

1 **TITLE:** Comparing traditional and Bayesian approaches to ecological meta-analysis

2

3 **RUNNING TITLE:** Comparing methods for ecological meta-analysis

4

5 **AUTHORS:** P. Pappalardo^{1*}, K. Ogle^{2,3}, E.A. Hamman¹, J.R. Bence⁵, B.A. Hungate³, and

6 C.W. Osenberg¹

7

8 **ORCID IDs:**

9 Paula Pappalardo: 0000-0003-0853-7681

10 Kiona Ogle: 0000-0002-0652-8397

11 Elizabeth A. Hamman: 0000-0002-3494-6641

12 James R. Bence: 0000-0002-2534-688X

13 Bruce A. Hungate: 0000-0002-7337-1887

14 Craig W. Osenberg: 0000-0003-1918-7904

15

16 **AFFILIATIONS:**

17 ¹Odum School of Ecology, University of Georgia, Athens, GA 30602, USA

18 ²School of Informatics, Computing, and Cyber Systems, Northern Arizona University,

19 Flagstaff, AZ 86011, USA

20 ³Center for Ecosystem Science and Society and Department of Biological Sciences, Northern

21 Arizona University, Flagstaff, AZ 86011, USA

22 ⁵Quantitative Fisheries Center, Department of Fisheries and Wildlife, Michigan State

23 University, East Lansing, MI 48824, USA

24 * Corresponding author contact information: Department of Invertebrate Zoology,
25 Smithsonian National Museum of Natural History, Washington, DC 20560, USA, email:
26 paulapappalardo@gmail.com, phone: 706-308-2979
27
28

29 **PAPER TYPE:** Research article

30

31 **KEY WORDS:** bias, confidence interval, coverage, credible interval, effect size, global
32 climate change, log response ratio, sample size

33

34 **ABSTRACT**

- 35 1. Despite the wide application of meta-analysis in ecology, some of the traditional
36 methods used for meta-analysis may not perform well given the type of data
37 characteristic of ecological meta-analyses.
- 38 2. We reviewed published meta-analyses on the ecological impacts of global climate
39 change, evaluating the number of replicates used in the primary studies (n_i) and the
40 number of studies or records (k) that were aggregated to calculate a mean effect size.
41 We used the results of the review in a simulation experiment to assess the
42 performance of conventional frequentist and Bayesian meta-analysis methods for
43 estimating a mean effect size and its uncertainty interval.
- 44 3. Our literature review showed that n_i and k were highly variable, distributions were
45 right-skewed, and were generally small (median $n_i = 5$, median $k = 44$). Our
46 simulations show that the choice of method for calculating uncertainty intervals was
47 critical for obtaining appropriate coverage (close to the nominal value of 0.95). When
48 k was low (< 40), 95% coverage was achieved by a confidence interval based on the
49 t -distribution that uses an adjusted standard error (the Hartung-Knapp-Sidik-
50 Jonkman, HKSJ), or by a Bayesian credible interval, whereas bootstrap or z -
51 distribution confidence intervals had lower coverage. Despite the importance of the
52 method to calculate the uncertainty interval, 39% of the meta-analyses reviewed did
53 not report the method used, and of the 61% that did, 94% used a potentially
54 problematic method, which may be a consequence of software defaults.
- 55 4. In general, for a simple random-effects meta-analysis, the performance of the best
56 frequentist and Bayesian methods were similar for the same combinations of factors
57 (k and mean replication), though the Bayesian approach had higher than nominal

58 (>95%) coverage for the mean effect when k was very low ($k < 15$). Our literature
59 review suggests that many meta-analyses that used z-distribution or bootstrapping
60 confidence intervals may have over-estimated the statistical significance of their
61 results when the number of studies was low; more appropriate methods need to be
62 adopted in ecological meta-analyses.

63

64 **RESUMEN**

- 65 1. A pesar del uso generalizado del meta-análisis en el área de Ecología, algunos de los
66 métodos de análisis tradicionalmente utilizados pueden dar resultados no ideales dado
67 el tipo de datos que los caracteriza.
- 68 2. En este trabajo se realizó una revisión de los meta-análisis publicados sobre los
69 impactos ecológicos del cambio climático global, evaluando el número de réplicas
70 utilizadas en las publicaciones originales (n_i) y el número de estudios o registros (k)
71 que fueron agrupados para calcular un tamaño de efecto promedio. Se utilizaron los
72 resultados de la revisión en un experimento de simulación para evaluar el desempeño
73 de métodos frecuentistas convencionales y métodos Bayesianos para estimar un
74 tamaño de efecto promedio y su intervalo de incertidumbre.
- 75 3. La revisión de la literatura demostró que n_i y k fueron muy variables, con
76 distribuciones sesgadas, y con valores en general bajos (mediana $n_i = 5$, mediana
77 $k = 44$). Nuestras simulaciones muestran que la elección del método para calcular un
78 intervalo de incertidumbre fue crítica para obtener una cobertura apropiada (alrededor
79 del valor nominal de 0.95). Cuando k fue bajo (< 40), obtuvimos una cobertura de
80 95% utilizando un intervalo de confianza basado en la distribución t de student que
81 usa un ajuste por el error estándar (llamada Hartung-Knapp-Sidik-Jonkman, HKSJ),

82 o utilizando un intervalo de credibilidad Bayesiano, mientras que los intervalos de
83 remuestreo o con una distribución Normal tuvieron cobertura baja. A pesar de la
84 importancia del método utilizado para calcular el intervalo de incertidumbre, 39% de
85 los meta-análisis revisados no reportaron el método utilizado y, de los 61% que si lo
86 reportaron, 94% usaron uno de los métodos potencialmente problemáticos, lo que
87 puede ser una consecuencia de la configuración por defecto de los programas
88 informáticos utilizados para meta-análisis.

89 4. En general, para un meta-análisis simple con efectos aleatorios, el desempeño del
90 mejor método frecuentista y el método Bayesiano fueron similares para las mismas
91 combinaciones de factores (k y número de réplicas promedio), aunque el método
92 Bayesiano tuvo cobertura mayor de la nominal (>95%) para el efecto promedio
93 cuando k fue muy bajo ($k < 15$). Nuestra revisión sugiere que muchos de los meta-
94 análisis que utilizaron una distribución Normal o intervalos de remuestreo pueden
95 haber sobreestimado la significancia estadística de sus resultados cuando el número
96 de estudios fue bajo. Otros métodos más apropiados deberían ser usados para meta-
97 análisis en Ecología.

98 **INTRODUCTION**

99 Meta-analysis uses statistical techniques to quantitatively summarize information from
100 different studies and is highly influential in the contemporary practice of science. To conduct
101 a meta-analysis an investigator gathers summary statistics from each study to calculate an
102 effect size, with the goal of computing an overall effect size (and its uncertainty) and
103 exploring the factors contributing to variation in effect sizes (Nakagawa, Noble, Senior, &
104 Lagisz, 2017). The use of meta-analysis in ecology has been growing rapidly since the 1990s,
105 and has proven particularly useful in discerning general patterns by comparing information
106 from different species, study sites, and systems (Cadotte, Mehrkens, & Menge, 2012). Advice
107 on best methodological practices for meta-analysis is widespread in disciplines with a longer
108 history of meta-analytic research (e.g. medical sciences) but is lagging behind in ecology
109 (Gates, 2002). This can be problematic because ecological meta-analyses have specific
110 challenges not necessarily typically in other disciplines.

