived from the IDEA
594 protocol with four-step question format (Hemming et al. 2018a, Hemming et al. 2018b).

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

598 The four-step question format needs to first be converted into quantiles of a continuous
599 probability distribution. It may be better to remove these assumptions by eliciting these
600 quantiles directly. However, the four-step question format is often used because it helps to
601 overcome overconfidence relative to eliciting fixed intervals (Speirs-Bridge et al. 2010), and
602 because experts who are unfamiliar with the language of statistical distributions are
603 comfortable in providing quantitative judgments of uncertainty (a problem not only
604 encountered in ecological domains (Walls and Quigley 2001, Hirsch et al. 2004)). These
605 trade-offs need to be considered when deciding how best to elicit estimates. If the four-step
606 question format is to be used with the Classical Model, then we suggest that the assumptions
607 about how the estimates will be interpreted are communicated to experts in introductory
608 material and through the feedback and discussion stages of the IDEA protocol.

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

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

621 Insights from our results emphasize that the Classical Model was designed to assess
622 probability distributions rather than point estimates or interval judgments (as some
623 interpretations suggest). Specifically, ‘statistical accuracy’ measures the degree to which an
624 expert’s multinomial distribution matches a theoretically optimal multinomial distribution,
625 and ‘information’ measures the departure from a uniform or log-uniform background
626 measure. As such the Classical Model is not focused primarily on avoiding surprises outside

627 of the 90% confidence intervals, or the precision of the intervals (as assessed in the IDEA
628 protocol) and may lead to counterintuitive outcomes in settings where this is a primary
629 concern.

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

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

643 While calibration and informativeness of interval judgments may be of interest they have not
644 yet been combined into a proper scoring rule (although telling experts they will be scored on
645 both should minimise gaming behaviour). Our results demonstrate that the Classical Model
646 does not by itself provide this information, which may be disappointing to those who seek to
647 apply the Classical Model to optimize or assess such judgments. However, if this information
648 was of interest the performance measures of the IDEA protocol may be used to provide this
649 information. Agreement as to which performance measures will be used should be made prior
650 to application.

651 Equal weighted aggregations are often used in ecology when combining expert judgments.
652 However, there are numerous methods by which an equal weighted aggregation can be
653 derived, and not all will perform equally well or have been validated. We contrasted two
654 forms of equal weighted aggregation, quantile aggregation (QuA, used in Hemming et al.
655 2018), and equal weighting via linear pooling of distributions (EW, used by the Classical
656 Model). We found that both forms of equal weighted aggregation were better than the median
657 ranked individual for each measure of statistical accuracy, calibration, and accuracy.
658 Furthermore, as was demonstrated in Hemming et al. (2018b), while some individuals could

659 outperform the group aggregation they could not be predicted by standard metrics of
660 expertise (years of experience, peer-recommendation, or self-rating). This suggests that
661 taking the equal weighted aggregation is a more robust method than trying to select a single
662 expert with good judgment based on their credentials and status.

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

668 We found that the degree of overconfidence was reduced when we accounted for questions
669 with zeros, and the way in which the Classical Model accounts for realisations which equal a
670 participant's estimates (i.e. if the realisation coincides with the lower bound it will be
671 considered as falling outside of the expert's 90% credible intervals). As these adjustments are
672 not made when the four-step question format is used in the IDEA protocol, the degree of
673 overconfidence from quantile aggregation may not typically be as severe for many
674 applications of the IDEA protocol. Nonetheless, we would suggest these findings warrant
675 further investigation on more case studies with the four-step question-format.

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

683 In our study, we demonstrated a modest improvement by performance weighted aggregation.
684 However, we note that there is no guarantee that performance weighted aggregation will lead
685 to improvements in all cases. However, a clear advantage of the Classical Model, and other
686 methods which utilise calibration questions is that they provide empirical evidence for the
687 legitimacy of final aggregations (often lacking in studies that use expert judgment). This is
688 especially important because decisions regarding who should be included in an elicitation and
689 how to aggregate these judgments may exclude potentially knowledgeable individuals, and
690 often lack validation.

691 Whether or not the decision context justifies (or can afford) the additional time and expense
692 ultimately depends on the context of the case study, the decision-maker and the value of
693 additional information. Wittmann et al. (2015) and Rothlisberger et al. (2009) justify their
694 application based on the immense value of fisheries to the Great Lakes and the possibility of
695 litigation following mismanagement. This suggests that there are contexts in ecology in
696 which this additional time and expense can be justified. If resources are not available to
697 deploy calibration questions and performance weighted aggregation, then our study shows
698 that an equal weighted aggregation (i.e. quantile aggregation or linear pooling of
699 distributions) provides an effective means to improve judgments relative to selecting a single
700 seemingly well-credentialed expert.

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

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

712 Calibration questions should be related to target variables, for which the answer is known or
713 will become known (Cooke and Goossens 2000). However, ascertaining whether or not a
714 question is *relevant* in many domains may be difficult because domains are often ill-defined,
715 making the selection of *relevant* questions a subjective decision (Colyvan and Ginzburg
716 2003). If datasets are difficult to obtain, then the analyst may need to rely on past questions
717 for which the data are available, or questions which are less relevant to the questions of
718 interest. It would be useful to understand at what point calibration questions become so
719 distantly related to target questions that in-sample validation is not a good predictor of
720 performance.

