

1 **Phylogenetic multilevel meta-analysis: A simulation study on the importance**
2 **of modeling the phylogeny**

3 Ozan Cinar, Department of Psychiatry and Neuropsychology, School for Mental Health and
4 Neuroscience, Faculty of Health, Medicine, and Life Sciences, Maastricht University,
5 Vijverdalseweg 1, 6226 NB Maastricht, the Netherlands, ozan.cinar@maastrichtuniversity.nl

6 Shinichi Nakagawa, Evolution & Ecology Centre, and School of Biological, Earth and
7 Environmental Sciences, BEES, University of New South Wales, Randwick NSW 2052,
8 Sydney, Australia, s.nakagawa@unsw.edu.au

9 Wolfgang Viechtbauer, Department of Psychiatry and Neuropsychology, School for Mental
10 Health and Neuroscience, Faculty of Health, Medicine, and Life Sciences, Maastricht
11 University, Vijverdalseweg 1, 6226 NB Maastricht, the Netherlands,
12 wolfgang.viechtbauer@maastrichtuniversity.nl

13 Corresponding author: Ozan Cinar, Address: Vijverdalseweg 1, 6226 NB Maastricht, the
14 Netherlands, Email: ozan.cinar@maastrichtuniversity.nl

15 Running Headline: Phylogenetic multilevel meta-analysis

Abstract

1. Meta-analyses in ecology and evolution require special attention due to certain study characteristics in these fields. First, the primary articles in these fields usually report results that are observed from studies conducted on different species, and the phylogeny among the species violates the independence assumption. Second, the primary articles frequently report multiple results which cannot be accounted for by conventional meta-analytic models. Although there is a model that accounts for these two problems in theory, its performance has not been examined extensively. In this article, we investigate the performance of this model in comparison with simpler models.

2. We conducted an extensive simulation study where data with different levels of complexities were generated and then various models were fitted to examine their performance. The models we used include the conventional random-effects and multilevel random-effects models along with more complex multilevel models that account for species-level variance with different variance components. Furthermore, we present an illustrative application of these models based on the data from a meta-analysis on size-assortative mating and comment on the results in light of the findings from the simulation study.

3. Our simulation results show that, when the phylogenetic relationships among the species are at least moderately strong, only the most complex model that decomposes the species-level variance into non-phylogenetic and phylogenetic components provides approximately unbiased estimates of the overall mean and variance components and yields confidence intervals with an approximately nominal coverage rate. On the other hand, removing the phylogenetic or non-phylogenetic component leads to biased variance component estimates and an increased risk for incorrect inferences about the overall mean. These findings are supported by the results derived from the illustrative application.

4. Based on our results, we suggest that meta-analyses in ecology and evolution should use the model that accounts for both the non-phylogenetic and phylogenetic species-level variance in addition to the multilevel structure of the data. Any attempts to simplify this

43 model, such as using only the phylogenetic variance component, may lead to erroneous
44 inferences from the data.

45

46 **Keywords:** comparative analysis, mixed-effects models, model efficiency, multilevel
47 models, phylogenetic meta-analysis, random-effects variance estimation.

1 Introduction

Meta-analysis encompasses an array of methods for synthesizing information from studies examining some phenomenon of interest and evaluating the consistency of their results (Glass, 1976; Hedges and Olkin, 1985; Cooper et al., 2009; Senior et al., 2016). Although these methods have been mostly developed in the medical and social sciences (Egger et al., 2001; Sutton and Higgins, 2008; Cooper et al., 2009), ecologists and evolutionary biologists have successfully adopted these techniques for conducting research syntheses in their respective fields (Gurevitch et al., 2001; Koricheva et al., 2013; Gurevitch et al., 2018). However, meta-analyses in ecology and evolution typically have several features that require special attention so that trustworthy evidence can be obtained.

To start, meta-analyses in these fields often incorporate data from multiple species which share an evolutionary history, known as phylogeny (Arnqvist and Wooster, 1995; Gurevitch and Hedges, 1999; Chamberlain et al., 2012). As a result, the samples (and the effect sizes obtained from these samples) are not independent which violates the independence assumption underlying conventional meta-analytic models. For example, the standard fixed- and random-effects models (Hedges and Olkin, 1985; Hedges and Vevea, 1998), often used for ecological meta-analyses (Nakagawa and Santos, 2012), assume independence among the effect sizes and therefore do not account for phylogeny (Chamberlain et al., 2012; Noble et al., 2017). This issue was first addressed by Adams (2008) and Lajeunesse (2009) who incorporated phylogenies into the fixed- and random-effects models, respectively.

Chamberlain et al. (2012) empirically investigated how the inclusion of phylogeny affects the estimate of the overall mean based on data from 30 meta-analyses in ecology and evolution. While the estimate of the overall mean did not change considerably in most cases (especially when using a random-effects model), a substantial portion of the meta-analyses, which reported significant results before, produced non-significant results when the phylogeny was incorporated into the model. Therefore, including phylogeny might be an important factor to reduce Type I error rates and to obtain an accurate reflection of the

75 uncertainty of meta-analytic estimates.

76 Although [Chamberlain et al. \(2012\)](#) is the most extensive study to date examining the
77 effects of phylogeny in meta-analysis, their work was based on available meta-analyses. To
78 investigate the issue of phylogeny more broadly, we require a simulation study to explore a
79 wider parameter space and under controlled conditions. Moreover, [Chamberlain et al. \(2012\)](#)
80 did not address the fact that ecological and evolutionary studies usually report multiple effect
81 sizes per study, which leads to another source of non-independence ([Nakagawa and Santos,](#)
82 [2012](#); [Noble et al., 2017](#)). Although past and current meta-analyses have sometimes avoided
83 this issue by selecting a single effect size from each study or by collapsing multiple effect sizes
84 into one, these procedures can lead to a severe loss of information ([Nakagawa and Santos,](#)
85 [2012](#)).

86 As an alternative, [Hadfield and Nakagawa \(2010\)](#) proposed a mixed-effects model that
87 accounts for the multilevel structure via a study-level random effect (i.e., multiple effect
88 sizes per study are nested within this random effect). In the same model, they include two
89 additional random effects to estimate the non-phylogenetic and the phylogenetic variance.
90 This way, among-species variance is decomposed into two components, the one resulting
91 from species similarities due to evolutionary history and the other from species similarities
92 due to shared ecology and other factors ([Lynch, 1991](#)). Although the model by [Hadfield](#)
93 [and Nakagawa \(2010\)](#) addresses two major statistical issues in ecological and evolutionary
94 meta-analyses, the complexity of the model poses certain challenges.

95 Partitioning the species variance into its two components is a challenging endeavor, be-
96 cause both components are modeled using random effects at the species level, with the only
97 difference that the phylogenetic component assumes that the random effects are correlated
98 according to a phylogenetic correlation matrix – which is derived from a phylogenetic tree
99 constructed based on the similarities and differences of species in terms of their (usually) ge-
100 netic (but sometimes also physical) characteristics ([Felsenstein, 2004](#)). This raises concerns
101 about the identifiability of the variance components and potential bias in their estimates,

102 issues that have also been raised outside the meta-analytic context when analyzing the data
103 of primary studies including multiple species (Paradis, 2012).

104 Moreover, the complexity of the model poses a threat to the convergence of optimization
105 algorithms (Bates et al., 2015). Accordingly, Nakagawa and Santos (2012) suggested that
106 model fitting may only be feasible with larger datasets, which would limit the applicability
107 of the model in practice. To avoid these problems, some ecological and evolutionary meta-
108 analyses have been carried out using a simplified model without the non-phylogenetic random
109 effect and that therefore accounts for species variance only via the phylogenetic component
110 (e.g., Garamszegi et al., 2012; Moore et al., 2016). However, the consequences of doing so,
111 and the performance of the more complex model, has yet to be evaluated in a simulation
112 study.

113 We therefore investigated the performance of models for conducting a phylogenetic mul-
114 tilevel meta-analysis in a comprehensive simulation study. We simulate studies that report
115 multiple effect sizes and use several models that vary in their complexity, starting from a
116 simple model (including only a random effect at the effect sizes level) to the most complex
117 model which incorporates a study-level and two among-species random effects. Further, we
118 generate specific conditions to examine the performance of the most complex model when
119 phylogenetic relationships are weak and the consequences of removing the non-phylogenetic
120 component. Finally, we present an illustrative application of these models based on the data
121 from a meta-analysis on size-assortative mating and comment on the results in light of the
122 findings from the simulation study.

