

1 **An assessment of statistical methods for non-independent data in ecological meta-analyses:**
2 **Comment**

3

4 Shinichi Nakagawa^{1*}, Alistair M. Senior², Wolfgang Viechtbauer³, and Daniel W. A. Noble⁴

5 1. Evolution & Ecology Research Centre and School of Biological, Earth and Environmental
6 Sciences, University of New South Wales, Sydney, NSW 2052, Australia.

7 2. Charles Perkins Centre and School of Life and Environmental Sciences, University of Sydney,
8 Camperdown, NSW 2006, Australia.

9 3. Department of Psychiatry and Neuropsychology, School for Mental Health and Neuroscience,
10 Faculty of Health, Medicine, and Life Sciences, Maastricht University, 6200 MD Maastricht,
11 The Netherlands

12 4. Division of Ecology and Evolution, Research School of Biology, The Australian National
13 University, Canberra, ACT, Australia

14 * correspondence: s.nakagawa@unsw.edu.au

15

16 Running title: Avoid averaging effect sizes per paper

17 Shinichi Nakagawa: 0000-0002-7765-5182

18 Alistair M. Senior: 0000-0002-7765-5182

19 Wolfgang Viechtbauer: 0000-0003-3463-4063

20 Daniel W. A. Noble: 0000-0001-9460-8743

21

22 **Key words:** meta-regression, Bayesian statistics, Satterthwaite approximation, multilevel
23 modeling, hierarchical models, degrees of freedom

24 **(Introduction)**

25 Recently, Song et al. (2020) conducted a simulation study using different methods to deal with
26 non-independence resulting from effect sizes originating from the same paper – a common
27 occurrence in ecological meta-analyses. The main methods that were of interest in their
28 simulations were: 1) a standard random-effects model used in combination with a weighted
29 average effect size for each paper (i.e., a two-step method), 2) a standard random-effects model
30 after randomly choosing one effect size per paper, 3) a multilevel (hierarchical) meta-analysis
31 model, modelling paper identity as a random factor, and 4) a meta-analysis making use of a
32 robust variance estimation method. Based on their simulation results, they recommend that meta-
33 analysts should either use the two-step method, which involves taking a weighted paper mean
34 followed by analysis with a random-effects model, or the robust variance estimation method.

35
36 Song et al.'s simulation results are an important and valuable contribution to the ecological
37 community. However, we disagree with their primary recommendation of calculating a weighted
38 average effect size for each study within a paper for two reasons. First, as we have stated
39 elsewhere (Nakagawa & Santos 2012, Noble et al. 2017), we recommend the use of multilevel
40 meta-analytic models because of improved power and the ability to answer richer biological
41 questions about the drivers underlying variation in published effects. Second, we do not
42 recommend the use of the two-step method with a weighted paper mean because other types of
43 within-study non-independence often co-occur that need to be considered but that are not
44 completely dealt with by Song et al. (2020)'s simulation. We fully agree that a robust variance
45 estimation method is useful, but from Song et al. (2020) paper it would appear to be limited in

46 applicability. However, we show that this method can easily be extended to multilevel meta-
47 analysis, making the best of both worlds.

48
49 In this Comment, we overview a previous simulation study with different conclusions to that of
50 Song et al. (2020) and put forward a strong case for why we need to make use of multilevel
51 meta-analysis in the field of ecology. We discuss how the results of this previous simulation,
52 along with our updated simulation results from Song et al. (2020), make different conclusions
53 that show multilevel meta-analysis can perform well when non-independence exists. In our
54 simulations, we demonstrate how a number of additional methods can provide solutions for any
55 increase in Type I error when fitting multilevel meta-analysis models (an issue noted by Song et
56 al., 2020).