721 We used *Excalibur* to generate aggregations and score experts, however, the program was
722 challenging to use. The analysis was time consuming and it was difficult to provide a

723 reproducible workflow for our analysis. The methods of aggregation and the scoring rules
724 should be simple enough to re-code in *R* and other freely available software (we note that
725 recently they have been re-coded in *MATLAB* (Leontaris and Morales-Nápoles 2018)). A
726 revision of *Excalibur* could help to increase adoption of the method.

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

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

741 **5 Conclusions**

742 Performance weighted aggregations with calibration questions has been proposed as a means
743 to improve expert judgments in ecology, however, applications have been scarce. We
744 explored how the Classical Model could be applied to the IDEA protocol with four-step
745 question format.

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

753 We demonstrated that equal weighted aggregations can achieve relatively well-calibrated
754 aggregated judgments. However, linear pooling of distributions may produce better calibrated
755 but less informative distributions than quantile aggregation as found by Bamber et al. (2016)
756 and Colson and Cooke (2017). We found that performance weighted aggregations can
757 outperform equal weighted aggregations, in our case by providing more informative
758 judgments, however, we emphasize that there is no guarantee they will do so in every case.
759 The main reason that the candidate alternatives for aggregation should be explored is to
760 ensure the final representation of uncertainty is the best possible (whether that be via equal
761 weights or performance weights).

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

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

771 **6 Acknowledgements**

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

778

779 **7 Literature citations**

780 Adams-Hosking, C., M. F. McBride, G. Baxter, M. Burgman, D. de Villiers, R.
781 Kavanagh, I. Lawler, D. Lunney, A. Melzer, P. Menkhorst, R. Molsher, B. D.

782 Moore, D. Phalen, J. R. Rhodes, C. Todd, D. Whisson, and C. A. McAlpine.
783 2016. Use of expert knowledge to elicit population trends for the koala
784 (*Phascolarctos cinereus*). *Diversity and Distributions* **22**:249-262.

785 Addison, P. F. E., K. de Bie, and L. Rumpff. 2015. Setting conservation management
786 thresholds using a novel participatory modeling approach. *Conservation Biology*
787 **29**:1411-1422.

788 Armstrong, J. S. 2001. Combining forecasts. Pages 417-439 in J. S. Armstrong, editor.
789 *Principles of forecasting: A handbook for researchers and practitioners*. Springer
790 US, Boston, United States of America.

791 Aspinall, W. P., and R. M. Cooke. 2013. Quantifying scientific uncertainty from expert
792 judgement elicitation. Pages 64-99 in J. Rougier, S. Sparks, and L. Hill, editors.
793 *Risk and Uncertainty Assessment for Natural Hazards*. Cambridge University
794 Press, Cambridge, United Kingdom.

795 Bamber, J., W. Aspinall, and R. Cooke. 2016. A commentary on “how to interpret expert
796 judgment assessments of twenty-first century sea-level rise” by Hylke de Vries
797 and Roderik SW van de Wal. *Climatic Change* **137**:321-328.

798 Barons, M. J., A. M. Hanea, S. K. Wright, K. C. Baldock, L. Wilfert, D. Chandler, S.
799 Datta, J. Fannon, C. Hartfield, and A. Lucas. 2018. Assessment of the response of
800 pollinator abundance to environmental pressures using structured expert
801 elicitation. *Journal of Apicultural Research*:1-12.

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

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

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

808 Burgman, M., A. Carr, L. Godden, R. Gregory, M. McBride, L. Flander, and L. Maguire.
809 2011a. Redefining expertise and improving ecological judgment. *Conservation*
810 *Letters* **4**:81-87.

- 811 Burgman, M. A. 2004. Expert frailties in conservation risk assessment and listing
812 decisions. Pages 20-29 *in* P. Hutchings, D. Lunney, and C. Dickman, editors.
813 Threatened species legislation: is it just an Act? Royal Zoological Society,
814 Mosman, NSW, Australia.
- 815 Burgman, M. A. 2015. Trusting Judgements: How to get the best out of experts.
816 Cambridge University Press, Cambridge, United Kingdom.
- 817 Burgman, M. A., M. McBride, R. Ashton, A. Speirs-Bridge, L. Flander, B. Wintle, F.
818 Fidler, L. Rumpff, and C. Twardy. 2011b. Expert status and performance. PLoS
819 One **6**:1-7.
- 820 Carwardine, J., T. G. Martin, J. Firn, R. P. Reyes, S. Nicol, A. Reeson, H. S. Grantham,
821 D. Stratford, L. Kehoe, and I. Chadès. 2019. Priority Threat Management for
822 biodiversity conservation: A handbook. *Journal of Applied Ecology* **56**:481-490.
- 823 Clemen, R. T. 2008. Comment on Cooke's classical method. *Reliability Engineering &*
824 *System Safety* **93**:760-765.
- 825 Clemen, R. T., and R. L. Winkler. 1999. Combining probability distributions from
826 experts in risk analysis. *Risk Analysis* **19**:187-203.
- 827 Colson, A. R., and R. M. Cooke. 2017. Cross validation for the classical model of
828 structured expert judgment. *Reliability Engineering & System Safety* **163**:109-
829 120.
- 830 Colyvan, M., and L. R. Ginzburg. 2003. Laws of nature and laws of ecology. *Oikos*
831 **101**:649-653.
- 832 Cooke, R. 2018a. Macro converting XL file to EXCALIBUR dtt file. *in* R. Cooke, editor.
833 <http://rogermcooke.net/>
- 834 Cooke, R., and L. Goossens. 2000. Procedures guide for structural expert judgement in
835 accident consequence modelling. *Radiation Protection Dosimetry* **90**:303-309.
- 836 Cooke, R. M. 1991. Experts in uncertainty: Opinion and subjective probability in
837 science. Oxford University Press, New York.