123 **2 Materials and Methods**

124 **2.1 Meta-Analytic Models**

125 To conduct a meta-analysis, the phenomenon of interest (e.g., the size of a treatment effect
126 or the strength of the association between two variables) needs to be quantified in terms of

127 an effect size estimate for each study to be included in the analysis. We use the term ‘study’
 128 broadly here, as a single study may contribute multiple estimates (e.g., for multiple species,
 129 subgroups, treatments), but for the moment we assume that each study contributes a single
 130 estimate to the meta-analysis.

131 The specific effect size measure to be used in a meta-analysis depends on the phenomenon
 132 of interest and the information reported in the studies (Nakagawa and Santos, 2012). For
 133 example, raw or standardized mean differences and response ratios (Hedges et al., 1999)
 134 are typically used to quantify group differences or treatment effects based on quantitative
 135 variables, correlation coefficients (or Fisher r-to-z transformed values thereof) reflect the
 136 (linear) relationship between two variables, while (log-transformed) odds/risk ratios and
 137 risk differences (calculated from 2×2 contingency tables) indicate group differences (e.g.,
 138 treated vs. untreated, exposed vs. non-exposed) with respect to dichotomous dependent
 139 variables (e.g., cured vs. not cured, diseased vs. not diseased). For all of these measures,
 140 we can also compute the sampling variances of the estimates, that is, the variability in each
 141 estimate that would be expected under repeated sampling of new study units under identical
 142 circumstances (Nakagawa and Cuthill, 2007; Cooper et al., 2009; Borenstein et al., 2011).

143 Regardless of the specific measure used in a meta-analysis, let y_i denote the effect size
 144 estimate for the i th study (with $i = 1, \dots, N_{studies}$) and v_i the corresponding sampling
 145 variance. The most basic model that can be considered for synthesizing the estimates is the
 146 fixed-effects model, which is given by

$$y_i = \mu + e_i, \tag{1}$$

147

$$\mathbf{e} \sim N(\mathbf{0}, \mathbf{V}), \tag{2}$$

148 where μ is the overall mean, e_i is the sampling error for the i th study, \mathbf{e} is a $1 \times N_{studies}$
 149 column vector with the e_i values (which are assumed to be normally distributed with mean
 150 0 and variance v_i), $\mathbf{0}$ is a column vector of zeros, and \mathbf{V} is an $N_{studies} \times N_{studies}$ matrix with

151 the v_i values along the diagonal.

152 The fixed-effects model assumes that the included studies share a single common true
153 effect. This assumption, however, is rarely met in multi-population and multi-species meta-
154 analyses of ecology and evolution studies (Gurevitch and Hedges, 1999; Higgins et al., 2009).
155 The random-effects model addresses this potential ‘heterogeneity’ among the true effects by
156 adding a random effect corresponding to each estimate and is given by

$$y_i = \mu + u_i + e_i \quad (3)$$

157

$$\mathbf{u} \sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}_u), \quad (4)$$

158 where u_i is the random effect corresponding to the i th estimate, \mathbf{u} is a $1 \times N_{studies}$ column
159 vector with the u_i values (which are assumed to be normally distributed with mean 0 and
160 variance σ_u^2), and \mathbf{I}_u is an $N_{studies} \times N_{studies}$ identity matrix.

161 Although the models above are suitable for conducting a meta-analysis in many cir-
162 cumstances, they do not account for the multilevel structure that arises when at least some
163 studies provide multiple effect size estimates (e.g., when the same experiment was conducted
164 under varying circumstances within the same study) and they do not account for phyloge-
165 netic dependence (when studies are conducted with multiple species that differ in similarity
166 due to differences in their shared evolutionary history).

167 To address the first issue, we can use a multilevel meta-analytic model (Konstantopoulos,
168 2011; Nakagawa and Santos, 2012) which includes a random effect at the effect size level (as in
169 model 3), but which now captures variability in the true effects within studies, and a random
170 effect at the study level, which captures between-study variability. Let y_{ij} denote the j th
171 effect in the i th study (with $j = 1, \dots, N_i$, where N_i is the number of effect sizes reported in
172 the i th study), v_{ij} the corresponding sampling variance, and let $N_{total} = \sum_{i=1}^{N_{studies}} N_i$ denote

173 the total number of effects. The model is then given by

$$y_{ij} = \mu + u_{ij} + s_i + e_{ij} \quad (5)$$

174

$$\mathbf{s} \sim N(\mathbf{0}, \sigma_s^2 \mathbf{I}_s), \quad (6)$$

175 where u_{ij} is a random effect corresponding to the j th effect size in the i th study, s_i is a
 176 random effect at the study level, \mathbf{u} is now a $1 \times N_{total}$ column vector with the u_{ij} values, \mathbf{s} is
 177 a $1 \times N_{studies}$ column vector with the s_i values (which are assumed to be normally distributed
 178 with mean 0 and variance σ_s^2), and \mathbf{I}_u and \mathbf{I}_s are $N_{total} \times N_{total}$ and $N_{studies} \times N_{studies}$ identity
 179 matrices, respectively. Finally, \mathbf{e} is now a $1 \times N_{total}$ column vector with the e_{ij} values and \mathbf{V}
 180 is the corresponding (diagonal) variance-covariance matrix with dimensions $N_{total} \times N_{total}$.

181 When the effect size estimates were computed based on a set of $N_{species}$ different species,
 182 we will need an additional index. Let y_{ijk} denote the j th effect in the i th study as before, but
 183 now let $k = 1, \dots, N_{species}$ be the index that indicates for which species a particular effect
 184 size estimate was computed. Model 5 can then be extended to account for species-level
 185 variability as follows:

$$y_{ijk} = \mu + u_{ij} + s_i + n_k + e_{ij}, \quad (7)$$

186

$$\mathbf{n} \sim N(\mathbf{0}, \sigma_n^2 \mathbf{I}_n), \quad (8)$$

187 where n_k is a species-specific random effect, \mathbf{n} is a $1 \times N_{species}$ column vector with the n_k
 188 values (which are assumed to be normally distributed with mean 0 and between-species
 189 variance σ_n^2), and \mathbf{I}_n has dimensions $N_{species} \times N_{species}$. Note that n_k is a crossed random
 190 effect (e.g., [Fernández-Castilla et al., 2019](#)) and not nested within studies and we therefore
 191 do not put subscript k on u_{ij} , s_i , or e_{ij} .

192 Model 7, however, does not account for phylogeny. For this, we further extend the
 193 model by including an additional species-level random effect ([Hadfield and Nakagawa, 2010](#)),
 194 but instead of assuming independence for different species (as for the n_k values), we allow

195 the values of the random effect to be correlated according to a phylogenetic correlation
 196 matrix, which in turn is derived from a phylogenetic tree based on some model of evolution
 197 (e.g., Brownian motion) prior to the analysis (e.g., Lajeunesse, 2009; Felsenstein, 1985, 2004;
 198 Freckleton et al., 2002). The model is given by

$$y_{ijk} = \mu + u_{ij} + s_i + n_k + p_k + e_{ij}, \quad (9)$$

199

$$\mathbf{p} \sim N(\mathbf{0}, \sigma_p^2 \mathbf{A}), \quad (10)$$

200 where p_k denotes the phylogenetic random effect for the k th species, \mathbf{p} is a $1 \times N_{species}$ column
 201 vector with the p_k values (which are assumed to follow a multivariate normal distribution with
 202 mean 0 and variance-covariance matrix $\sigma_p^2 \mathbf{A}$, where \mathbf{A} is the $N_{species} \times N_{species}$ phylogenetic
 203 correlation matrix). Hence, the model includes a non-phylogenetic species-level random effect
 204 (i.e., the n_k values) to account for heterogeneity in the effects sizes due to differences between
 205 species unrelated to phylogeny (e.g., the influence of differences in the environments they
 206 live in) and a phylogenetic random effect (i.e., the p_k values) that captures dependencies in
 207 the effect sizes according to the similarities between species due to phylogenetic relatedness.