57

58 **Similar Simulations, Different Conclusions**

59 Moeyaert et al. (2017) conducted a similar simulation study to Song et al. (2020) with some
60 minor differences. First, Moeyaert et al. (2017) did not include a condition involving randomly
61 choosing one effect per paper and used the standardized mean difference (aka Cohen's d or
62 Hedges' g) as their effect size, instead of the log response ratio used by Song et al. (2020)
63 (Hedges et al. 1999). Second, Moeyaert et al. (2017) did not model different correlations within
64 papers (they referred to papers as studies) and heteroscedasticity among papers (different
65 between-paper variances). Finally, Moeyaert et al. (2017) used PROC GLM in SAS 9.3 (SAS
66 Institute Inc, 2011-2014) while Song et al. used R's *metafor* (Viechtbauer 2010) and *robumeta*
67 (Fisher et al. 2017) for multilevel meta-analysis and robust variance estimation, respectively.

68

69 Like Song et al. (2020), Moeyaert et al. (2017) found that all of the three methods examined
70 produced unbiased estimates of the overall (meta-analytic) mean. A striking difference was that
71 in Moeyaert et al. (2017), multilevel meta-analyses performed as well as a robust variance
72 estimation method in terms of 95% confidence interval (CI) coverage. In contrast, Song et al.
73 (2020) reported consistently higher Type I error rates (i.e., greater than 5%) for multilevel meta-
74 analyses. The highest Type I error rate of multi-level meta-analysis models achieved across all
75 scenarios was about 8.2% [Mean (Median) Error Rates: 6.42% (6.42%)], which seems marginal
76 in absolute terms, but relative to the nominal rate of 5% constitutes an increase of 64%. Further,
77 Moeyaert et al. (2017) noted that when effect sizes from the same studies are not correlated, the
78 two-step method with a weighted paper mean provided confidence intervals that were too wide
79 (inefficient), which was also the case in Song et al.'s simulation. Based on their results,
80 Moeyaert et al. (2017) recommend both multilevel meta-analysis and robust variance estimation
81 methods but advised against the averaging method.

82
83 The differences in recommendations between Moeyaert et al. (2017) and Song et al. (2020) may
84 have originated from a well-known issue in linear mixed-effects models, of which multilevel
85 meta-analysis is a special type (Nakagawa & Santos 2012); that is, for linear mixed-effects
86 (multilevel) models, it is difficult to determine the appropriate degrees of freedom, which is
87 required for CI calculations. The SAS procedure used by Moeyaert et al. (2017) implements a
88 method for calculating the degrees of freedom that is more appropriate for smaller sample sizes,
89 while R's *metafor* used by Song et al. (2020) is yet to do so (at the moment, it simply sets the
90 degrees of freedom equal to the total number of effect sizes minus 1). Indeed, Song et al. (2020)
91 suspected this shortcoming by stating "this issue is addressed by adjusting the degrees of

92 freedom... ”, but did not explore the possible corrections. If this is the issue, there are several
93 solutions that already exist to correct Type I error rates toward the nominal value. Another
94 potential cause of the difference between Song et al. (2020) and Moeyaert et al. (2017) is that the
95 former modeled heteroscedasticity among papers, while the latter did not. Here, we expand Song
96 et al. (2020)’s simulations to show how currently available tools can resolve many of the issues
97 they identified without the need to resort to averaging methods. Before doing so, we would first
98 like to review the reasons for why, in the past, we have strongly recommend the use of
99 multilevel/hierarchical meta-analytic models over a method that averages multiple effect sizes
100 per paper.

101

102 **Why Multilevel/Hierarchical Models over Averaging?**

103 Nakagawa and Santos (2012) recommended the use of the following meta-analytic model for
104 datasets which can include effect sizes across different species as well as different papers:

$$105 \quad y_i = \mu + a_k + s_k + p_j + e_i + m_i, \quad (1)$$

106 where:

- 107 A. y_i is the i th effect size ($i = 1, \dots, N_{effect-size}$; the number of effect sizes),
108 B. μ is the meta-analytic mean,
109 C. a_k is the phylogenetic effect for the k th species ($k = 1, \dots, N_{species}$), which is distributed
110 with $N(0, \sigma_a^2 \mathbf{A})$, where \mathbf{A} is a correlation matrix derived from a phylogenetic tree for
111 species included in a meta-analysis,
112 D. s_k is the non-phylogenetic (species) effect for the k th species, distributed according to
113 $N(0, \sigma_s^2)$,
114 E. p_j is the j th paper effect ($j = 1, \dots, N_{paper}$), distributed according to $N(0, \sigma_p^2)$,