111 One pervasive characteristic of ecological meta-analyses is the high heterogeneity
112 (i.e., large among-study variation in effect sizes). Senior et al. (2016) analyzed 86 meta-
113 analyses in ecology and evolution and found that the among-study variation averaged 92%
114 of the total variance. In contrast, a review of 509 meta-analyses in medicine found that there
115 was no detectable among-study variation in 50% of the studies (Higgins, Thompson, &
116 Spiegelhalter, 2009). Ecological studies also differ from many other disciplines in the typical
117 level of within-study replication, which is fewer than 10 replicates per study (Hillebrand &
118 Gurevitch, 2014). Such low levels of replication will influence the precision of the estimates
119 of effect size from the primary studies (Langan, Higgins, & Simmons, 2016). Importantly,
120 the low level of replication typical of ecological studies is outside the range used in most
121 simulation studies designed to assess meta-analytic methods, which typically range from

122 dozens to hundreds (Langan et al., 2016). Thus differences between ecology and other
123 disciplines potentially limit the insights ecologists can gain from existing simulations that
124 compare different meta-analytic methods.

125 Specific advice for conducting ecological meta-analyses include suggestions on the
126 type of meta-analytic model and effect size calculation to use (Gurevitch & Hedges, 1999;
127 Osenberg, Sarnelle, Cooper, & Holt, 1999; Lajeunesse, 2015), and how to deal with non-
128 independence (Gurevitch & Hedges, 1999; Noble, Lagisz, O’dea, & Nakagawa, 2017; Song,
129 Peacor, Osenberg, & Bence, 2020). For example, a random-effects model is often
130 recommended for ecological meta-analysis over a fixed-effects model (Gurevitch & Hedges,
131 1999), and multi-level models are increasingly being used to incorporate the non-
132 independence commonly found in ecological meta-analyses (Nakagawa & Santos, 2012). A
133 topic addressed in the medical literature that has received little attention in ecology (but see
134 Adams, Gurevitch, & Rosenberg, 1997) is the choice of confidence interval (CI) used to
135 estimate the mean effect size in a meta-analysis (Hartung & Knapp, 2001; Sidik & Jonkman,
136 2003, Sánchez-Meca & Marín-Martínez, 2008).

137 Simulation studies have shown that when the number of studies (k) in the meta-
138 analysis is low, the CIs for a mean effect size calculated using a normal approximation are
139 too narrow, leading to coverage below the nominal level (i.e., a 95% CI should include the
140 true value 95% of the time) (Brockwell & Gordon, 2001; Sánchez-Meca & Marín-Martínez,
141 2008). To avoid this problem, meta-analyses in the medical literature often use the HKSJ
142 (Hartung-Knapp-Sidik-Jonkman; Hartung & Knapp, 2001; Sidik & Jonkman, 2003) method,
143 which is based on a t -distribution and can achieve good coverage even when k is small
144 (Inthout, Ioannidis, & Borm, 2014). Bootstrap techniques have been recommended for
145 estimating CIs for means in ecological meta-analyses, due to its robustness to departures

146 from normality (Adams et al., 1997). On the other hand, boot-strapped CIs can lead to poor
147 coverage when estimating the among-study variance (Viechtbauer, 2007).

148 Bayesian methods, and the credible interval, offer an alternative approach to
149 estimating uncertainty in meta-analyses. Although Bayesian methods may have a steep
150 learning curve, they offer advantages in handling hierarchical models, for incorporating prior
151 information, and for dealing with missing data (Ogle, Barber, & Sartor, 2013). Bayesian
152 meta-analytic techniques produce a posterior distribution of the mean effect size and
153 associated variance terms. Estimates of uncertainty, including credible intervals, can be
154 directly obtained from the posterior distributions, offering an easier to interpret alternative to
155 the frequentist-based CI (Kruschke & Liddell, 2008).

156 Our main goal is to compare the performance of traditional and Bayesian methods to
157 measure the uncertainty around the estimation of a mean effect in the context of ecological
158 meta-analysis. To achieve this goal, we conducted a two-pronged study. First, we reviewed
159 published ecological meta-analyses to characterize the types of confidence interval used in
160 ecological meta-analyses, the number of replicates used in the primary studies (n_i) included
161 in published meta-analyses, and the number of studies (k) that were aggregated to calculate
162 a mean effect size. Second, we used the n_i and k found in our literature review to inform the
163 range of parameter values to use in conducting simulation experiments relevant to ecological
164 meta-analyses. In particular, we determined the typical levels of n_i , k , and the among-study
165 variance and then varied them systematically in our simulation studies. We then evaluated
166 performance of frequentist and Bayesian meta-analysis methods when applied to the
167 simulated data, especially with respect to their ability to estimate the true mean effect and
168 among-study variance, and their quantification of uncertainty intervals (i.e., confidence or
169 credible intervals). Based on our findings, we generate recommendations on the methods to

170 measure uncertainty that perform best for ecological meta-analysis and highlight how simple
171 choices (sometimes overlooked by the investigators) can affect the results of meta-analyses.
172

173 **MATERIALS AND METHODS**

174 **Literature review to assess characteristics of ecological datasets**

175 **Literature search.** We searched the Core Collection of the ISI Web of Science database in
176 March 2017; the search string for TOPIC included ([“meta-analy*” OR “metaanaly*” OR
177 “meta analy*”] AND [“climate change” OR “global change”]). We only included articles
178 and reviews within the “Ecology”, “Environmental Sciences”, “Biodiversity Conservation”
179 and “Plant Sciences” categories. The search resulted in 581 citations; the PRISMA diagram
180 detailing the screening process is provided in Figure S1. After abstract screening, we checked
181 the full text of the 205 articles published between 2013 and 2016. Of these, 96 papers satisfied
182 the inclusion criteria for the final analysis.

183 **Criteria for inclusion.** We focused on narrow sense meta-analyses: i.e., those that used a
184 quantitative meta-analytic method to combine effect sizes that compared a control and a
185 treatment group. We excluded studies that 1) only cited published meta-analyses, 2) reviewed
186 meta-analytic methods, but did not perform a meta-analysis, 3) were identified as meta-
187 analysis by the authors but did not use a meta-analytic model or did not calculate effect sizes,
188 4) used the correlation between two variables as an effect size, and 5) were not “biological
189 meta-analyses” (as defined in Nakagawa et al., 2017), such as studies related to human health
190 or human behavior.

191 **Information extracted.** For each paper we extracted the number of studies (k) from the text,
192 figure captions, figures, and supplementary materials. Here we define a “study” as yielding
193 an estimate of an effect, so that a given primary paper could generate multiple effects and
194 thus multiple studies. The k values were determined at three levels, 1) overall: i.e., the total
195 k collected by the authors (e.g., if they conducted meta-analyses on different response
196 variables, then we summed the k across these variables); 2) analysis: i.e., the total k used in
197 a particular analysis (e.g., if an analysis examined variation among four levels of a moderator,
198 then we summed up the number of studies in each level); and 3) category: i.e., the k included
199 in each category of a categorical analysis. In some cases, authors calculated mean effect sizes
200 for different categories separately and only compared the categories using confidence
201 intervals (i.e., there was no integrated analysis incorporating a category effect). In this case,
202 we considered each’s categories’ k to apply at the “analysis” level.

203 When available, we also recorded the number of replicates (n_i) in the original studies.
204 If the level of replication was unequal for the control and treatment groups, we recorded the
205 average. Finally, from each meta-analysis, we also recorded the inferential paradigm used
206 (frequentist vs. Bayesian) and the method used to obtain confidence intervals for the
207 frequentist approaches (e.g., non-parametric bootstrap, normal-based, KHSJ, etc.).