208 A concern with model 9 arises when phylogenetic relationships are weak. In that case,
 209 \mathbf{A} starts to resemble \mathbf{I}_n and hence σ_p^2 and σ_n^2 are confounded and may not be uniquely
 210 identifiable. This may lead to bias in the estimates of the variance components. This concern,
 211 or the complexity of model 9 in general, has led some researchers to adopt a simplified model
 212 in their meta-analyses where the non-phylogenetic variance component is removed. This
 213 leads to the model

$$y_{ijk} = \mu + u_{ij} + s_i + p_k + e_{ij}, \quad (11)$$

214 with all terms as explained before. Whether this simplified version is an adequate substitute
 215 for model 9 is currently unknown.

216 The models described above can be fitted within a Bayesian or likelihood framework
 217 (Hadfield and Nakagawa, 2010). For the latter, the `metafor` package (Viechtbauer, 2010)

218 for R (R Core Team, 2021) is particularly attractive as it is freely available and was written
219 specifically for the purposes of conducting meta-analyses. Maximum likelihood (ML) or
220 restricted maximum likelihood (REML) estimation can be used for model fitting (the latter
221 usually being the preferred choice; see Patterson and Thompson, 1971), providing estimates
222 of the variance components included in a particular model, the estimate of μ (i.e., $\hat{\mu}$), and its
223 standard error (i.e., $SE[\hat{\mu}]$). Likelihood ratio tests and profile likelihood confidence intervals
224 provide inferences for the variance components. An approximate 95% Wald-type confidence
225 interval for μ can be obtained with $\hat{\mu} \pm 1.96SE[\hat{\mu}]$. Analogously, $H_0: \mu = 0$ can be tested
226 by comparing $z = \hat{\mu}/SE[\hat{\mu}]$ against the critical values (i.e., ± 1.96) of a standard normal
227 distribution.

228 Although fitting the models and deriving inference from them is feasible, the consequences
229 of using the various models have not been examined systematically. We therefore conducted
230 an extensive simulation study to investigate the performance of the various model under
231 varying circumstances.

232 2.2 Simulation Setup

233 In our setup, the primary studies could provide one or multiple effect size estimates for one
234 or multiple species. We set $(N_{studies}, N_{species})$ either to (20, 40) or (50, 100) to examine the
235 difference between a smaller versus larger meta-analysis. Furthermore, we set σ_u^2 , σ_s^2 , σ_n^2 , and
236 σ_p^2 to either 0, 0.05, or 0.3 (plus an additional parameter α to be described below to either
237 0.5, 1, or 2) to define a particular condition within the simulation study. Table 1 provides
238 an overview of the 158 conditions that were studied in this manner. Note that we used
239 a ‘conditional factorization’ of the variance components to keep the number of conditions
240 manageable and to generate scenarios where one of the models described in equations 3, 5, 7,
241 and 9 corresponds to the true data generating mechanism (see Table 1). Within a particular
242 condition, the following steps were repeated 1000 times.

243 First, the number of effect sizes provided by the studies (i.e., the N_i values) were simulated

$N_{studies}$	$N_{species}$	σ_u^2	σ_s^2	σ_n^2	σ_p^2	α	Conditions	True model
20	40	0, 0.05, 0.30	0	0	0	1	3	Model 3
20	40	0.05, 0.30	0.05, 0.30	0	0	1	4	Model 5
20	40	0.05, 0.30	0.05, 0.30	0.05, 0.30	0	0.5, 1, 2	24	Model 7
20	40	0.05, 0.30	0.05, 0.30	0.05, 0.30	0.05, 0.30	0.5, 1, 2	48	Model 9
50	100	0, 0.05, 0.30	0	0	0	1	3	Model 3
50	100	0.05, 0.30	0.05, 0.30	0	0	1	4	Model 5
50	100	0.05, 0.30	0.05, 0.30	0.05, 0.30	0	0.5, 1, 2	24	Model 7
50	100	0.05, 0.30	0.05, 0.30	0.05, 0.30	0.05, 0.30	0.5, 1, 2	48	Model 9

Table 1: Overview of the conditions examined in the simulation study. The first two columns show the number of studies and species, respectively. The next four columns indicate the true values of the variance components. The α column represent the power parameter. All values were crossed within a particular row of the table. The last two columns respectively represent the number of conditions generated in each row and the model that corresponds to the true data generating mechanism for the conditions in a particular row.

244 from a right-skewed distribution, as typically observed in practice. For this, we generated
245 $N_{studies}$ random values from a Beta(1.5, 3) distribution, which were then multiplied by 39,
246 rounded to the closest integer, and increased by 1. Therefore, the number of estimates per
247 study could vary between 1 and 40 (with a mean, median, and mode of approximately 14,
248 13, and 9, respectively).

249 In the next step, we simulated the species indices (i.e., the k values) by generating N_{total}
250 random values from a Beta(2, 2) distribution, which were multiplied by $N_{species} - 1$, rounded
251 to the closest integer, and then increased by 1. Accordingly, the number of times that the
252 various species were studied followed a symmetric unimodal distribution (with mean equal
253 to $(N_{species} + 1)/2$). In order to guarantee that all species appear at least once in each meta-
254 analysis, a randomly chosen $N_{species}$ random numbers generated this way were replaced with
255 the integers from 1 to $N_{species}$.

256 Next, we generated a phylogenetic tree for the species using the `rtree()` function from
257 the R package `ape` (Paradis and Schliep, 2019), which uses a recursive random splitting
258 algorithm to simulate a phylogeny (Paradis, 2012). The branch lengths were then computed

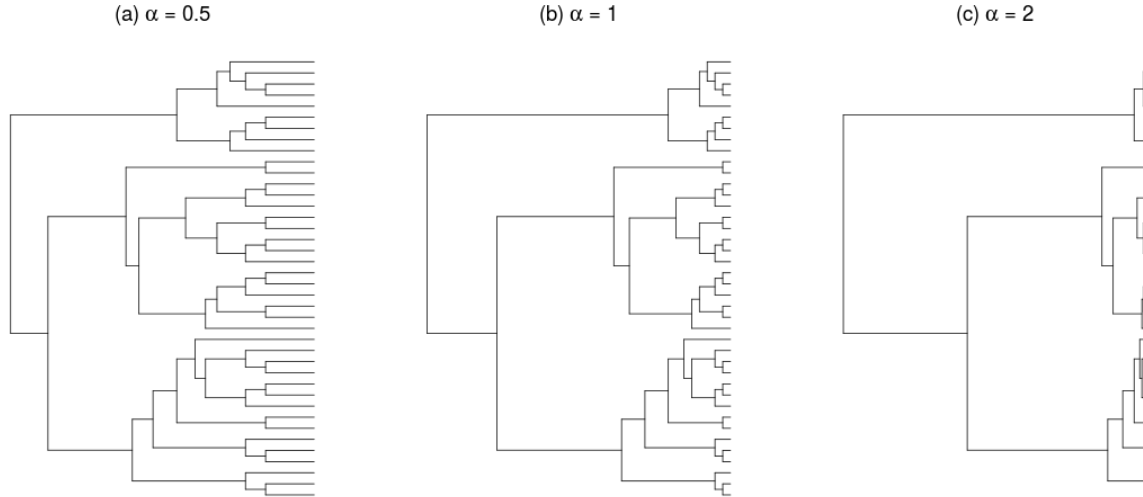


Figure 1: An example of a simulated phylogenetic tree for 40 species modified with different values of the power parameter α (i.e., 0.5, 1, and 2).

259 using the `compute.brLen()` function based on the method by Grafen (1989), using the power
 260 parameter α to adjust the ‘height’ of branch lengths at the tips of the phylogenetic tree,
 261 leading to phylogenetic relationships that are generally stronger when branches are shorter
 262 at the tips or weaker when branches are longer at the tips. Fig. 1 shows an example of such
 263 a simulated tree for 40 species modified by different α values. Finally, the correlation matrix
 264 that represents the phylogenetic relationships (matrix \mathbf{A} in equation 10) was calculated
 265 from the tree by using the `vcv()` function based on a Brownian model of evolution (i.e.,
 266 $\mathbf{A}_{k,k'} = 1 - b_{k,k'}$, where $b_{k,k'}$ is the branch length for a pair of species to their most recent
 267 common ancestor).