- 838 Cooke, R. M. 2018b. Validation in the Classical Model. Pages 37-59 *in* L. C. Dias, A.
839 Morton, and J. Quigley, editors. Elicitation: The science and art of structuring
840 judgement. Springer International Publishing, Cham, Switzerland.
- 841 Cooke, R. M., and L. L. Goossens. 2008. TU Delft expert judgment data base. Reliability
842 Engineering & System Safety **93**:657-674.
- 843 Drescher, M., A. Perera, C. Johnson, L. Buse, C. Drew, and M. Burgman. 2013. Toward
844 rigorous use of expert knowledge in ecological research. Ecosphere **4**:1-26.
- 845 Einhorn, H. J., R. M. Hogarth, and E. Klempner. 1977. Quality of group judgment.
846 Psychological Bulletin **84**:158.
- 847 Estévez, R. A., F. O. Mardones, F. Álamos, G. Arriagada, J. Carey, C. Correa, J.
848 Escobar-Dodero, Á. Gaete, A. Gallardo, and R. Ibarra. 2019. Eliciting expert
849 judgements to estimate risk and protective factors for Piscirickettsiosis in Chilean
850 salmon farming. Aquaculture.
- 851 Flandoli, F., E. Giorgi, W. P. Aspinall, and A. Neri. 2011. Comparison of a new expert
852 elicitation model with the Classical Model, equal weights and single experts,
853 using a cross-validation technique. Reliability Engineering & System Safety
854 **96**:1292-1310.
- 855 GBRMPA. 2014. Great Barrier Reef Outlook Report. Townsville, Australia.
- 856 Gneiting, T., and A. E. Raftery. 2007. Strictly proper scoring rules, prediction, and
857 estimation. Journal of the American Statistical Association **102**:359-378.
- 858 Goossens, L. H. J., R. M. Cooke, A. R. Hale, and L. Rodić-Wiersma. 2008. Fifteen years
859 of expert judgement at TUDelft. Safety Science **46**:234-244.
- 860 Gregory, R., L. Failing, M. Harstone, G. Long, T. McDaniels, and D. Ohlson. 2012.
861 Structured Decision Making: a practical guide to environmental management
862 choices, Chichester, West Sussex.
- 863 Gregory, R., and R. L. Keeney. 2017. A Practical Approach to Address Uncertainty in
864 Stakeholder Deliberations. Risk Analysis **37**:487-501.

- 865 Hanea, A., M. McBride, M. Burgman, B. Wintle, F. Fidler, L. Flander, B. Manning, and
866 S. Mascaro 2016. InvestigateDiscussEstimateAggregate for structured expert judgement.
867 International journal of forecasting **33**:267-269.
- 868 Hemming, V. 2019. Code: Weighting and Aggregating Expert Ecological Judgements.
869 The Open Science Framework. DOI 10.17605/OSF.IO/FXQVK.
870 <http://osf.io/fxqvk>
- 871 Hemming, V., N. Armstrong, M. A. Burgman, and A. M. Hanea. 2019a. Improving
872 expert forecasts in reliability: Application and evidence for structured elicitation
873 protocols. Quality and Reliability Engineering International **n/a**.
- 874 Hemming, V., M. A. Burgman, A. M. Hanea, M. F. McBride, and B. C. Wintle. 2018a.
875 A practical guide to structured expert elicitation using the IDEA protocol.
876 Methods in Ecology and Evolution **9**:169-181.
- 877 Hemming, V., S. Lane, and A. Hanea. 2019b. Classical Model Calculator. The Open
878 Science Framework. DOI 10.17605/OSF.IO/BGYZU. <http://osf.io/bgyzu>
- 879 Hemming, V., T. V. Walshe, A. M. Hanea, F. Fidler, and M. A. Burgman. 2018b.
880 Eliciting improved quantitative judgements using the IDEA protocol: A case
881 study in natural resource management. PLoS One **13**:e0198468.
- 882 Hirsch, K. G., J. J. Podur, R. F. Janser, R. S. McAlpine, and D. L. Martell. 2004.
883 Productivity of Ontario initial-attack fire crews: results of an expert-judgement
884 elicitation study. Canadian Journal of Forest Research **34**:705-715.
- 885 Hogarth, R. M. 1978. A note on aggregating opinions. Organizational Behavior and
886 Human Performance **21**:40-46.
- 887 Hora, S. C. 2004. Probability judgments for continuous quantities: Linear combinations
888 and calibration. Management Science **50**:597-604.
- 889 Hudson, E. G., V. J. Brookes, and M. P. Ward. 2017. Assessing the risk of a canine
890 rabies incursion in Northern Australia. Frontiers in Veterinary Science **4**:141.
- 891 Krzywinski, M., and N. Altman. 2014. Points of Significance: Visualizing samples with
892 box plots. Nat Meth **11**:119-120.