208 **Simulation Experiments**

209 Our literature review showed that 67% of the reported primary studies had less than ten
210 replicates. In addition, the review of meta-analyses in ecology and evolution by Senior et al.
211 (2012) showed that among-study variation was important, and typically large, in ecological
212 studies. Given these characteristics of ecological data, we simulated data in a full-factorial
213 design that considered the following levels: mean number of replicates $n = \{3, 5, 10, 15, 20,$
214 $30\}$, number of studies $k = \{5, 10, 15, 25, 35, 50\}$, and among-study variance $\sigma_{among}^2 = \{0.1,$

215 0.25, 0.5, 1, 2, 5}. We simulated 2,000 replicated meta-analyses for each combination of n ,
 216 k , and σ_{among}^2 . We then evaluated the performance of four meta-analytic methods applied to
 217 the simulated data: three frequentist approaches that differed in how they calculated
 218 confidence intervals for a mean effect and a Bayesian approach.

219 **Simulating raw data for a study.** We first determined the number of replicates for study i
 220 (n_i) based on a random draw from a Poisson distribution:

$$n_i^* \sim \text{Poisson}(n - 2) \quad (\text{Eq. 1})$$

$$n_i = n_i^* + 2 \quad (\text{Eq. 2})$$

221 where n is the mean number of replicates representative of ecological meta-analyses. We
 222 subtracted 2 to sample from the Poisson and added 2 to the simulated n_i^* to make the
 223 minimum number of replicates for each simulated study equal 2 rather than 0. For each study,
 224 we assumed equal number of replicates for the control and treatment groups.

225 Individual observations ($j = 1, 2, \dots, n_i$) for the control and treatment groups were
 226 generated from a lognormal distribution (LN) such that for study i and observation j :

$$y_{Cij} \sim LN(0, \sigma_{rep}^2) \quad (\text{Eq. 3})$$

$$y_{Tij} \sim LN(0 + \mu + \varepsilon_i, \sigma_{rep}^2) \quad (\text{Eq. 4})$$

229 where σ_{rep}^2 is the among-replicates variation, μ is the true overall effect, and y_{Cij} and y_{Tij}
 230 are the simulated observations for study i and observation j of the control and treatment
 231 group, respectively. We set the among-replicate variation equal to 1 for both the control and
 232 treatment. For convenience, we set the location parameter for the control group equal to zero,
 233 resulting in median (y_C) = 1. For the treatment group in study i , we set median (y_T) = $\mu + \varepsilon_i$,

234 where μ is the overall true treatment effect (hereafter, true effect size) and ε_i is the random
 235 effect associated with study i . We simulated ε_i as:

$$236 \quad \varepsilon_i \sim N(0, \sigma_{among}^2) \quad (\text{Eq. 5})$$

237 Thus, the true effect size from any given study departs from μ due to its random effect
 238 (determined by ε_i), while the estimated effect size differs from the true effect size due to
 239 within-study sampling error (i.e., as influenced by n_i and σ_{rep}^2). The range of values used for
 240 σ_{among}^2 were chosen to produce a similar distribution of I^2 (the proportion of variation among
 241 effect sizes not explained by sampling error) to that reported by Senior et al. (2016) for meta-
 242 analyses in ecology and evolution (I^2 simulation results are presented in Figure S2).

243 **Estimating the effect size and within-study variance.** Using the raw data simulated from
 244 each study, we computed the observed effect size for study i as the log response ratio ($lnRR_i$),
 245 which is widely used in ecology (Nakagawa & Santos, 2012) and it is often a reasonable
 246 approximation of ecological phenomena (Osenberg, Sarnelle, & Cooper, 1997):

$$247 \quad lnRR_i = \ln\left(\frac{\bar{y}_{T_i}}{\bar{y}_{C_i}}\right) \quad (6)$$

248 where \bar{y}_{T_i} and \bar{y}_{C_i} are the sample means of the treatment and control groups, respectively.

249 The expected sample means for each treatment in a simulated study are $E(y_{C_{ij}}) =$
 250 $\exp\left(\frac{\sigma_{rep}^2}{2}\right)$ and $E(y_{T_{ij}}) = \exp\left(\mu + \varepsilon_i + \frac{\sigma_{rep}^2}{2}\right)$. Thus, the log of the ratio of the expected
 251 values for the treatment and control groups is $\mu + \varepsilon_i$, corresponding to what we call the true
 252 study-specific effect size.

253 We calculated the estimated within-study variance of the log ratio (Eq. 1 in Hedges,
 254 Gurevitch, & Curtis, 1999) ($\sigma_{within_i}^2$) as:

255
$$\sigma_{within_i}^2 = \frac{SD_{T_i}^2}{n_{T_i} \bar{y}_{T_i}^2} + \frac{SD_{C_i}^2}{n_{C_i} \bar{y}_{C_i}^2} \quad (7)$$

256 where SD_T and SD_C are the sample standard deviations of the treatment and control groups,
 257 respectively, and $n_{T_i} = n_{C_i} = n_i$ are the simulated number of replicates in study i .

258

259 **Meta-analytic approaches**

260 Given that we simulated independent data to highlight how the choice of uncertainty interval
 261 affects the estimation of a mean effect, we used a standard random-effects model (Gurevitch
 262 & Hedges, 1999). We comment on how our results may change with a multi-level
 263 (hierarchical) model in the Discussion section. We assume the simulated effect size for study
 264 i ($\ln RR_i$, calculated from Eq. 6) follows a normal distribution with mean θ_i (the true effect
 265 for study i) and within-study variance $\sigma_{within_i}^2$:

266
$$\ln RR_i \sim N(\theta_i, \sigma_{within_i}^2) \quad (8)$$

267
$$\theta_i \sim N(\mu, \sigma_{among}^2) \quad (9)$$

268 We assume $\sigma_{within_i}^2$ is known, as calculated via Eq. 7. Likewise, the true study-specific effect
 269 size, θ_i , is assumed to follow a normal distribution with mean μ (the true overall effect) and
 270 among-study variance, σ_{among}^2 (which is sometimes referred to as τ^2 in other meta-analytic
 271 papers).

272 We compared different methods to construct confidence intervals (CIs) for a mean
 273 effect (at the analysis level) within the frequentist methods versus Bayesian credible
 274 intervals. For the frequentist-based analyses, we compared: a) a CI based on a z -distribution,
 275 which is a large sample approximation, b) a weighted CI based on the Hartung-Knapp-Sidik-
 276 Jonkman (HKSJ) method, which does not assume a large sample and instead uses a t -

277 distribution, and c) bootstrap methods. For the Bayesian-based analysis, we calculated the
278 highest posterior density (HPD) credible interval.

279 **Frequentist approaches.** We applied the random-effects model described by Eqs. 8 and 9
280 with inverse variance weights using the “rma” function in the R package *metafor*
281 (Viechtbauer, 2010), and estimated σ_{among}^2 with the default REML method. To calculate the
282 z-distribution CI, we used the default settings for the random-effects model in *metafor*, which
283 returns a 95% CI for μ based on the normal distribution. To apply the HKSJ CI, we set the
284 option knha=T in *metafor*. The resulting CI for μ is based on both a refined estimate of
285 σ_{among}^2 and a Student’s t-distribution (Hartung & Knapp, 2001; Sidik & Jonkman, 2003),
286 which accounts for the fact that σ_{among}^2 is estimated and not known. For the bootstrapped CI,
287 we estimated the bias-corrected non-parametric bootstrapped 95% CI for both μ and σ_{among}^2
288 via the *boot* package in R (Canty & Ripley, 2017). Since the choice of HKSJ or z-distribution
289 for the μ CI does not affect the estimation of σ_{among}^2 , in both cases we used *metafor*’s
290 function “confint” to obtain the CI for σ_{among}^2 (“confint” applies a Q-profile method in
291 combination with REML).