268 We then generated the values for the four random effects, corresponding to the variance
 269 components σ_u^2 , σ_s^2 , σ_n^2 , and σ_p^2 , either as independent draws from normal distributions for
 270 the first three components or from a multivariate normal distribution for the last one. In
 271 conditions where a variance component is equal to 0, the corresponding random effect values
 272 are then just a series of 0s of the appropriate length. To complete the data generating step,
 273 the sampling variances (i.e., the v_{ij} values) were simulated from a right-skewed Beta(2,20)

274 distribution (and hence had a value of .091 on average) which were then used to generate
275 the N_{total} sampling errors from a normal distribution with mean 0 and variance v_{ij} . We then
276 summed up the random effects and sampling errors as shown in equation 9, setting $\mu = 0$
277 without loss of generality.

278 After generating the data, we fitted the four models shown in equations 3, 5, 7, and 9,
279 using REML estimation as implemented in the `rma.mv()` function from the `metafor` package.
280 For model 3, we simply treated each estimate as a separate study (one can also think of this
281 as model 5 without the addition of the study-level random effect). For each model, we then
282 saved the estimate of μ , the variance component estimates, the bounds of the 95% Wald-type
283 confidence interval for μ , and the model fitting time. In case any one of the four models
284 did not converge within a particular iteration (with the default settings of the `rma.mv()`
285 function), the iteration was discarded and a new iteration was run to guarantee that a 1000
286 successful model fits were available for all four models.

287 After the 1000 iterations, we computed the mean of the $\hat{\mu}$ values for each model, the mean
288 of the variance component estimates, the proportion of iterations where 0 was included in
289 the confidence interval (i.e., the empirical coverage rate for μ), the mean confidence interval
290 width, the convergence rate, and the mean model fitting time. The simulation was run on a
291 workstation with an Intel Xeon E5-2630v4 processor utilizing 15 cores in parallel. Completion
292 time for the simulation was approximately 7 days (roughly 2520 core hours).

293 We generated two other sets of conditions to investigate specific questions. First, we
294 examined conditions where the phylogenetic relationships could also be weaker than in the
295 main scenarios to test the performance of model 9 under such conditions. These conditions
296 were generated by setting α to (0.1, 0.2, 0.3, 0.4, 0.5, 1, 2) when $(N_{studies}, N_{species}) = (50, 100)$,
297 the estimate- and study-level variance components were both large (0.3), and the levels of
298 the remaining variance components were factorized with values of 0.05 and 0.3 (for a total of
299 28 different conditions). Second, we compared the performance of model 9 and the simplified
300 model 11 (that leaves out the non-phylogenetic species-level random effect). For this, we set

301 $(N_{studies}, N_{species}) = (50, 100)$, $\sigma_u^2 = 0.05$, $\sigma_s^2 = 0.05$, and $\alpha = 1$, and then generated different
302 conditions by factorizing different values of only σ_n^2 and σ_p^2 , where the former was set to
303 values from 0 to 0.3 with increments of 0.05, whereas the latter was set to either 0, 0.05, or
304 0.3 (for a total of 21 different conditions).

305 **3 Results**

306 **3.1 Simulation Results**

307 [Fig. 2a](#) displays boxplots of the mean $\hat{\mu}$ values for each of the four models across the 158
308 conditions, separated by which model was the true data generating mechanism. Generally,
309 the means were clustered tightly around 0, indicating little to no bias in $\hat{\mu}$, although in a small
310 set of conditions there was some slight positive bias in the estimates of the overall mean.
311 These conditions were characterized by non-zero values for all four variance components
312 (i.e., when model 9 was the true model), $(N_{studies}, N_{species}) = (20, 40)$, a weak phylogenetic
313 relationship ($\alpha = 0.5$), and a large phylogenetic variance ($\sigma_p^2 = 0.3$).

314 In contrast to the results for the overall mean, the coverage rates of the 95% confidence
315 interval for μ differed markedly across models ([Fig. 2b](#)). For conditions where model 3 was
316 the true data generating mechanism, all models achieved coverage rates close to or slightly
317 above the nominal 95% confidence level regardless of the condition. As the other variance
318 components were introduced into the data, however, the coverage rates of models that did
319 not account for these additional sources of variability started to decrease, at times severely
320 so. Only model 9 was able to achieve rates close to the nominal level across the majority
321 of conditions, although the rates also fell somewhat below the nominal level for certain
322 conditions when all variance components were larger than zero.

323 Given that estimates of μ were relatively unbiased for all models, the closer to nominal
324 coverage rates of model 9 would be expected to be mainly a consequence of wider confidence
325 intervals (that consequently have a better chance of capturing the true value of μ). [Fig. 2c](#)

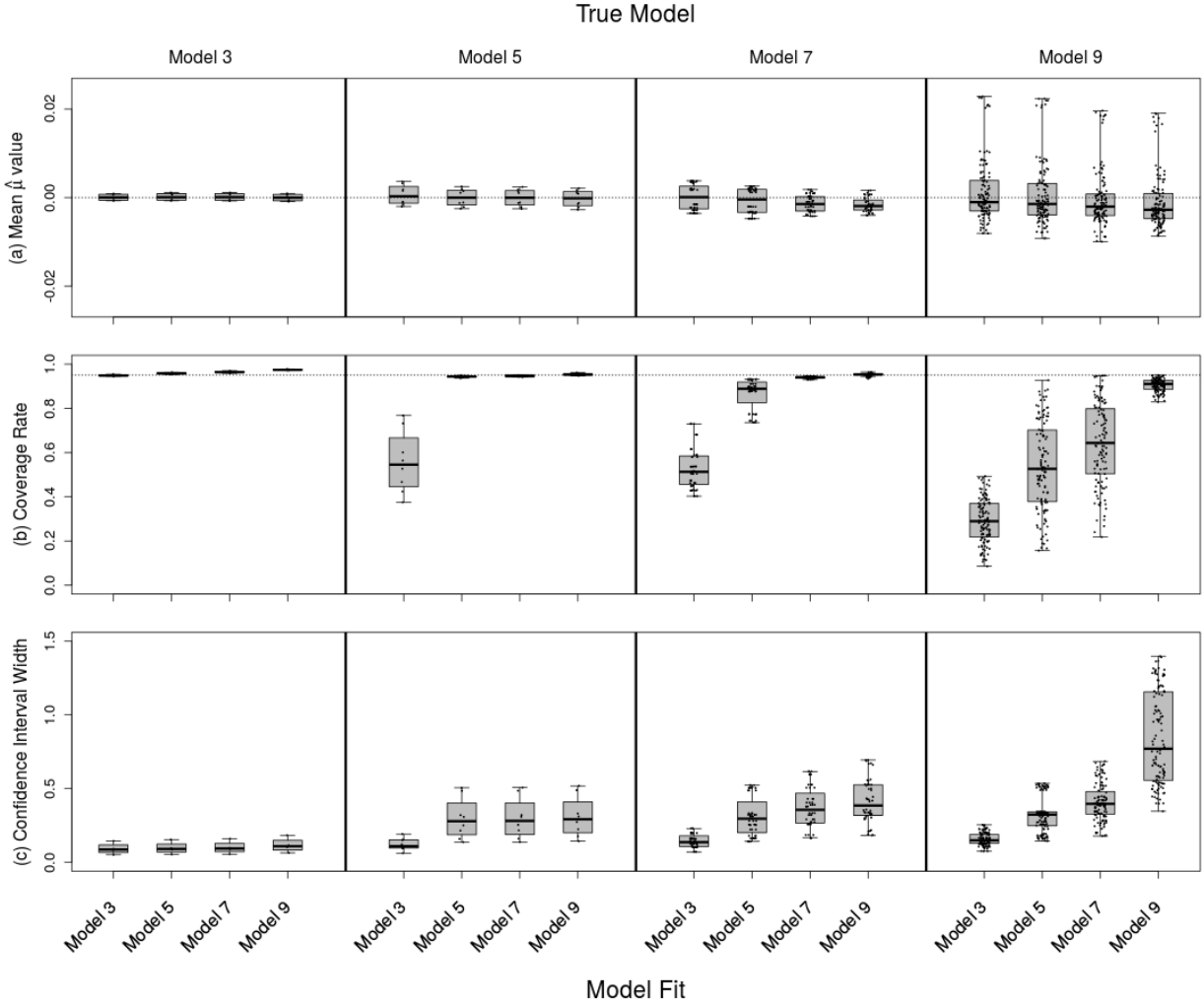


Figure 2: Boxplots based on the (a) mean $\hat{\mu}$ values, (b) coverage rates of the 95% confidence interval for μ , and (c) mean confidence interval widths for each of the four models across the 158 conditions, separated by which model was the true data generating mechanism.