- 893 Kuhnert, P. M., T. G. Martin, and S. P. Griffiths. 2010. A guide to eliciting and using
894 expert knowledge in Bayesian ecological models. *Ecology Letters* **13**:900-914.
- 895 Larrick, R. P., and J. B. Soll. 2006. Intuitions about combining opinions:
896 Misappreciation of the averaging principle. *Management Science* **52**:111-127.
- 897 Leontaris, G., and O. Morales-Nápoles. 2018. ANDURIL — A MATLAB toolbox for
898 ANalysis and Decisions with UnceRtaInty: Learning from expert judgments.
899 *SoftwareX* **7**:313-317.
- 900 Lichtendahl Jr, K. C., Y. Grushka-Cockayne, and R. L. Winkler. 2013. Is it better to
901 average probabilities or quantiles? *Management Science* **59**:1594-1611.
- 902 Lichtenstein, S., and B. Fischhoff. 1977. Do those who know more also know more
903 about how much they know? *Organizational Behavior and Human Performance*
904 **20**:159-183.
- 905 Lightwist. 2013. Excalibur. <http://www.lighttwist.net/wp/excalibur>
- 906 Lin, S.-W., and V. M. Bier. 2008. A study of expert overconfidence. *Reliability*
907 *Engineering & System Safety* **93**:711-721.
- 908 Lorenz, J., H. Rauhut, F. Schweitzer, and D. Helbing. 2011. How social influence can
909 undermine the wisdom of crowd effect. *Proceedings of the National Academy of*
910 *Sciences* **108**:9020-9025.
- 911 Low Choy, S., R. O'Leary, and K. Mengersen. 2009. Elicitation by design in ecology:
912 using expert opinion to inform priors for Bayesian statistical models. *Ecology*
913 **90**:265-277.
- 914 MacDonald, J. A., M. J. Small, and M. Morgan. 2008. Explosion probability of
915 unexploded ordnance: expert beliefs. *Risk Analysis* **28**:825-841.
- 916 Martin, T. G., M. A. Burgman, F. Fidler, P. M. Kuhnert, S. Low-Choy, M. McBride, and
917 K. Mengersen. 2012. Eliciting expert knowledge in conservation science.
918 *Conservation Biology* **26**:29-38.
- 919 Martin, T. G., B. A. Wintle, J. R. Rhodes, P. M. Kuhnert, S. A. Field, S. J. Low-Choy, A.
920 J. Tyre, and H. P. Possingham. 2005. Zero tolerance ecology: improving

- 921 ecological inference by modelling the source of zero observations. *Ecology*
922 *Letters* **8**:1235-1246.
- 923 McBride, M. F., F. Fidler, and M. A. Burgman. 2012a. Evaluating the accuracy and
924 calibration of expert predictions under uncertainty: predicting the outcomes of
925 ecological research. *Diversity and Distributions* **18**:782-794.
- 926 McBride, M. F., S. T. Garnett, J. K. Szabo, A. H. Burbidge, S. H. Butchart, L. Christidis,
927 G. Dutson, H. A. Ford, R. H. Loyn, and D. M. Watson. 2012b. Structured
928 elicitation of expert judgments for threatened species assessment: a case study on
929 a continental scale using email. *Methods in Ecology and Evolution* **3**:906-920.
- 930 Meek, M. H., C. Wells, K. M. Tomalty, J. Ashander, E. M. Cole, D. A. Gille, B. J.
931 Putman, J. P. Rose, M. S. Savoca, and L. Yamane. 2015. Fear of failure in
932 conservation: the problem and potential solutions to aid conservation of extremely
933 small populations. *Biological Conservation* **184**:209-217.
- 934 Mellers, B., E. Stone, T. Murray, A. Minster, N. Rohrbaugh, M. Bishop, E. Chen, J.
935 Baker, Y. Hou, and M. Horowitz. 2015. Identifying and cultivating
936 superforecasters as a method of improving probabilistic predictions. *Perspectives*
937 *on Psychological Science* **10**:267-281.
- 938 Metcalf, S. J., and K. J. Wallace. 2013. Ranking biodiversity risk factors using expert
939 groups – Treating linguistic uncertainty and documenting epistemic uncertainty.
940 *Biological Conservation* **162**:1-8.
- 941 Morgan, M. G., and M. Henrion. 1990. *Uncertainty: A guide to dealing with uncertainty*
942 *in quantitative risk and policy analysis* Cambridge University Press. New York,
943 NY, United States of America.
- 944 Quigley, J., A. Colson, W. Aspinall, and R. M. Cooke. 2018. Elicitation in the Classical
945 Model. Pages 15-36 in L. C. Dias, A. Morton, and J. Quigley, editors. *Elicitation:*
946 *The science and art of structuring judgement*. Springer International Publishing,
947 Cham, Switzerland.
- 948 Rothlisberger, J. D., D. M. Lodge, R. M. Cooke, and D. C. Finnoff. 2009. Future
949 declines of the binational Laurentian Great Lakes fisheries: the importance of