292 **Bayesian approach.** We used a “hybrid” Bayesian framework to implement the random-
293 effects model (Eqs. 8 and 9) in which we treat σ_{within}^2 as known; whereas a fully Bayesian
294 model may treat σ_{within}^2 as unknown (this hybrid model is comparable to the “empirical
295 Bayes” method discussed in Schmid & Mengersen, 2013). Initial explorations with full and
296 hybrid models gave qualitatively similar results and we only include the hybrid model in our
297 analysis.

298 We specified relatively non-informative priors for the unknown quantities (e.g., μ and
299 σ_{among}^2). For the mean effect size, μ , we specified a conjugate normal prior with a mean of
300 zero and large variance: $N(0, 10000)$. Given that even diffuse priors for σ_{among}^2 can influence
301 the posterior for σ_{among}^2 , particularly under small group size (Gelman, 2006), we explored
302 five different priors for σ_{among}^2 (Supporting Information Figures S12-15). For the final
303 analysis, convergence statistics and computational speed led us to focus on the *Uniform*(0,10)
304 prior for the standard deviation (σ_{among}).

305 The Bayesian meta-analyses were implemented in JAGS with the *rjags* R package
306 (Plummer, 2018). For each model, we ran three parallel Markov chain Monte Carlo (MCMC)
307 sequences for 200,000 iterations, and discarded the first 100,000 iterations as the burn-in
308 period. We used the \hat{R} convergence diagnostic (Gelman & Rubin, 1992) to evaluate
309 convergence of the MCMC sequences to the posterior. For the final simulations, we only
310 included runs that had $\hat{R} < 1.1$, and checked that the proportion of discarded runs was lower
311 than 1%. Using post-burn-in MCMC samples, we computed posterior means for quantities
312 of interest (e.g., μ and σ_{among}^2) as point estimates. We computed 95% credible intervals as
313 HPD intervals for both μ and σ_{among}^2 using the “HPDinterval” function in the *coda* package
314 (Plummer, 2006).

315 **Implementation and Assessment of the Meta-analysis Approaches**

316 We ran all the analyses and simulations in the R environment (R Core Team, 2019); code is
317 provided in the Supporting Information. For each simulated dataset, we estimated μ and
318 σ_{among}^2 via the frequentist and Bayesian methods described above. We summarized the
319 results from the 2,000 replicated meta-analyses for each combination of factors (n, k, σ_{among}^2)

320 and modeling approaches (i.e., frequentist and Bayesian methods to measure uncertainty).
 321 The results for the model performance associated with estimating σ_{among}^2 are presented in
 322 Figures S7-10.

323 We evaluated model performance using: coverage, width of the uncertainty intervals,
 324 bias, and efficiency. We estimated *coverage* for both μ and σ_{among}^2 as the proportion (out of
 325 the 2,000 simulation replicates) of calculated 95% uncertainty intervals (CIs for the
 326 frequentist methods and credible interval for the Bayesian approach) that included the
 327 corresponding true value. Ideally, coverage should equal the nominal value of 0.95 (95%).
 328 CIs for these “coverage proportions” were computed using the “binom.confint” function in
 329 the R *binom* (Sundar, 2014) package, with the method “wilson” (Agresti & Coull, 1998).

330 We summarized the perceived uncertainty for μ and σ_{among}^2 as the mean *width of the*
 331 *95% uncertainty intervals* for the 2,000 intervals for each scenario, and assessed how well
 332 the mean width was estimated using a 95% CI based on a *t*-distribution. All else being equal,
 333 smaller uncertainty is a desirable feature, but not if it is accompanied by a reduction in
 334 coverage below the nominal level.

335 To evaluate *bias*, we calculated the mean difference between the point estimates for
 336 μ and σ_{among}^2 and their true values based on the 2,000 simulation replicates, and report a
 337 95% CI for this estimate based on the *t*-distribution. Ideally, bias should be centered on zero.

338 Finally, to quantify the *efficiency* of the point estimates, we calculated the root mean
 339 squared error (RMSE) between the estimated and true values for μ and σ_{among}^2 as:

$$340 \quad RMSE = \sqrt{\frac{\sum_{s=1}^{N_{sim}} (\hat{a}_s - a_{true_s})^2}{N_{sim}}}, \quad (10)$$

341 where $a = \mu$ or σ_{among}^2 , \hat{a} is the point estimate from each model, a_{true} is the true value used
342 in the simulations, and N_{sim} is the number of simulations.

343

344 **RESULTS**

345 **Literature review to assess characteristic of ecological datasets**

346 Of the 96 meta-analyses that satisfied our criteria (Table S1), 95 and 26 provided information
347 on the number of studies (k) and number of replicates (n_i) associated with the original dataset,
348 respectively. Only three meta-analyses used a Bayesian approach. The majority of meta-
349 analyses were published in *Global Change Biology* (23), followed by *Agriculture Ecosystems*
350 *& Environment* (7) and *Ecology* (6) (Figure S3 displays the full list). The quality of reporting
351 varied, and is discussed in more detail in the Supporting Information. We also provide
352 additional information on k and n_i (by taxa, environment, and topic) in the Supporting
353 Information (Table S2, Figures S4-S5).

354 **Number of studies.** The number of studies (k) used to estimate an effect was highly skewed
355 at the three levels we considered: overall, analysis, and category (Figure 1). The overall k
356 ranged from 25 to 32,567 (Figure 1A upper panel), with a median of 273 and with relatively
357 few (12%) including more than 1,000 studies. For most papers, however, analyses were
358 performed for different response variables or different moderators, and the k used for a
359 particular analysis was considerably lower (Figure 1A middle panel), ranging from $k = 1$ (for
360 a paper that presented all possible comparisons, even when one potential analysis was
361 represented by only a single study) to $k = 8,474$, with a median of $k = 44$ (i.e., 50% of meta-
362 analysis included 44 or fewer studies); 16% had $k \leq 10$. The number of studies included

363 within categories ranged from $k = 1$ to 1,430, with a median of 16; 36% had $k \leq 10$ (Figure
364 1A lower panel).

365 **Number of replicates.** The distribution of the reported number of replicates in the original
366 studies (n_i) cited by the climate change meta-analyses was highly skewed, ranging from $n_i =$
367 1 to 21,600, with most studies having only a few replicates; the median was 5 (Figure 1B).
368 The strong skewness in these data led us to inspect some of the original publications from
369 which exceptionally large n_i values were reported. We found publications in which n_i values
370 were likely misreported or greatly inflated by pseudoreplication (details in Table S3 and
371 Figure S6).

372 **Analytic method to estimate the uncertainty interval for a mean effect.** In 38.5% of the
373 papers reviewed, the method used to calculate the frequentist-based CI for the mean effect
374 was not mentioned (Figure 2). Of the papers reporting how the CI was calculated, the
375 majority used bootstrapped or z -distribution CIs; only three papers used credible intervals
376 (Bayesian method), and a few used a combination of methods (Figure 2). No papers reported
377 using HKSJ method. Of the papers that did not specify the method, nine used Metawin (which
378 defaults to a t -distribution for the parametric CI, without the KHSJ refinement); 12 papers
379 used the packages *meta* or *metafor* in R (which default to a z -distribution); and two used the
380 Comprehensive Meta-Analysis software (which defaults to a z -distribution). Assuming these
381 23 papers used the software defaults, then 31 papers used a z -distribution, and nine used a t -
382 distribution but without the KHSJ refinement. Thus, bootstrapped and z -distribution CIs
383 likely comprise the vast majority of approaches, with KHSJ CIs being entirely absent from
384 our dataset.