326 confirms this, showing the mean confidence interval widths for the various models across
 327 the various conditions. However, what is particularly noteworthy is that the use of model
 328 9 under conditions where actually a simpler model is the true data generating mechanism
 329 only leads to a relatively minor increase in the mean interval width.

330 Fig. 3 displays the bias in the variance component estimates of model 9 under the 28
 331 different conditions generated by varying α , σ_n^2 , and σ_p^2 (while holding σ_u^2 and σ_s^2 constant
 332 at 0.3). The results show no bias in the estimates of σ_u^2 and σ_s^2 . Furthermore, the model

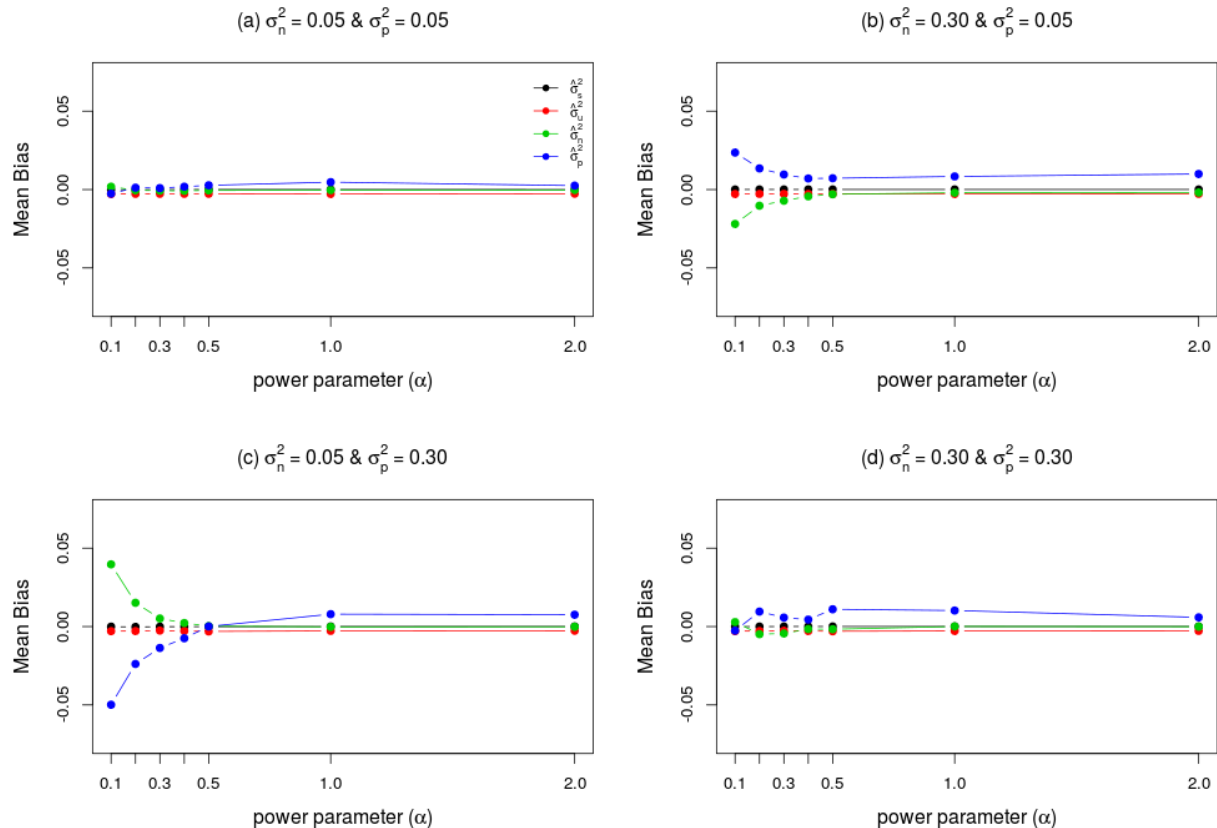


Figure 3: Mean bias of the variance component estimates of model 9 under different combinations of the power parameter (α) and the non-phylogenetic and phylogenetic variance components (σ_n^2 and σ_p^2 , respectively).

333 is able to estimate σ_n^2 and σ_p^2 with little to no bias, except when the strength of the phylo-
 334 genetic relationships decreased. As expected, under such conditions, the model struggles to
 335 provide unbiased estimates of the non-phylogenetic and phylogenetic species-level variance
 336 components (especially when $\sigma_n^2 \neq \sigma_p^2$).

337 Fig. 4a shows the coverage rates of the confidence interval for μ for models 9 and 11 as
 338 the size of the non-phylogenetic species-level variance component (i.e., σ_n^2) was increased.
 339 While model 9 provided rates close to or somewhat below the nominal level, the rates for
 340 model 11 were often equal to 100% and hence the confidence interval tended to be too wide.
 341 Furthermore, Fig. 4b demonstrates that the bias in the phylogenetic variance component of
 342 model 11 inflated rapidly as the value of σ_n^2 increased (the value of σ_p^2 had no noteworthy

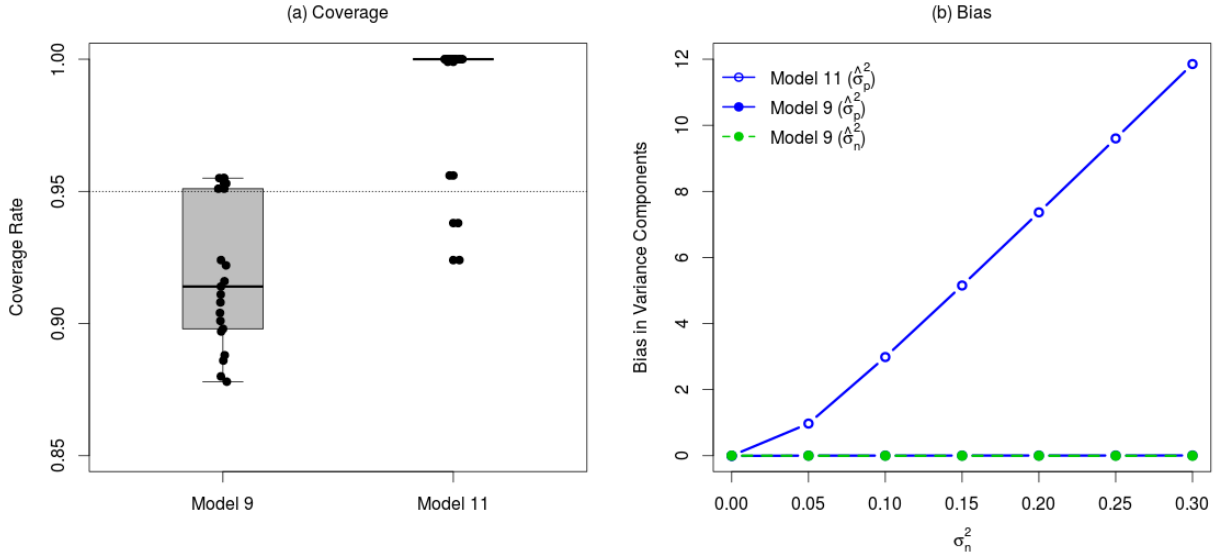


Figure 4: Comparison of models 9 and 11 as the size of the non-phylogenetic species-level variance component (i.e., σ_n^2) was systematically increased. (a) Coverage rates of the 95% confidence intervals for μ , (b) bias in the non-phylogenetic and phylogenetic variance components.