- 950 environmental and cultural change. *Frontiers in Ecology and the Environment*
951 **8**:239-244.
- 952 Satopää, V. A., J. Baron, D. P. Foster, B. A. Mellers, P. E. Tetlock, and L. H. Ungar.
953 2014. Combining multiple probability predictions using a simple logit model.
954 *International journal of forecasting* **30**:344-356.
- 955 Shanteau, J., D. J. Weiss, R. P. Thomas, and J. C. Pounds. 2002. Performance-based
956 assessment of expertise: How to decide if someone is an expert or not. *European*
957 *Journal of Operational Research* **136**:253-263.
- 958 Speirs-Bridge, A., F. Fidler, M. McBride, L. Flander, G. Cumming, and M. Burgman.
959 2010. Reducing overconfidence in the interval judgments of experts. *Risk*
960 *Analysis* **30**:512-523.
- 961 Surowiecki, J. 2004. *The wisdom of crowds: Why the many are smarter than the few and*
962 *how collective wisdom shapes business, economies, societies, and nations.* Little,
963 Brown, London, United Kingdom.
- 964 Tetlock, P. 2005. *Expert political judgment: How good is it? How can we know?*
965 Princeton University Press, Princeton, New Jersey, USA.
- 966 Tetlock, P., and D. Gardner. 2015. *Superforecasting: The art and science of prediction.*
967 Random House, New York.
- 968 Walls, L., and J. Quigley. 2001. Building prior distributions to support Bayesian
969 reliability growth modelling using expert judgement. *Reliability Engineering &*
970 *System Safety* **74**:117-128.
- 971 Wallsten, T. S., D. V. Budescu, A. Rapoport, R. Zwick, and B. Forsyth. 1986. Measuring
972 the vague meanings of probability terms. *Journal of Experimental Psychology:*
973 *General* **115**:348.
- 974 Ward, T. 2014. *The rapid assessment workshop to elicit expert consensus to inform the*
975 *development of the Great Barrier Reef Outlook Report 2014.* Townsville.
- 976 Weiss, D. J., and J. Shanteau. 2004. The vice of consensus and the virtue of consistency.
977 *Psychological investigations of competent decision making*:226-240.

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

980 Wittmann, M. E., R. M. Cooke, J. D. Rothlisberger, E. S. Rutherford, H. Zhang, D. M.
981 Mason, and D. M. Lodge. 2015. Use of structured expert judgment to forecast
982 invasions by bighead and silver carp in Lake Erie. Conservation Biology **29**:187-
983 197.

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

986

987 **8 Data availability statement**

988 All data and code for the analyses presented in this paper are available on the Open Science
989 Framework (Hemming 2019): <https://doi.org/10.17605/OSF.IO/FXQVK>

990

991

992

993

994

995 **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.

996

997

998 Figure 1 Key steps of the IDEA protocol (figure from Hemming et al. (2018b)). The
999 four-step question format (Speirs-Bridge et al. 2010) (depicted in Step 2) is commonly used
1000 to derive a best estimate and credible interval in Round 1 and Round 2.

1001 Figure 2 Judgments provided by two hypothetical experts over 10 questions. The blue
1002 lines represent their 90% credible intervals, the blue dots their 'best estimate' or their
1003 'median'. The crosses represent where the realisation fell in relation to their estimates. To
1004 calculate statistical accuracy (SA) according to the Classical Model, the proportion of
1005 questions answered where realisations fell, 1) below their lowest interval (i.e. 5th quantile), 2)
1006 between their lowest estimate and their best estimate / median, 3) between their best estimate
1007 / median and their upper estimate / 95th quantile, and 4) above their upper / 95th quantile is

1008 calculated and compared to a theoretically optimal distribution $p=(0.05, 0.45, 0.45, 0.05)$. CA
1009 refers to calibration as calculated according to the IDEA protocol, which is defined as the
1010 proportion of credible intervals capturing the realisation.

1011 Figure 3 The statistical accuracy and information of $n = 58$ participants. A trade-off
1012 exists between the two measures used by the Classical Model. Those who are statistically
1013 accurate (above 0.05, red horizontal line) often have a lower information score than the
1014 median score for individuals (grey vertical line). The blue dashed line shows the highest
1015 statistical accuracy score possible for 13 questions (0.93), and the black line shows the
1016 highest score obtained by individuals in the elicitation (0.53).

1017 Figure 4 Statistical accuracy of the Classical Model (CM) compared to IDEA
1018 calibration for $n = 58$ participants. The graph shows that participants with perfect calibration
1019 when assessed by the IDEA protocol, can have poor statistical accuracy for the Classical
1020 Model. On the righthand side, we show where the realisations fell in each of the expert's
1021 multinomial distributions (used to calculate statistical accuracy), and contrast this with how
1022 many realisations fell within the participant's 90% credible intervals (calibration). Bold
1023 numbers indicate the highest scores possible for statistical accuracy and calibration.