- 115 F. e_i is the i th effect-size specific effect, distributed according to $N(0, \sigma_e^2)$, and
116 G. m_i is the i th sampling error effect, distributed according to $N(0, \sigma_i^2)$ where σ_i^2 is the
117 sampling error variance for the i th effect size (note when sampling errors are correlated, a
118 variance-covariance matrix can replace σ_i^2 ; see below).

119

120 Although Song et al. (2012) did not mention phylogenetic non-independence (shown in Equation
121 1), this issue commonly arises in ecological meta-analysis, and is similar in manner to non-
122 independence due to effects being derived from the same source paper. What is more, it is often
123 important to appropriately take phylogeny into consideration in a meta-analysis (Chamberlain et
124 al. 2012). If we follow the logic of averaging, and we want to avoid using multilevel meta-
125 analysis, we need to average per species. Nakagawa and Santos (2012) put forward three main
126 arguments against averaging (similar arguments were independently put forward by Cheung
127 2014): 1) the potential loss of statistical power and needlessly large standard errors for the
128 overall effect, 2) the loss of information resulting from not being able to estimate within-paper
129 (within-study) variance, and 3) perhaps most importantly, not being able to estimate ecologically
130 important moderator effects given that aggregation will reduce the information content
131 dramatically (i.e. removes within-species variation in estimated effects).

132

133 Both simulation studies suggest the first argument may not apply unless correlations among
134 effect sizes are close to zero. Importantly, Song et al. (2012) discuss two scenarios where
135 dependence among effect sizes could arise: 1) “because they were observed in the same
136 experiment or may have been based on the same subjects” and 2) “even if they arose from
137 separate experiments because experiments likely share common methods, contexts, or other

138 characteristics that influence the effect size”. However, Song et al. (2020) only focused on the
139 latter scenario in their simulation not for the former. For the former, where sampling errors are
140 correlated, we need to use the following formula (Borenstein et al. 2009) to obtain a sampling
141 error variance (or a sampling standard error) to accompany a weighted mean, rather than the
142 fixed-effect model used by Song et al. (2012):

$$143 \quad \text{var}\left(\frac{1}{n}\sum_{i=1}^n y_i\right) = \left(\frac{1}{n}\right)^2 \left(\sum_{i=1}^n \sigma_i^2 + \sum_{i \neq g}^n r_{ig} \sqrt{\sigma_i^2 \sigma_g^2}\right) \quad (2),$$

144 where y_i is the i th effect size ($i = 1, \dots, n$ and $g = 1, \dots, n$, where n is the number of effect size
145 within a paper to be combined), σ_i^2 and σ_g^2 are the sampling error variances for y_i and y_g , and r_{ig}
146 is the correlation between the sampling errors of y_i and y_g (note that one can use the function,
147 *aggregate* in *metafor* to calculate a weighted mean and accompanying sampling variance as in
148 Equation 2). We believe that both types of non-independence frequently co-occur and need to be
149 accounted for. For the multilevel meta-analysis, we can model the variance-covariance matrix of
150 the sampling errors for the former type of non-independence as well as model a random effect
151 for paper, although r_{ig} is often not known and needs to be assumed (detailed in Noble et al. 2017;
152 see also, Lajeunesse, 2009; 2011. Further, beyond these two types of multilevel meta-analytic
153 models we can model different sources of non-independence (e.g., phylogenetic relatedness and
154 species relatedness not due to phylogeny; see Equation 1) simultaneously and flexibly, although
155 more data is required for more complex models (Nakagawa and Santos 2012).