343 influence on the bias and hence we averaged these results over the three possible values of
 344 σ_p^2). In contrast, model 9 estimated these two variance components essentially without bias
 345 under these scenarios.

346 Model fitting times differed between the various models (Table 2), with model 9 requiring
 347 the most amount of time on average, regardless of the true data generating mechanism.
 348 The most challenging conditions for the more complex models were those scenarios where
 349 model 3 corresponded to the true data generating mechanism. In this case, a single fit of
 350 model 9 took around 26 seconds on average when $(N_{studies}, N_{species}) = (50, 100)$. In these
 351 conditions, convergence rates were also the lowest, although even model 9 then converged in
 352 approximately 96% of the iterations.

353 3.2 Illustrative Example

354 We use the data from the meta-analysis by Rios Moura et al. (2021) on size-assortative mating
 355 (SAM) to illustrate an application of the models. Each study included in the meta-analysis

(a) ($N_{studies}, N_{Species}$) = (20, 40)					(b) ($N_{studies}, N_{Species}$) = (50, 100)				
Model Fit	True Model				Model Fit	True Model			
	Model 3	Model 5	Model 7	Model 9		Model 3	Model 5	Model 7	Model 9
Model 3	0.939 (100.00%)	0.668 (100.00%)	0.700 (100.00%)	0.700 (100.00%)	Model 3	1.589 (100.00%)	1.313 (100.00%)	1.307 (100.00%)	1.294 (100.00%)
Model 5	2.653 (99.81%)	1.104 (100.00%)	1.151 (100.00%)	1.162 (100.00%)	Model 5	3.934 (99.78%)	1.986 (100.00%)	1.999 (100.00%)	1.959 (100.00%)
Model 7	2.484 (97.53%)	1.876 (100.00%)	0.868 (100.00%)	0.858 (100.00%)	Model 7	19.823 (96.86%)	14.752 (100.00%)	7.364 (100.00%)	7.393 (100.00%)
Model 9	3.316 (96.56%)	3.053 (99.73%)	2.288 (99.99%)	1.463 (99.99%)	Model 9	25.980 (95.63%)	23.540 (99.60%)	18.641 (100.00%)	11.005 (100.00%)

Table 2: Average model fitting times in seconds and convergence rates (in parentheses) of all models under the different data generating mechanisms.

356 provided one or multiple correlation coefficients describing the similarity in some measure of
357 body size in mating couples. For the analysis, the correlation coefficients were transformed
358 with Fisher’s r-to-z transformation. We focus here on the estimate of the overall mean
359 (transformed) correlation coefficient, leaving aside the issue of differences between studies
360 where correlations were computed with or without pooling of data across different timepoints
361 or areas (i.e., temporal/spatial pooling). Also, using the method by Grafen (1989), we turned
362 the phylogenetic tree used by Rios Moura et al. (2021) into an ultrametric tree before fitting
363 models 9 and 11, to bring these analyses more in line with how our simulation study was
364 conducted. The dataset includes 1828 effect size estimates (i.e., transformed correlations)
365 collected from 457 studies and 341 species.

366 Table 3 presents the results obtained from each model. Interestingly, the estimate of the
367 overall mean tended to be somewhat larger in the more complex models, although differences
368 between models 7, 9, and 11 were relatively small. More importantly, we see a substantial
369 increase in the standard error of the estimated overall mean for the more complex models.
370 As a result, the confidence intervals become wider, the values of the test statistics smaller,
371 while the respective p -values increase. Although each model suggests that the overall mean
372 significantly differs from 0 (at the conventional 0.05 level of significance), the p -value for

	$\hat{\mu}$	$SE[\hat{\mu}]$	95% CI	Z	p	$\hat{\sigma}_u^2$	$\hat{\sigma}_s^2$	$\hat{\sigma}_n^2$	$\hat{\sigma}_p^2$	AIC
Model 3	0.24	0.007	0.23, 0.25	34.15	<0.0001	0.0641	–	–	–	1082.8
Model 5	0.30	0.015	0.27, 0.33	20.42	<0.0001	0.0149	0.0806	–	–	429.0
Model 7	0.34	0.020	0.30, 0.38	17.37	<0.0001	0.0143	0.0195	0.0815	–	386.3
Model 9	0.37	0.130	0.11, 0.62	2.83	0.0046	0.0145	0.0192	0.0555	0.0512	344.7
Model 11	0.36	0.172	0.02, 0.70	2.07	0.0383	0.0149	0.0557	–	0.0914	367.2

Table 3: Results derived from fitting the various models to the example dataset. The first five columns show the estimated overall mean, its standard error, the 95% confidence interval, the test statistic, and the p -value for testing $H_0: \mu = 0$, respectively. The next four columns show the estimates of the variance components in the respective models. The last column shows the Akaike Information Criteria (AIC) values.

373 model 11 was approaching the rejection threshold.

374 The estimates of the variance components also show some interesting patterns. While the
375 simple random-effects model 3 cannot distinguish between different sources of variability and
376 attributes all of the heterogeneity to differences between the individual effect size estimates,
377 model 5 suggests that the variance in the effects is more related to differences between
378 studies than particular estimates within studies. However, once species-level variability is
379 considered, it becomes apparent that this is actually the dominant source of heterogeneity.
380 Moreover, model 9 shows that this variability is to approximately equal parts attributable to
381 non-phylogenetic and phylogenetic species-level differences. In contrast, when ignoring the
382 non-phylogenetic variance component in the simplified model 11, part of the variance from
383 that component is forced back into the study-level variance component. Furthermore, $\hat{\sigma}_p^2$ in
384 the simplified model is substantially inflated compared to model 9 which may be an example
385 of the inflation in this component when σ_n^2 is excluded (see Fig. 4b). Based on these findings
386 and the Akaike Information Criteria (AIC) values of the various models, we would strongly
387 favor model 9 in this comparison, illustrating that both non-phylogenetic and phylogenetic
388 variance components should be considered in the analysis.

4 Discussion

Meta-analyses in the fields of ecology and evolution typically need to address the fact that multiple effect size estimates can be extracted from at least some of the studies and that the estimates are based on various species that are related to each other due to their shared evolutionary history. In this paper, we investigated the performance of the phylogenetic multilevel meta-analytic model by Hadfield and Nakagawa (2010) and Nakagawa and Santos (2012) that captures these intricacies along with some simpler models. Despite the concerns raised in the introduction, the model can successfully estimate the overall mean and its uncertainty. It also provides approximately unbiased estimates of all variance components, including the non-phylogenetic and phylogenetic species-level variances, as long as there are at least moderately strong phylogenetic relationships among the species. In addition, despite its complexity, the model does not appear to suffer from convergence problems and model fitting does not require excessive computational times.

4.1 Estimating the Overall Mean and its Uncertainty

Not only the phylogenetic multilevel meta-analytic model, but also the simpler models that leave out certain variance components provide essentially unbiased estimates of the overall mean, regardless of the nature of the true model that underlies the data (Fig. 2a). However, the uncertainty in the overall mean will only be estimated accurately when the fitted model includes the variance components that contribute to the heterogeneity and the dependencies among the underlying true effects. Fitting underspecified models typically led to severe undercoverage of the confidence interval for the overall mean and hence anticonservative inferences. In fact, subtracting the coverage rates shown in Fig. 2b from 1 yields the Type I error rates for the test of the overall mean, which could go as high as 91% when using a simple random-effects model that ignores the multilevel structure and the species-level variance components.

414 These findings are in line with those by [Chamberlain et al. \(2012\)](#), who demonstrated,
415 based on 30 published meta-analyses, that the inclusion of phylogeny into a random-effects
416 model usually only led to minor changes in the pooled effect size, but had a more substantial
417 impact on the statistical significance of the finding (turning significant findings into non-
418 significant ones in the majority of cases where changes occurred).

419 Our findings can also be used to alleviate concerns with using the phylogenetic multilevel
420 meta-analytic model when it is actually an overspecified model (i.e., when the actual data
421 generating mechanism is simpler). In those cases, the mean confidence interval width of the
422 model was just barely wider than that of the simpler models, indicating little to no loss in
423 efficiency by fitting an overly complex model ([Fig. 2c](#)). The superfluous variance components
424 then converge towards 0 (or close to it), which appears to be slightly more challenging for
425 the optimization algorithm, leading to longer model fitting times and occasional convergence
426 problems, but not to any worrisome degree ([Table 2](#)). Moreover, in practice, for any particu-
427 lar dataset, convergence problems can typically be resolved by selecting a different optimizer
428 or making changes to the settings for the optimization routine, so the convergence rates as
429 given only apply to the default settings.