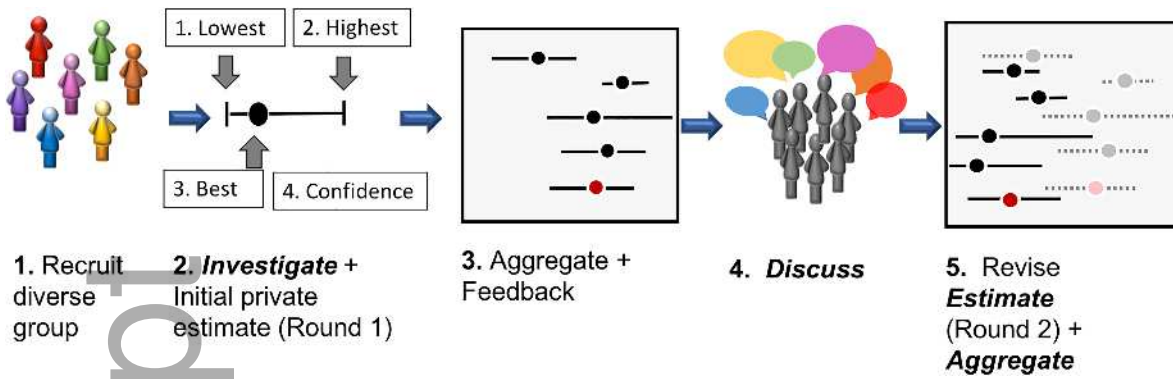
1024 Figure 5 The spearman correlation between information calculated for the Classical
1025 Model, and informativeness calculated for the IDEA protocol for $n = 58$ participants. The
1026 shaded area represents a 95% confidence interval.

1027 Figure 6 CM Scores derived for each aggregation.

1028 Figure 7 Component performance measures of the Classical Model (CM) and IDEA
1029 protocol for $n = 8$ groups under six alternative procedures for aggregation. a) Statistical
1030 accuracy, the red-dashed line represents the 0.05 threshold for statistically inaccurate scores,
1031 (Classical Model), the blue dashed line represents a perfect statistical accuracy score for 13
1032 questions, and the black dashed line represents the highest score obtained by any individual,
1033 b) information score (Classical Model), the red line represents the median information of an
1034 individual c) calibration, (IDEA) the red line represents perfect calibration (0.90), d)
1035 informativeness (IDEA), the red line represents the informativeness of the median individual,
1036 e) accuracy (IDEA) of the best estimate, the red line represents the accuracy of the median
1037 individual.

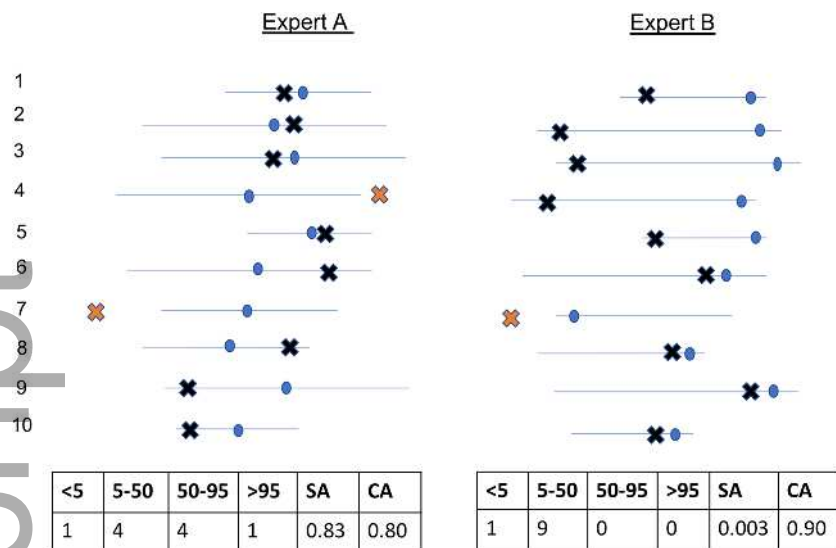
1038 Figure 8 The scores of aggregations under the Classical Model and the IDEA protocol
1039 when adjustments are made to correct for questions for which the realised truth had been
1040 zero. a) Statistical accuracy, the red-dashed line represents the 0.05 threshold for statistically
1041 inaccurate scores, (Classical Model), the blue dashed line represents a perfect statistical
1042 accuracy score for 13 questions, and the black dashed line represents the highest score
1043 obtained by any individual prior to the adjustment; b) calibration, (IDEA) the red line
1044 represents perfect calibration (0.90).