385

386 **Simulation experiments: estimation of a mean effect**

387 The number of studies, k , used to estimate a mean effect size, μ , substantially affected the
388 coverage of the frequentist methods, but this effect of k depended on the type of method used
389 to estimate the 95% CIs (Figure 3A). For example, z -distribution CIs for μ had coverage
390 lower than the nominal level when $k < 40$, and coverage was appreciably lower for $k < 20$
391 (Figure 3A). Similarly, bootstrapped CIs had lower than nominal coverage when $k < 40$
392 (Figure 3A). In contrast, KHSJ CIs had close to nominal coverage over all values of k (Figure
393 3A). The Bayesian credible interval generally showed coverages around 95%, but when $k =$
394 5, coverage was $>95\%$ (Figure 3A).

395 Coverage can be smaller than nominal levels either because of bias or because the
396 width of the uncertainty interval is inappropriately narrow (i.e., uncertainty is
397 underestimated). The three frequentist methods for computing CIs for μ used the same
398 approach for obtaining point estimates and had minimal bias centered on zero (Figures S11
399 A,C,E). Thus, the observed differences in coverage for μ resulted from differences in the
400 width of the uncertainty interval (Figure 3B). The Bayesian credible interval was generally
401 wider than the frequentist-based CIs, and of the frequentist CIs, the KHSJ CI tended to be
402 the widest; when k was small, the z -distribution and boot-strapped CIs were $\sim 1/3$ smaller than
403 they should be based upon the more appropriate KHSJ CI (Fig. 3B).

404 Increasing the mean number of replicates (n) in the primary studies did not greatly
405 affect coverage (Figure 3B), the width of the uncertainty interval (Figure 3E), bias (Figure
406 S11C), or RMSE (Figure S11D) for μ . Our results were likely produced because the among-
407 study variation dominated within-study variation over the range of levels considered for the
408 simulation factors (as determined by the review by Senior et al., 2016).

409 Increasing the among-study variance (σ_{among}^2) increased the width of the uncertainty
410 interval for μ (Figure 3F), but had only small effects on coverage (Figure 3C). Bias in the
411 estimation of μ was negligible and unaffected by an increase in σ_{among}^2 (Figure S11E), but
412 the error in the estimation increased with the increase in heterogeneity (RMSE, Figure S11F).

413

414 **DISCUSSION**

415 Our literature review shows that ecological meta-analyses are highly variable in terms
416 of how many studies (k) are included in the meta-analysis and the number of replicates
417 reported in the original publications (n_i). Despite this high variability, both across and within
418 meta-analyses, k and n_i tend to be low. The high frequency of meta-analyses with
419 comparatively few studies ($k \leq 44$ in 50% of meta-analyses reviewed) is not unique to
420 ecology; even lower number of studies are pervasive in medical research (Kontopantelis,
421 Springate, & Reeves, 2013) where there has been an effort to develop methods that improve
422 the performance of meta-analyses in such scenarios (Inthout et al., 2014). Furthermore, our
423 simulations show that the method used to calculate an uncertainty interval greatly influences
424 how often the interval includes the true mean effect and is very important for producing
425 intervals with close to correct coverage when k is low. Despite its importance, a large
426 proportion of the ecological meta-analyses we reviewed (38%) did not report the type of
427 uncertainty interval used, and the ones that did report their methods used intervals that are
428 problematic when k is low.

429 Low coverage of the z-distribution confidence interval (CI) when the number of
430 observations (in the meta-analysis context, the number of studies, k) are low is well known
431 in classical statistical contexts as well as in meta-analyses (Hedges et al., 1999; Brockwell &

432 Gordon, 2001; IntHout et al., 2014). In meta-analyses, however, approaches typically default
433 to assuming large k and thus justify the application of the z -distribution. In ecology, this
434 large-sample approach is often unwarranted (Figure 1A). Furthermore, bootstrapped CIs are
435 also well known to be problematic with small k (Hesterberg, 2015), although ecological meta-
436 analyses tend to prioritize the potential for non-normal distributions over concerns about
437 small k (Adams et al., 1997) – based upon our results, such prioritization may be misplaced.

438 When k is low, the CI for a mean effect size (μ) based on the z -distribution is too
439 narrow. Some practitioners have addressed this problem by not calculating CIs when k is
440 very small (e.g.: Augusto, Delerue, Gallet-Budynek, & Achat, 2013). Others have resorted
441 to using bootstrapped CIs (e.g.: Thébault, Mariotte, Lortie, & MacDougall, 2014). Given that
442 bootstrapped CIs also had poor coverage when $k < 40$, this approach appears to be ill-advised.
443 In our review, nearly half of the mean effect sizes used in an individual analysis were
444 calculated with $k < 40$ effect sizes, where the choice of method for computing uncertainty
445 intervals matters. As a result, many effects declared as significant probably should not have
446 been. This is exemplified in a review of medical meta-analyses from the Cochrane Database,
447 where of the 315 meta-analyses that yielded significant effects with the z -distribution CI,
448 only 79 were significant using the HKSJ CI (IntHout et al., 2014).

449 The default option for frequentist CIs for μ varies among software packages. For
450 example, a t -distribution CI (but without the HKSJ refinement) is Metawin's default, whereas
451 the z -distribution is the default in the Comprehensive Meta-Analysis software and in the R
452 packages *meta* and *metafor* (*metafor* is one of the most common software packages currently
453 in use by ecologists). For those planning to conduct a random-effects meta-analysis using
454 frequentist methods, we advise use of the HKSJ CI, which employs both a weighted estimator
455 of the variance for the overall effect size and a t -distribution for its associated CI (this can be

456 set up in *metafor* using the option `knha= T`). Sánchez-Meca and Marín-Martínez (2008)
457 report that the HKSJ method outperforms the simple CI-based on the t -distribution. However,
458 in some scenarios, coverage could be as low as 90% even using the HKSJ CI, for example,
459 when heterogeneity is high, $k < 10$, and the number of replicates varies greatly among studies
460 (Inthout et al., 2014). In our simulations that did not include highly uneven number of
461 replicates, we showed that HKSJ CI's and the Bayesian credible intervals provide accurate
462 (or at least conservative, >95%) coverage and performed best. We encourage researchers to
463 be aware of the software defaults when calculating an uncertainty interval, and to report the
464 method used.

465 The climate change meta-analyses showed exceedingly high variation in the number
466 of replicates reported (n_i), spanning five orders of magnitude, but the majority of values were
467 low. In fact, $n_i < 10$ in 67% of the cases, and $n_i \leq 5$ in 51% of the cases we reviewed. This
468 pattern may be similar in other fields of ecology (Table S2, Figures S4, S5). For example, a
469 competition meta-analysis found n_i ranging from 1 to 1,455, with a median of 10 (Gurevitch
470 et al., 1992). To obtain a more accurate estimate of μ , some authors specify a minimum n_i to
471 calculate mean effect sizes (Gurevitch et al., 1992; Schirmel et al., 2016). Such censoring
472 might improve confidence interval performance by reducing variation in replication among
473 studies (Inthout et al. 2014) but at the high cost of discarding important information. While
474 one would in general expect better estimates with more replication, our simulation
475 experiment did not show important effects of the mean number of replicates on the estimation
476 of and inferences about μ . A similar insensitivity to the number of replicates has been
477 observed in other studies (Sánchez-Meca & Marín-Martínez 2008), although we included
478 fewer replicates than most other simulations. Variation in replication among studies, should
479 produce variation in within-study variance, especially when the number of replicates is small.

480 However, in our simulations among-study variation was much larger than within-study
481 variation, consistent with the characteristics of ecological meta-analyses (Senior et al., 2016),
482 minimizing the role of variation in the number of replicates.