430 At the same time, we should point out that the coverage rate of the model did fall slightly
431 below the nominal 95% level in the majority of conditions when all variance components
432 were in fact non-zero (see [Fig. 2b](#), rightmost panel). This undercoverage stems from using
433 an overly simple Wald-type confidence interval using critical values based on a standard
434 normal distribution that ignores the uncertainty in the estimates of the variance components
435 (especially in the study and the two species-level components when $N_{studies}$ and $N_{species}$ are
436 low). A similar issue, but for a simpler model with only between- and within-study variance
437 components (i.e., model 5 in our simulation) was also recently pointed out by [Song et al.](#)
438 ([2020](#)). Improved methods based on the t-distribution, with various approximations for
439 the degrees of freedom, have been proposed and studied extensively in the context of the
440 standard random-effects model (e.g., [Sanchez-Meca and Marin-Martinez, 2008](#)) and mixed-

441 effects models in general (e.g., [Luke, 2017](#)), but these methods have not been generalized
 442 to the present context. As a simple approximation, using the smaller of $N_{studies} - 1$ and
 443 $N_{species} - 1$ as the degrees of freedom for a confidence interval based on a t-distribution is
 444 likely to bring the coverage rate quite close to the nominal rates in the majority of conditions.

445 4.2 Including and Testing the Phylogenetic Effect

446 Phylogenies play a central role in the context of phylogenetic comparative studies ([Freckleton
 447 et al., 2002](#); [Blomberg et al., 2003](#); [Ives et al., 2007](#)). An important step in such studies is
 448 testing the significance of the ‘phylogenetic signal’ in some trait of interest. This test is
 449 often performed through a statistic such as λ ([Pagel, 1999](#)) or K ([Blomberg et al., 2003](#)).
 450 Although model 9 does not parameterize the phylogenetic effect in this manner, one can
 451 derive information from its output that shows its relationship to the λ statistic. In particular,
 452 Pagel’s λ is a multiplicative factor that is applied to the off-diagonal values of the correlation
 453 matrix that represents the phylogenetic relationships (i.e., the \mathbf{A} matrix). For example, the
 454 variance-covariance matrix for three species would be given by

$$\sigma^2 \begin{bmatrix} 1 & \lambda a_{12} & \lambda a_{13} \\ & 1 & \lambda a_{23} \\ & & 1 \end{bmatrix}$$

455 while the decomposition of the species-level heterogeneity in model 9 implies the variance-
 456 covariance matrix

$$\sigma_n^2 \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} + \sigma_p^2 \begin{bmatrix} 1 & a_{12} & a_{13} \\ & 1 & a_{23} \\ & & 1 \end{bmatrix} = (\sigma_n^2 + \sigma_p^2) \begin{bmatrix} 1 & \left(\frac{\sigma_p^2}{\sigma_n^2 + \sigma_p^2}\right) a_{12} & \left(\frac{\sigma_p^2}{\sigma_n^2 + \sigma_p^2}\right) a_{13} \\ & 1 & \left(\frac{\sigma_p^2}{\sigma_n^2 + \sigma_p^2}\right) a_{23} \\ & & 1 \end{bmatrix}$$

457 and hence $\sigma^2 = \sigma_n^2 + \sigma_p^2$ and $\lambda = \sigma_p^2 / (\sigma_n^2 + \sigma_p^2)$ (see also Lynch, 1991; Freckleton et al., 2002).
458 Hence, $\sigma_p^2 / (\sigma_n^2 + \sigma_p^2)$ indicates the degree of the phylogenetic signal in the overall variance
459 sourced from the species. A likelihood ratio test of $H_0: \sigma_p^2 = 0$ can be easily performed by
460 comparing $X^2 = -2(ll_7 - ll_9)$ against a chi-squared distribution with one degree of freedom,
461 where ll_7 and ll_9 are the (restricted) log likelihoods of models 7 and 9, respectively.

462 4.3 Estimating the Non-Phylogenetic and Phylogenetic Variance

463 Given the informative nature of these two variance components, it is essential to estimate
464 their true values accurately. We found that model 9 was usually able to estimate these compo-
465 nents unbiasedly, but should note that the model struggles to separate the non-phylogenetic
466 and phylogenetic species effects when phylogenetic relationships are weak. In essence, the
467 two sources of variability then start to collapse into one, with a total variance of $\sigma_n^2 + \sigma_p^2$. This
468 total variance is then distributed in approximately equal parts into the two estimates, which
469 explains the apparent low bias when (coincidentally) $\sigma_n^2 = \sigma_p^2$ (Fig. 3a and d). However,
470 when $\sigma_n^2 \neq \sigma_p^2$, the bias in the two estimates becomes apparent (Fig. 3b and c). Therefore,
471 we would caution against the use of model 9 when phylogenetic relationships are weak. As a
472 rough guideline, for $\alpha = 0.5$, the mean correlation in the **A** matrix (excluding the diagonal)
473 is around 0.2 and hence a lower mean correlation would call into question the trustworthiness
474 of the estimates of σ_n^2 and σ_p^2 .

475 Some meta-analyses in ecology and evolution have used model 11 to reduce model com-
476 plexity (e.g., Garamszegi et al., 2012; Moore et al., 2016). Our results indicate that this
477 approach cannot be recommended. As we increased the value of σ_n^2 , the bias in the phyloge-
478 netic variance component inflated massively in this simplified model (Fig. 4b). As a result,
479 the relevance of the phylogeny could be greatly overestimated. In addition, the confidence
480 interval for the overall mean then becomes extremely conservative with coverage rates at
481 or very close to 100%. This in turn implies a loss of efficiency for estimating the overall
482 mean and a loss of power for testing $H_0: \mu = 0$. The illustrative example also shows this

483 phenomenon.

484 4.4 Caveats and Conclusions

485 For the simulation study, we used a ‘generic’ effect size measure, that is, we directly simulated
486 the sampling errors from a normal distribution and treated the sampling variances (i.e.,
487 the v_{ij} values) as known. These conditions only apply asymptotically to measures typically
488 used in practice (e.g., standardized mean differences, response ratios, correlation coefficients,
489 risk/odds ratios). The present results therefore reflect the performance of the various models
490 under idealized conditions (i.e., when the sample sizes of the individual studies are sufficiently
491 large, such that the sampling distributions of the estimates are indeed approximately normal
492 and any inaccuracies in the estimated sampling variances are negligible). The advantage of
493 using a generic measure is that we were able to identify problems that are inherent to certain
494 models and not (potentially) a consequence of violations to the model assumptions. On the
495 other hand, it remains to be determined how well the phylogenetic multilevel model performs
496 when the effect sizes are generated based on the exact distributional assumptions underlying
497 specific measures.

498 Therefore, at least for the moment, the present results suggest that model 9 is the most
499 appropriate tool for conducting a multi-species meta-analysis in ecology and evolution. For
500 the vast majority of conditions examined, it provides approximately unbiased estimates of
501 the variance components and the overall mean and a confidence interval for the latter with a
502 close to nominal coverage rate. Therefore, we recommend that meta-analysts in ecology and
503 evolution use the phylogenetic multilevel model as the de facto standard when analyzing
504 multi-species datasets.

505

506 **Conflict of interest statement:** The authors declare that they have no competing
507 interests.

508

509 **Author contributions:** SN provided contextual and literature review support, WV
510 provided the code to run the simulation, all authors contributed to the manuscript.

511

512 **Data accessibility statement:** No new data were used in this study. The material to
513 reproduce the results are available at: <https://osf.io/ms8eq/>.

514 References

515 Dean C. Adams. Phylogenetic meta-analysis. *Evolution*, 62(3):567–572, 2008. doi: 10.1111/
516 j.1558-5646.2007.00314.x.