Author Manuscript

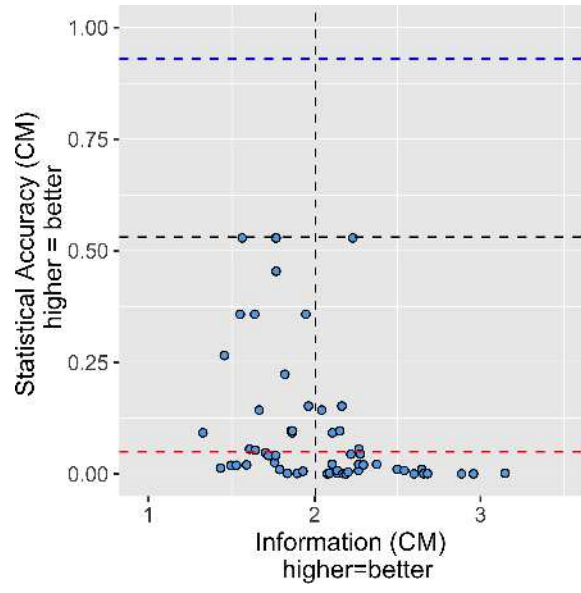


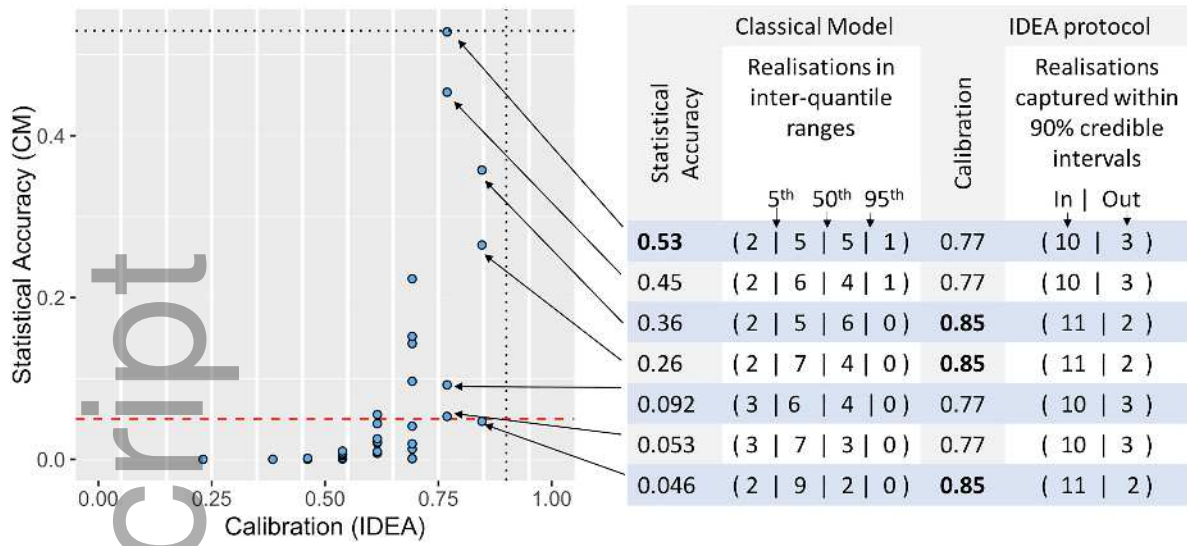
eap_2075_f1.tif

Author Manuscript



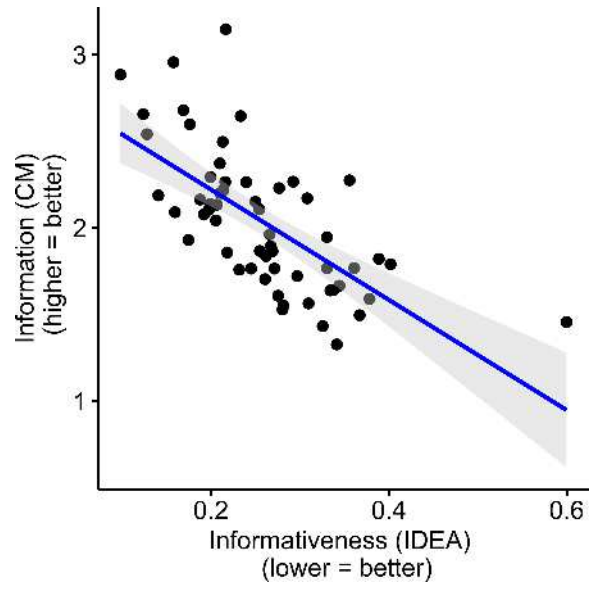
eap_2075_f2.tif



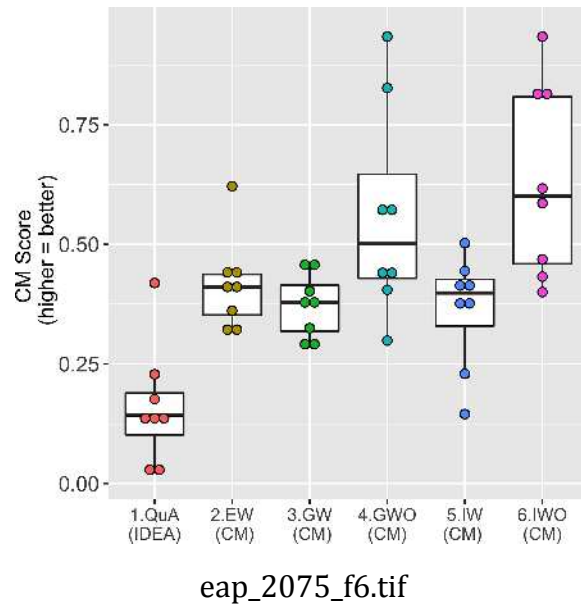