483 When the number of replicates reported (n_i) was unusually high, we checked a few
484 of the original papers cited in each meta-analysis. Upon revisiting 17 of the original
485 publications, we found at least 15 cases in which n_i was misreported (Table S3). This
486 manifested in different ways. Some meta-analyses reported the total n_i in an experiment
487 instead of the number of replicates per treatment. In other cases, authors reported the total n_i
488 from repeated measurements or the numbers of individuals rather than the number of true
489 replicates (e.g., plots or cages). There were also cases in which we were unable to verify the
490 origin of the number reported in the meta-analysis. An incorrect n_i decreases the sampling
491 variance for that effect size, which affects the weights and also the estimation of the overall
492 heterogeneity (Noble et al., 2017). Researchers conducting a meta-analysis should be
493 cautious when extracting data from the original studies to avoid misreporting (or inflating)
494 the number of replicates. Publication of the data and code used to conduct a meta-analysis
495 would also be useful to inform research on best practices for meta-analysis.

496 In our simulations using a random-effects model, the performance in the estimation
497 of the among-study variance (σ_{among}^2) was better when the true σ_{among}^2 was high (Figures
498 S4-7). In agreement with Viechtbauer (2007), we observed that the Q-profile CI method for
499 σ_{among}^2 performed better than the bootstrap method (Figures S7-10). The Bayesian method
500 performed best, but had coverage above the nominal level when the number of studies was
501 low ($k < 20$). Bayesian methods led to higher perceived uncertainty in such cases, which
502 could be real, but this could also be a consequence of positive bias in the σ_{among}^2 estimates,
503 which was more pronounced for the Bayesian methods when $k < 20$. In this scenario, one

504 approach to improve coverage is to use priors for σ_{among}^2 that perform better when k is low
505 (Gelman, 2006). Another solution is to specify more informative priors for σ_{among}^2 based on
506 a synthesis of past publications (Higgins et al., 2009). One reason to desire good estimation
507 of σ_{among}^2 is because overestimation of this variance component can lead to higher perceived
508 uncertainty in the estimate of μ . An additional reason is that the estimates of σ_{among}^2 represent
509 real variation in effects and could be of importance in risk assessment.

510 In the initial explorations with the full Bayesian model, the MCMC chains for μ
511 converged quickly, but they converged more slowly for σ_{among}^2 , often falling into a “zero
512 variance trap” (Gelman, 2004) when the true among-study variance was close to zero. In
513 general, convergence and mixing problems were most frequent for low k and low σ_{among}^2 .
514 While low σ_{among}^2 is rare in ecology, low k is not. Of the priors we explored (Supporting
515 Information Figures S12-15), the folded- t and the uniform prior for the standard deviation
516 performed best when k was low (we chose the uniform prior for the final simulations because
517 it ran slightly faster). In our simulations, the hybrid Bayesian model exhibited the practical
518 advantages of the Bayesian methods (e.g., produces full posteriors and direct evaluation of
519 uncertainty without approximating assumptions, among others), and was easy (and faster) to
520 implement than the full model. On the other hand, a full Bayesian approach may be more
521 useful for multi-level models that include missing data, hierarchical structures, and/or
522 covariate effects (Ogle et al., 2013), and could benefit from informative priors for σ_{among}^2 ,
523 particularly when k is low.

524 Our study simulated independent effect sizes. Often though, observed effect sizes are
525 not independent (e.g., multiple observed effect sizes might be obtained from a single
526 published article). As observed effect sizes within a group might respond similarly (due to

527 similar methods, or similar environmental conditions), some of the among-study variation
528 could be common to all members of a group or subgroup. Multi-level (hierarchical) models
529 can be used to account for this. We believe that our results, including the insensitivity of our
530 results to n , would not be materially altered in such situations, assuming the among-study
531 variation still dominates the within-study variation. There are some challenges to be faced,
532 however, when applying our results to more complex multi-level models. In particular,
533 although the R package *metafor* has a function that handles multi-level models (`rma.mv`), the
534 KHSJ adjustment is not available in this context, and the best that can be done with *metafor*
535 is to construct t -based confidence intervals of the mean (also referred to as conditional t -test).
536 For multi-level models, these t -based confidence intervals have inflated error rates (Luke,
537 2017; Song et al., in press), although they do outperform normal-based confidence intervals
538 (Song, personal communication). Song et al. (in press) speculated that the inflated error rates
539 of t -based confidence intervals resulted from not accounting for uncertainty in estimated
540 variances. Methods exist for adjusting tests and confidence intervals to account for
541 uncertainty in estimated variances in multi-level models, such as the Kenward-Rogers
542 adjustment, or simulation of null distributions (Halekoh & Hojsgaard, 2014), but to our
543 knowledge these have not been implemented in any readily available software for conducting
544 meta-analyses.

545

546 **AUTHORS CONTRIBUTIONS**

547 All authors conceived the idea; PP collected and analyzed the data with contributions from
548 CWO, EAH, JRB, and KO. PP led the writing; all authors contributed critically to the drafts
549 and gave final approval for publication.

550

551 **DATA AVAILABILITY**

552 The data compiled in the literature review, the R code for the simulation experiment, and the
553 results from the simulation experiments are deposited in Dryad repository:
554 <https://doi.org/10.5061/dryad.zw3r22863>.

555

556 **ACKNOWLEDGMENTS**

557 This research was funded by the U.S. Department of Energy (DE-SC-0010632), National
558 Science Foundation (DEB-1655426 and DEB-1655394), and utilized Georgia Advanced
559 Computing Resource Center resources. This is publication 2020-17 of the Quantitative
560 Fisheries Center. We thank W. Viechtbauer for feedback on *metafor* computations, Natasja
561 vanGestel and Kees Jan vanGroenigen for feedback during early stages of the project, Chao
562 Song for helpful discussion, and Sergio Estay for his feedback on the abstract in Spanish.

563

564 **REFERENCES**

565 Adams, D. C., Gurevitch, J., & Rosenberg, M. S. (1997). Resampling tests for meta-
566 analysis of ecological data. *Ecology*, 78(4), 1277. <https://doi.org/10.2307/2265879>

567 Agresti, A., & Coull, B. A. (1998). Approximate is better than “exact” for interval
568 estimation of binomial proportions. *The American Statistician*, 52(2), 119–126.
569 <https://doi.org/10.1080/00031305.1998.10480550>

570 Augusto, L., Delerue, F., Gallet-Budynek, A., & Achat, D. L. (2013). Global assessment of
571 limitation to symbiotic nitrogen fixation by phosphorus availability in terrestrial
572 ecosystems using a meta-analysis approach. *Global Biogeochemical Cycles*, 27(3),
573 804–815. <https://doi.org/10.1002/gbc.20069>

- 574 Brockwell, S. E., & Gordon, I. R. (2001). A comparison of statistical methods for meta-
575 analysis. *Statistics in Medicine*, 20(6), 825–840. <https://doi.org/10.1002/sim.650>
- 576 Cadotte, M. W., Mehrkens, L. R., & Menge, D. N. L. (2012). Gauging the impact of meta-
577 analysis on ecology. *Evolutionary Ecology*, 26(5), 1153–1167.
578 <https://doi.org/10.1007/s10682-012-9585-z>
- 579 Canty, A., & Ripley, A. (2017). boot: Bootstrap R (S-Plus) Functions (Version R package
580 version 1.3-20).
- 581 Gates, S. (2002). Review of methodology of quantitative reviews using meta-analysis in
582 ecology. *Journal of Animal Ecology*, 71(4), 547–557.
583 <https://doi.org/10.1046/j.1365-2656.2002.00634.x>
- 584 Gelman, A. (2004) Parameterization and Bayesian modeling. *Journal of the American*
585 *Statistical Association*, 99(466), 537–545
- 586 Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models.
587 *Bayesian Analysis*, 1(3), 515–533.
- 588 Gelman, A., & Rubin, D. B. (1992). Inference from iterative simulation using multiple
589 sequences. *Statistical Science*, 7(4), 457–472.
590 <https://doi.org/10.1214/ss/1177011136>
- 591 Gurevitch, J., & Hedges, L. V. (1999). Statistical issues in ecological meta-analyses.
592 *Ecology*, 80(4), 1142–1149.