156

157 The loss of information is a more serious issue, especially the loss of moderator information. The
158 high heterogeneity observed in ecological meta-analyses (Senior et al. 2016) often implies that
159 ecologists must use meta-regression models, which use moderators (or ‘predictors’) to explain

160 variation among effect sizes. In many cases, meta-regression models are likely to be more useful
161 and informative in ecology than simple meta-analytic models (Gurevitch et al. 2017). Indeed,
162 meta-regression can provide us with review- or synthesis-generated evidence which cannot be
163 obtained via single studies (Nakagawa et al. 2017). If we extend Song et al. (2020)'s
164 recommendation of not using multilevel meta-analyses to 'multilevel meta-regression', this
165 would severely limit our ability to test moderator effects. For example, it is common to obtain
166 separate effect sizes for males and females from one paper. If we aggregate these effect sizes per
167 paper then we would not be able to test sex-specific effects, which runs counter to recent
168 movements to test ubiquitous sex effects (Tannenbaum et al. 2019; Zajitschek et al. 2020).

169

170 **Solutions for Type I Errors in Multilevel Meta-analysis without the Need for Averaging**

171 Alongside the methods (referred to as Methods 1-5) used by Song et al. (2020), we explored four
172 further methods that are known to overcome the slight excess in Type I error rates observed
173 when using multi-level meta-analytic models. Our simulations reproduced the simulations by
174 Song et al. (2020) but added: 1) a simple correction to the degrees of freedom used to calculate
175 the overall effect size confidence intervals. This involved simply using one less than the total
176 number of papers instead of the typical degrees of freedom that uses the total number of effect
177 sizes (i.e., $df = \text{total papers} - 1$); 2) a Satterthwaite approximation to the effective degrees of
178 freedom, which is commonly applied in the linear mixed effect model literature (Satterthwaite,
179 1946); 3) a second cluster-robust estimation method implemented in the *clubSandwich* package
180 in R (Pustejovsky, 2020) that uses a bias-reduced linearization method (Pustejovsky and Tipton,
181 2018). The R package *clubSandwich* uses a similar robust-variance estimation method as
182 *robumeta* (Fisher et al. 2017) used in Song et al. (2020), but can be applied to *metafor's rma.mv*

183 model objects; and 4) a Bayesian modelling approach that uses an MCMC algorithm (using the
184 R package *MCMCglmm* – Hadfield, 2010), instead of restricted maximum likelihood (REML)
185 estimation, as MCMC algorithms are known to have robust coverage, albeit are slightly
186 conservative with small sample sizes (Pappalardo et al. 2020). We also explored other modelling
187 approaches, but present these four as they are simple solutions that can be easily implemented.
188 We focus exclusively on coverage / error rates given that bias was unaffected by the different
189 modelling approaches in Song et al. (2020)’s simulations. For each method (the five existing
190 methods from Song et al. 2020) plus the four new approaches we describe above, we ran 5,000
191 iterations across all the scenarios detailed in Song et al. (2020). An updated set of scripts from
192 Song et al. (2020), including a coding correction, that implements these new methods can be
193 found at [https://\(Fisher et al. 2017\).](https://(Fisher et al. 2017).)

194
195 Our new simulation results (Figure 1) show that the four proposed solutions perform quite well
196 across all the scenarios described by Song et al. (2020). The overall performance of each method
197 for each specific simulation scenario is provided in Figure 1B, which reproduces Figure 3
198 (Experiment 1 and 2) from Song et al. (2020). Overall, the simple approaches we implemented
199 corrected the excess in Type I error rates in the multi-level meta-analytic models implemented in
200 *metafor* (Figure 1A). In particular, Bayesian methods, while having inflated Type II error under
201 small sample situations, perform extremely well across a variety of conditions (Figure 1A and
202 B), with average Type I error rates converging on the 5% level but being slightly conservative
203 overall [Mean (SD) = 4.82% (0.0053)]. A Satterthwaite approximation to the effective degrees of
204 freedom also performs quite well under a variety of conditions as expected [Figure 1A & B –
205 5.02% (0.0046)]; even the simplest degrees of freedom correction that uses total papers minus

206 one performs quite well [Figure 1A & B – 5.39% (0.0059)]. Considering these results above, and
207 the ease of implementation, we recommend fitting a multilevel model with a robust variance
208 estimator because it can easily be applied to multilevel-meta-analytic models in *metafor*. Also,
209 one can certainly use Bayesian modelling, as long as the dataset is large enough (e.g., > 100
210 effect sizes). A step-by-step guide to implement both of these methods can be found at
211 https://github.com/daniel1noble/ecology_comment.