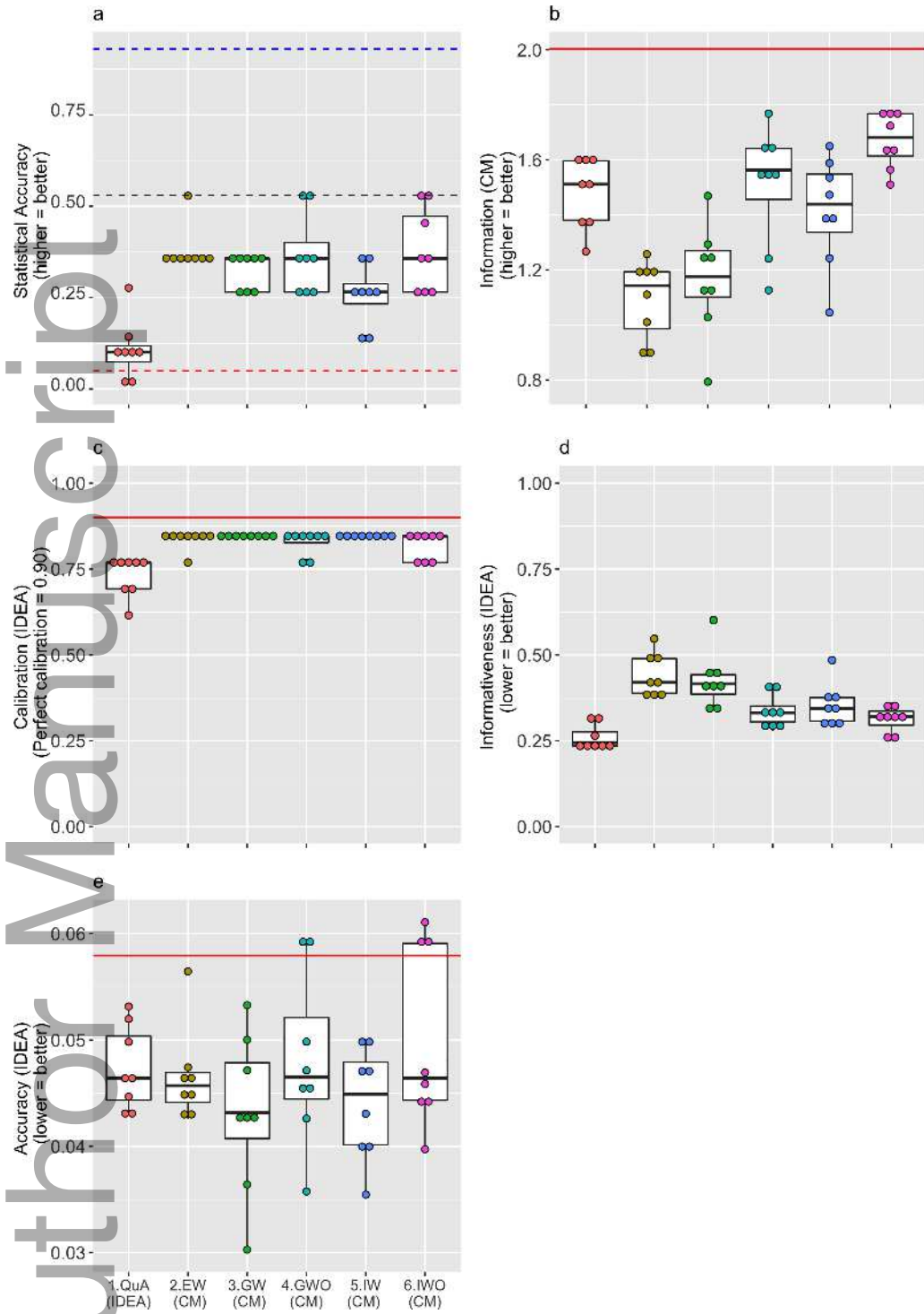
eap_2075_f4.tif

Author Manuscript

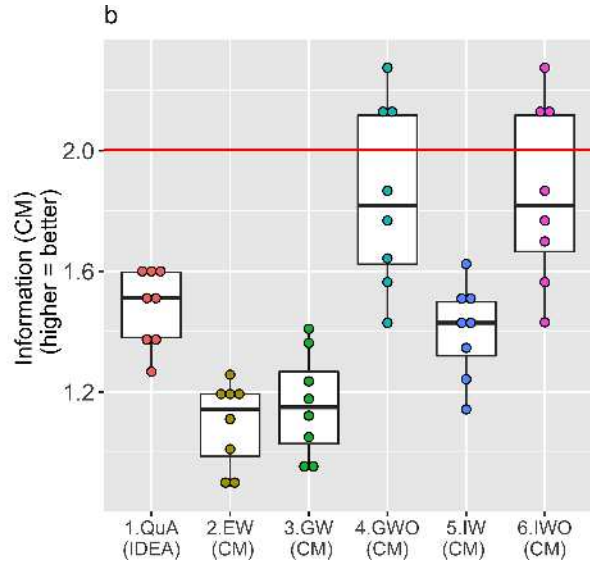
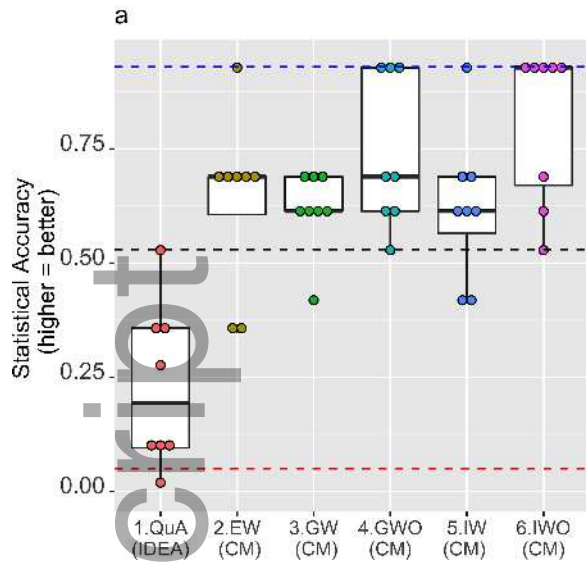


eap_2075_f5.tif





eap_2075_f7.tif



eap_2075_f8.tif