- 593 Gurevitch, J., Morrow, L. L., Wallace, A., & Walsh, J. S. (1992). A meta analysis of
594 competition in field experiments. *The American Naturalist*, *140*(4), 539–572.
595 <https://doi.org/10.1086/285428>
- 596 Halekoh, U. & Hojsgaard, S. (2014). A Kenward-Roger approximation and parametric
597 bootstrap methods for tests in linear mixed models – the R package pbrktest.
598 *Journal of Statistical Software*, *59*(9), 1-32. <https://doi.org/10.18637/jss.v059.i09>
- 599 Hartung, J., & Knapp, G. (2001). A refined method for the meta-analysis of controlled
600 clinical trials with binary outcome. *Statistics in Medicine*, *20*(24), 3875–3889.
601 <https://doi.org/10.1002/sim.1009>
- 602 Hedges, L. V., Gurevitch, J., & Curtis, P. S. (1999). The meta-analysis of response ratios in
603 experimental ecology. *Ecology*, *80*(4), 1150. <https://doi.org/10.2307/177062>
- 604 Hesterberg, T. C. (2015). What teachers should know about the bootstrap: resampling in the
605 undergraduate statistics curriculum, *The American Statistician*, *69*(4), 371-386
- 606 Higgins, J. P. T., Thompson, S. G., & Spiegelhalter, D. J. (2009). A re-evaluation of
607 random-effects meta-analysis. *Journal of the Royal Statistical Society: Series A*
608 *(Statistics in Society)*, *172*(1), 137–159. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-985X.2008.00552.x)
609 [985X.2008.00552.x](https://doi.org/10.1111/j.1467-985X.2008.00552.x)
- 610 Hillebrand, H., & Gurevitch, J. (2014). Meta-analysis results are unlikely to be biased by
611 differences in variance and replication between ecological lab and field studies.
612 *Oikos*, *123*(7), 794–799. <https://doi.org/10.1111/oik.01288>
- 613 IntHout, J., Ioannidis, J. P., & Borm, G. F. (2014). The Hartung-Knapp-Sidik-Jonkman
614 method for random effects meta-analysis is straightforward and considerably

615 outperforms the standard DerSimonian-Laird method. *BMC Medical Research*
616 *Methodology*, 14(1). <https://doi.org/10.1186/1471-2288-14-25>

617 Kontopantelis, E., Springate, D. A., & Reeves, D. (2013). A re-analysis of the Cochrane
618 library data: the dangers of unobserved heterogeneity in meta-analyses. *PLoS ONE*,
619 8(7), e69930. <https://doi.org/10.1371/journal.pone.0069930>

620 Kruschke, J.K. & Liddell, T.M. (2018). The Bayesian new statistics: hypothesis testing,
621 estimation, meta-analysis, and power analysis from a Bayesian perspective.
622 *Psychonomic Bulletin & Review*, 25, 178-206. [https://doi.org/10.3758/s13423-016-](https://doi.org/10.3758/s13423-016-1221-4)
623 1221-4

624 Lajeunesse, M. J. (2015). Bias and correction for the log response ratio in ecological meta-
625 analysis. *Ecology*, 96(8), 2056–2063. <https://doi.org/10.1890/14-2402.1>

626 Langan, D., Higgins, J. P. T., & Simmonds, M. (2016). Comparative performance of
627 heterogeneity variance estimators in meta-analysis: a review of simulation studies.
628 *Research Synthesis Methods* 8(2), 181-198. <https://doi.org/10.1002/jrsm.1198>

629 Luke, S.G. (2017). Evaluating significance in linear mixed-effect models in R. *Behavior*
630 *Research Methods*, 49, 1494-1502. <https://doi.org/10.3758/s13428-016-0809-y>

631 Nakagawa, S., & Santos, E. S. A. (2012). Methodological issues and advances in biological
632 meta-analysis. *Evolutionary Ecology*, 26(5), 1253–1274.
633 <https://doi.org/10.1007/s10682-012-9555-5>

634 Nakagawa, S., Noble, D. W., Senior, A.M. & Lagisz, M. (2017). Meta-evaluation of meta-
635 analysis: ten appraisal questions for biologists. *BMC Biology*, 15:18.
636 <https://doi.org/10.1186/s12915-017-0357-7>

637 Noble, D. W. A., Lagisz, M., O’dea, R. E., & Nakagawa, S. (2017). Nonindependence and
638 sensitivity analyses in ecological and evolutionary meta-analyses. *Molecular*
639 *Ecology*, 26(9), 2410–2425. <https://doi.org/10.1111/mec.14031>

640 Ogle, K., Barber, J., & Sartor, K. (2013). Feedback and Modularization in a Bayesian
641 Meta-analysis of Tree Traits Affecting Forest Dynamics. *Bayesian Analysis*, 8(1),
642 133–168. <https://doi.org/10.1214/13-BA806>

643 Osenberg, C. W., Sarnelle, O., & Cooper, S. D. (1997). Effect size in ecological
644 experiments: the application of biological models in meta-analysis. *The American*
645 *Naturalist*, 150(6), 798–812.

646 Osenberg, C. W., Sarnelle, O., Cooper, S. D., & Holt, R. D. (1999). Resolving ecological
647 questions through meta-analysis: goals, metrics, and models. *Ecology*, 80(4), 1105–
648 1117.

649 Pappalardo, P. K. Ogle, E.A. Hamman, J.R. Bence, B.A. Hungate, & C.W. Osenberg.
650 (2020). Data from: Comparing traditional and Bayesian approaches to ecological
651 meta-analysis. *Methods in Ecology and Evolution* doi:/10.5061/dryad.zw3r22863

652 Plummer, M. (2018). rjags: Bayesian Graphical Models using MCMC (Version R package
653 version 4-8.).

654 R Core Team. (2019). *R: A language and environment for statistical computing*. Vienna,
655 Austria.: R Foundation for Statistical Computing.

656 Sánchez-Meca, J., & Marín-Martínez, F. (2008). Confidence intervals for the overall effect
657 size in random-effects meta-analysis. *Psychological Methods*, 13(1), 31–48.
658 <https://doi.org/10.1037/1082-989X.13.1.31>

- 659 Schirmel, J., Bundschuh, M., Entling, M. H., Kowarik, I., & Buchholz, S. (2016). Impacts
660 of invasive plants on resident animals across ecosystems, taxa, and feeding types: a
661 global assessment. *Global Change Biology*, 22(2), 594–603.
662 <https://doi.org/10.1111/gcb.13093>
- 663 Schmid, C. H., & Mengersen, K. (2013). Bayesian meta-analysis. In *Handbook of meta-*
664 *analysis in ecology and evolution* (pp. 145–173). Princeton, New Jersey: Princeton
665 University Press.
- 666 Senior, A. M., Grueber, C. E., Kamiya, T., Lagisz, M., O’Dwyer, K., Santos, E. S. A., &
667 Nakagawa, S. (2016). Heterogeneity in ecological and evolutionary meta-analyses:
668 its magnitude and implications. *Ecology*, 97(12), 3293–3299.
669 <https://doi.org/10.1002/ecy.1591>
- 670 Sidik, K., & Jonkman, J. N. (2003). On constructing confidence intervals for a standardized
671 mean difference in meta-analysis. *Communications in Statistics - Simulation and*
672 *Computation*, 32(4), 1191–1203. <https://doi.org/10.1081/SAC-120023885>
- 673 Sundar, D.-R. (2014). binom: binomial confidence intervals for several parameterizations
674 (Version R package version 1.1-1).
- 675 Song, C., Peacor, S.D., Osenberg, C.W., & Bence, J.R. (2020). An assessment of statistical
676 methods for non-independent data in ecological meta-analyses. *Ecology* (under
677 review).
- 678 Thébault, A., Mariotte, P., Lortie, C. J., & MacDougall, A. S. (2014). Land management
679 trumps the effects of climate change and elevated CO₂ on grassland functioning.
680 *Journal of Ecology*, 102(4), 896–904. <https://doi.org/10.1111/1365-2745.12236>