212

213 **Conclusion**

214 We appreciate the thorough simulations conducted by Song et al. (2020) in an attempt to better
215 understand the ways in which meta-analysts can overcome one of the most common challenges
216 of meta-analysis; dealing with non-independent effect sizes. While we agree with their
217 recommendation of using robust variance estimation methods (with caveats), we disagree with
218 their recommendation that averaging effect sizes within studies is a solution. While we recognise
219 that there may be times when averaging effect sizes is easier (e.g., when there are very few
220 studies with repeated effects), one most likely needs to use Equation 2 above, not the method of
221 averaging suggested by Song et al. (2020). Regardless, averaging effect sizes within studies
222 comes with a number of significant disadvantages that include: 1) not being able to control for
223 additional sources of non-independence, such as phylogenetic non-independence, which will be
224 commonplace in ecological meta-analyses and 2) not being able to understand the drivers of
225 effect size heterogeneity given that moderator information, which could be included in meta-
226 regression models, is lost. As we have shown, there are a number of very simple, and easily
227 implemented solutions to correct any inflated Type I error rates to their nominal level. Indeed,
228 even robust variance estimators can readily be incorporated into multilevel meta-analytic models,

229 which we recommend ecologists employ. Ignoring these elements prevents meta-analysts from
230 answering a richer set of biologically relevant questions about the drivers underlying effect size
231 variability. As such, we argue strongly against averaging effect sizes within a paper whenever
232 possible.

234 **Acknowledgement**

235 We thank James Pustejovsky for his advice on the robust variance estimator implemented in
236 *clubSandwich*. SN and DWAN are supported by an ARC (Australian Research Council)
237 Discovery grant (DP200100367). AMS is supported by an ARC fellowship (DE180101520).

239 **Literature Cited**

- 240 Borenstein, M., Hedges, L.V., Higgins, J.P.T. & Rothstein, H.R. 2009. Introduction to meta-
241 analysis. Wiley, Oxford.
- 242 Chamberlain, S. A., S. M. Hovick, C. J. Dibble, N. L. Rasmussen, B. G. Van Allen, B. S.
243 Maitner, J. R. Ahern, L. P. Bell-Dereske, C. L. Roy, M. Meza-Lopez, J. Carrillo, E.
244 Siemann, M. J. Lajeunesse, and K. D. Whitney. 2012. Does phylogeny matter? Assessing
245 the impact of phylogenetic information in ecological meta-analysis. *Ecology Letters*
246 **15**:627-636.
- 247 Cheung, M. W. L. 2014. Modeling Dependent Effect Sizes With Three-Level Meta-Analyses: A
248 Structural Equation Modeling Approach. *Psychological Methods* **19**:211-229.
- 249 Fisher, Z., E. Tipton, and H. Zhipeng. 2017. *robumeta*: Robust variance meta-regression. R
250 package version 2.0.

251 Hadfield, J. D. 2010. MCMC methods for multi-response Generalised Linear Mixed Models: the
252 MCMCglmm R package. *Journal of Statistical Software* **33**:1-22.

253 Lajeunesse, M. J. 2009. Meta-Analysis and the Comparative Phylogenetic Method. *American*
254 *Naturalist* **174**:369-381.

255 Lajeunesse, M. J. 2011. On the meta-analysis of response ratios for studies with correlated and
256 multi-group designs. *Ecology* **92**:2049-2055.

257 Moeyaert, M., M. Ugille, S. N. Beretvas, J. Ferron, R. Bunuan, and W. Van den Noortgate. 2017.
258 Methods for dealing with multiple outcomes in meta-analysis
259 a comparison between averaging effect sizes, robust variance estimation and multilevel meta-
260 analysis. *International Journal of Social Research Methodology* **20**:559-572.