681 Viechtbauer, W. (2007). Confidence intervals for the amount of heterogeneity in meta-
682 analysis. *Statistics in Medicine*, 26(1), 37–52. <https://doi.org/10.1002/sim.2514>

683 Viechtbauer, W. (2010). Conducting Meta-Analyses in R with the metafor Package.
684 *Journal of Statistical Software*, 36(3). <https://doi.org/10.18637/jss.v036.i03>

685

686 **FIGURE LEGENDS**

687 **Figure 1.** Results from the literature review of ecological meta-analyses: A) distribution of
688 the number of studies (k) reported for overall, analysis, and category levels; the median k is
689 indicated in each panel; B) distribution of the number of replicates used in the original studies
690 (n_i), as reported in each meta-analysis; the median n_i is indicated with a dashed line. Note
691 that the x-axes are on a log scale.

692

693 **Figure 2.** Types of uncertainty intervals reported by the ecological meta-analyses. In some
694 cases, more than one type of uncertainty interval was reported.

695

696 **Figure 3.** Coverage and the width of the 95% uncertainty interval for different methods used
697 to estimate the mean effect size (μ) in a meta-analysis as a function of the number of studies
698 (A, D), the mean number of replicates (B, E), and the among-study variance (C, F). The
699 dashed horizontal line in panels A, B, and C indicates the nominal value of 95%. Different
700 colors denote the method used to estimate the uncertainty interval. Error bars provide the
701 95% CI.

702

703 **SUPPLEMENTARY FIGURE LEGENDS**

704

705 **Figure S1.** PRISMA diagram.

706 **Figure S2.** Mean I^2 as a function of the true (simulated) among-study variance for different
707 combinations of the mean number of replicates, n_i , and number of studies, k , in the simulated
708 datasets.

709 **Figure S3.** Number of climate-change meta-analyses reviewed, summarized by journal in
710 which each was published, between 2013 and 2016.

711 **Figure S4.** Results from the exploratory literature search on sub-disciplines of ecological
712 meta-analyses. A) distribution of the number of studies (k) by sub-discipline; B) distribution
713 of the number of replicates (n_i) used in the primary papers, as reported in each meta-analysis.
714 Replication was not reported in any meta-analyses for ocean acidification. Note that the x-
715 axes are on a log scale.

716 **Figure S5.** Additional results for the climate/global change meta-analysis. Variability on the
717 median number of studies at the analysis level (A) and the median number of replicates (B)
718 by type of organism (or variable) measured, type of environment, and meta-analysis topic.

719 **Figure S6.** Distribution of the number of replicates, n_i , in the original studies for each of the
720 26 meta-analysis publications in our review that provided the original data. The boxplots
721 represent the median (thick vertical line), the 25th and 75th percentiles (box), the upper
722 whisker extends from the box to the larger value no further than 1.5xIQR, and the lower
723 whisker extends from the box to the smallest value at most 1.5xIQR. Extreme values that
724 exceed the whiskers are plotted individually as solid points.

725 **Figure S7.** Performance measures of the estimation of the among-study variance as a
726 function of the number of studies (left column), the number of replicates in the original
727 studies (middle column) and the simulated among-study variance (right column).
728 Performance was assessed using coverage (A, B, C), perceived uncertainty (width of the
729 uncertainty interval) (D, E, F), bias (G, H, I), and RMSE (J, K, L). Error bars provide the
730 95% CI for panels A-I. Please note different scales in the y-axis for bias and width of the
731 uncertainty interval. Simulation parameters: $n = 5, k = 25, \sigma_{\text{among}}^2 = 0.5$, except for the
732 cases in which that parameter was varied.

733 **Figure S8.** Performance measures of the estimation of the among-study variance as a
734 function of the number of studies (left column), the number of replicates in the original
735 studies (middle column) and the simulated among-study variance (right column).
736 Performance was assessed using coverage (A, B, C), perceived uncertainty (width of the
737 uncertainty interval) (D, E, F), bias (G, H, I), and RMSE (J, K, L). Error bars provide the
738 95% CI for panels A-I. Please note different scales in the y-axis for bias and width of the
739 uncertainty interval. Simulation parameters: $n = 5, k = 25, \sigma_{\text{among}}^2 = 2$, except for the cases
740 in which that parameter was varied.

741 **Figure S9.** Performance measures of the estimation of the among-study variance as a
742 function of the number of studies (left column), the number of replicates in the original
743 studies (middle column) and the simulated among-study variance (right column).
744 Performance was assessed using coverage (A, B, C), perceived uncertainty (width of the
745 uncertainty interval) (D, E, F), bias (G, H, I), and RMSE (J, K, L). Error bars provide the
746 95% CI for panels A-I. Please note different scales in the y-axis for bias and width of the

747 uncertainty interval. Simulation parameters: $n = 20, k = 25, \sigma_{\text{among}}^2 = 2$, except for the
748 cases in which that parameter was varied.

749 **Figure S10.** Performance measures of the estimation of the among-study variance as a
750 function of the number of studies (left column), the number of replicates in the original
751 studies (middle column) and the simulated among-study variance (right column).
752 Performance was assessed using coverage (A, B, C), perceived uncertainty (width of the
753 uncertainty interval) (D, E, F), bias (G, H, I), and RMSE (J, K, L). Error bars provide the
754 95% CI for panels A-I. Please note different scales in the y-axis for bias and width of the
755 uncertainty interval. Simulation parameters: $n = 20, k = 25, \sigma_{\text{among}}^2 = 0.5$, except for the
756 cases in which that parameter was varied.

757 **Figure S11.** Bias and RMSE from the estimation of a mean effect in 2000 replicated meta-
758 analyses as a function of the number of studies (A, B), the mean number of replicates in the
759 original studies (C, D), and the among-study variance (E, F). Simulation parameters: $n =$
760 $5, k = 25, \sigma_{\text{among}}^2 = 2$, except for the cases in which that parameter was varied. Error bars
761 provide the 95% CI for panels A-E.

762 **Figure S12.** Number of replicates yielding bad \hat{R} ($\hat{R} \geq 1.1$) for different combinations of
763 priors, true among-study variance, mean number of replicates, and number of studies.

764 **Figure S13.** Median of the posterior distribution of the among-study variance for all the
765 different priors tested, number of replicates, number of studies, and true among-study
766 variance. A) $n = 5$; B) $n = 25$. The vertical dashed line in each panel indicates the true
767 among-study variance.

768 **Figure S14.** Median of the posterior distribution of the among-study variance for the four
769 priors with the best performance (i.e., Uniform (0, 10), Uniform (0, 100), Gamma, Folded-
770 t), number of replicates, number of studies, and true among-study variance. A) $n = 5$; B)
771 $n = 25$. The vertical dashed line in each panel indicates the true among-study variance.

772 **Figure S15.** Median of the posterior distribution of the among-study variance for the four
773 priors with the best performance (i.e., Uniform (0, 10), Uniform (0, 100), Gamma, Folded-
774 t), when the number of studies was low ($k = 5$). A) $n = 5$; B) $n = 25$. The vertical dashed
775 line in each panel indicates the true among-study variance.

776