261 Nakagawa, S., D. W. Noble, A. M. Senior, and M. Lagisz. 2017. Meta-evaluation of meta-
262 analysis: ten appraisal questions for biologists. *BMC Biology* **15**:18.

263 Nakagawa, S., and E. S. A. Santos. 2012. Methodological issues and advances in biological
264 meta-analysis. *Evolutionary Ecology* **26**:1253-1274.

265 Noble, D. W. A., M. Lagisz, E. O'Dea R, and S. Nakagawa. 2017. Nonindependence and
266 sensitivity analyses in ecological and evolutionary meta-analyses. *Molecular Ecology*
267 **26**:2410-2425.

268 Pappalardo, P., K. Ogle, E. A. Hamman, J. R. Bence, B. A. Hungate, and C. W. Osenberg. 2020.
269 Comparing traditional and Bayesian approaches to ecological meta-analysis. *Methods in*
270 *Ecology and Evolution* **11**:1286-1295.

271 Pustejovsky, J. 2017. clubSandwich: Cluster-robust (sandwich) variance estimators with small-
272 sample corrections. R package version 0.2. 3. R Found. Stat. Comput., Vienna.

273 Pustejovsky, J. E., and E. Tipton. 2018. Small-sample methods for cluster-robust variance
274 estimation and hypothesis testing in fixed effects models. *Journal of Business &*
275 *Economic Statistics* **36**:672-683.

276 Satterthwaite, F. E. 1946. An Approximate Distribution of Estimates of Variance Components.
277 *Biometrics Bulletin* **2**:110-114.

278 Song, C., S. D. Peacor, C. W. Osenberg, and J. R. Bence. 2020. An assessment of statistical
279 methods for nonindependent data in ecological meta-analyses. *Ecology* **n/a**:e03184.

280 Tannenbaum, C., R. P. Ellis, F. Eyssel, J. Zou, and L. Schiebinger. 2019. Sex and gender
281 analysis improves science and engineering. *Nature* **575**:137-146.

282 Viechtbauer, W. 2010. Conducting meta-analyses in R with the metafor package. *Journal of*
283 *Statistical Software* **36**:1-48.

284 Zajitschek, S. R., F. Zajitschek, R. Bonduriansky, R. C. Brooks, W. K. Cornwell, D. S. Falster,
285 M. Lagisz, J. Mason, A. M. Senior, D. W. A. Noble, and S. Nakagawa. 2020. Sex and
286 Power: sexual dimorphism in trait variability and its eco-evolutionary and statistical
287 implications. *eLife*.

288

289

290

291 **Figure captions**

292

293 **Figure 1** – A) Density distribution of average Type I error rates (%) across all 51 scenarios
294 simulated by Song et al. (2020). The first four methods reproduce Song et al. (2020)’s
295 simulations and include: “One” = Choosing a single effect size; “AV” = two-step method
296 that averages effects within a study; “MLM” = Multi-level meta-analytic model; “RVE”
297 = Robust variance estimation method with *robumeta*. In addition to these, we
298 implemented four new methods to correct the slight increase in Type I error rates for the
299 MLM method. These included: “CS” = club sandwich robust variance estimation;
300 “Papers_df” = degrees of freedom equal to the total number of papers minus one to adjust
301 confidence intervals from MLM; “SW_df” = Satterthwaite degrees of freedom to adjust
302 the confidence intervals of the MLM, and “Bayes” = Bayesian estimation methods. See
303 details in text. Raw error rates across all simulated scenarios described by Song et al.
304 (2020) are depicted by black points. Grey dashed line represents the nominal 5% error
305 rate. Note that the method ignoring non-independence is not included here (see Figure
306 S1). B) Average Type I error rates (%) across a sub-sample of scenarios simulated by
307 Song et al. (2020). Note that the sub-sample of simulation scenarios matches those
308 presented in Song et al. (2020) and does not include all 51 simulation scenarios presented
309 in panel A. Colors match methods described in panel A, except we also present the
310 original Method 1 as denoted in black, which completely ignores non-independence.

311

312