517 Göran Arnqvist and David Wooster. Meta-analysis: Synthesizing research findings in
518 ecology and evolution. *Trends in Ecology & Evolution*, 10(6):236–240, 1995. doi:
519 10.1016/s0169-5347(00)89073-4.

520 Douglas Bates, Reinhold Kliegl, Shravan Vasishth, and Harald Baayen. Parsimonious mixed
521 models. *arXiv preprint arXiv:1506.04967*, 2015.

522 Simon P Blomberg, Theodore Garland Jr, and Anthony R Ives. Testing for phylogenetic
523 signal in comparative data: Behavioral traits are more labile. *Evolution*, 57(4):717–745,
524 2003. doi: 10.1554/0014-3820(2003)057[0717:TFPSIC]2.0.CO;2.

525 Michael Borenstein, Larry V Hedges, Julian PT Higgins, and Hannah R Rothstein. *Intro-*
526 *duction to meta-analysis*. Wiley, Chichester, UK, 2011.

527 Scott A Chamberlain, Stephen M Hovick, Christopher J Dibble, Nick L Rasmussen, Ben-
528 jamin G Van Allen, Brian S Maitner, Jeffrey R Ahern, Lukas P Bell-Dereske, Christopher L
529 Roy, Maria Meza-Lopez, et al. Does phylogeny matter? Assessing the impact of phyloge-
530 netic information in ecological meta-analysis. *Ecology Letters*, 15(6):627–636, 2012. doi:
531 10.1111/j.1461-0248.2012.01776.x.

- 532 Harris Cooper, Larry Vernon Hedges, and Jeffrey C Valentine. *The handbook of research*
533 *synthesis and meta-analysis*. Russell Sage Foundation, New York, 2nd edition, 2009.
- 534 Matthias Egger, George Davey-Smith, and Douglas Altman. *Systematic reviews in health*
535 *care: Meta-analysis in context*. Wiley, London, 2nd edition, 2001.
- 536 J. Felsenstein. *Inferring phylogenies*. Sinauer Associates, Sunderland, MA, 2nd edition, 2004.
- 537 Joseph Felsenstein. Phylogenies and the comparative method. *The American Naturalist*, 125
538 (1):1–15, 1985. doi: 10.1086/284325.
- 539 B. Fernández-Castilla, M. Maes, L. Declercq, L. Jamshidi, S. N. Beretvas, P. Onghena, and
540 W. Van den Noortgate. A demonstration and evaluation of the use of cross-classified
541 random-effects models for meta-analysis. *Behavior Research Methods*, 51(3):1286–1304,
542 2019. doi: 10.3758/s13428-018-1063-2.
- 543 Rob P Freckleton, Paul H Harvey, and Mark Pagel. Phylogenetic analysis and comparative
544 data: A test and review of evidence. *The American Naturalist*, 160(6):712–726, 2002. doi:
545 10.1086/343873.
- 546 László Zsolt Garamszegi, Gábor Markó, and Gábor Herczeg. A meta-analysis of correlated
547 behaviours with implications for behavioural syndromes: Mean effect size, publication
548 bias, phylogenetic effects and the role of mediator variables. *Evolutionary Ecology*, 26(5):
549 1213–1235, 2012. doi: 10.1007/s10682-012-9589-8.
- 550 Gene V Glass. Primary, secondary, and meta-analysis of research. *Educational Researcher*,
551 5(10):3–8, 1976. doi: 10.3102/0013189x005010003.
- 552 Alan Grafen. The phylogenetic regression. *Philosophical Transactions of the Royal Society*
553 *of London, Series B*, 326(1233):119–157, 1989.
- 554 Jessica Gurevitch and Larry V Hedges. Statistical issues in ecological meta-analyses. *Ecology*,
555 80(4):1142–1149, 1999. doi: 10.1890/0012-9658(1999)080[1142:siiema]2.0.co;2.

- 556 Jessica Gurevitch, Peter S Curtis, and Michael H Jones. Meta-analysis in ecology. *Advances*
557 *in Ecological Research*, 32:199–247, 2001.
- 558 Jessica Gurevitch, Julia Koricheva, Shinichi Nakagawa, and Gavin Stewart. Meta-analysis
559 and the science of research synthesis. *Nature*, 555(7695):175, 2018. doi: 10.1038/
560 nature25753.
- 561 Jarrod D Hadfield and Shinichi Nakagawa. General quantitative genetic methods for
562 comparative biology: Phylogenies, taxonomies and multi-trait models for continuous
563 and categorical characters. *Journal of Evolutionary Biology*, 23(3):494–508, 2010. doi:
564 10.1111/j.1420-9101.2009.01915.x.
- 565 Larry Hedges and Ingram Olkin. *Statistical models for meta-analysis*. Academic Press, New
566 York, 1985.
- 567 Larry V Hedges and Jack L Vevea. Fixed- and random-effects models in meta-analysis.
568 *Psychological Methods*, 3(4):486–504, 1998. doi: 10.1037/1082-989x.3.4.486.
- 569 Larry V Hedges, Jessica Gurevitch, and Peter S Curtis. The meta-analysis of response ratios
570 in experimental ecology. *Ecology*, 80(4):1150–1156, 1999. doi: 10.1890/0012-9658(1999)
571 080[1150:tmaorr]2.0.co;2.
- 572 Julian PT Higgins, Simon G Thompson, and David J Spiegelhalter. A re-evaluation of
573 random-effects meta-analysis. *Journal of the Royal Statistical Society: Series A*, 172(1):
574 137–159, 2009. doi: 10.1111/j.1467-985x.2008.00552.x.
- 575 Anthony R Ives, Peter E Midford, and Theodore Garland Jr. Within-species variation and
576 measurement error in phylogenetic comparative methods. *Systematic Biology*, 56(2):252–
577 270, 2007. doi: 10.1080/10635150701313830.
- 578 S. Konstantopoulos. Fixed effects and variance components estimation in three-level meta-
579 analysis. *Research Synthesis Methods*, 2(1):61–76, 2011. doi: 10.1002/jrsm.35.

580 Julia Koricheva, Jessica Gurevitch, and Kerrie Mengersen, editors. *Handbook of meta-*
581 *analysis in ecology and evolution*. Princeton University Press, Princeton, NJ, 2013.

582 Marc J Lajeunesse. Meta-analysis and the comparative phylogenetic method. *The American*
583 *Naturalist*, 174(3):369–381, 2009. doi: 10.2307/40306065.

584 Steven G Luke. Evaluating significance in linear mixed-effects models in R. *Behavior Re-*
585 *search Methods*, 49(4):1494–1502, 2017. doi: 10.3758/s13428-016-0809-y.

586 Michael Lynch. Methods for the analysis of comparative data in evolutionary biology. *Evo-*
587 *lution*, 45(5):1065–1080, 1991. doi: 10.1111/j.1558-5646.1991.tb04375.x.

588 Fhionna R Moore, David M Shuker, and Liam Dougherty. Stress and sexual signaling:
589 A systematic review and meta-analysis. *Behavioral Ecology*, 27(2):363–371, 2016. doi:
590 10.1093/beheco/arv195.

591 Shinichi Nakagawa and Innes C Cuthill. Effect size, confidence interval and statistical sig-
592 nificance: A practical guide for biologists. *Biological Reviews*, 82(4):591–605, 2007. doi:
593 10.1111/j.1469-185x.2007.00027.x.

594 Shinichi Nakagawa and Eduardo SA Santos. Methodological issues and advances in bi-
595 ological meta-analysis. *Evolutionary Ecology*, 26(5):1253–1274, 2012. doi: 10.1007/
596 s10682-012-9555-5.

597 Daniel WA Noble, Malgorzata Lagisz, Rose E O’dea, and Shinichi Nakagawa. Nonindep-
598 endence and sensitivity analyses in ecological and evolutionary meta-analyses. *Molecular*
599 *Ecology*, 26(9):2410–2425, 2017. doi: 10.1111/mec.14031.

600 Mark Pagel. Inferring the historical patterns of biological evolution. *Nature*, 401(6756):877,
601 1999. doi: 10.1038/44766.

602 Emmanuel Paradis. *Analysis of phylogenetics and evolution with R*. Springer, New York,
603 2nd edition, 2012.

604 Emmanuel Paradis and Klaus Schliep. ape 5.0: An environment for modern phylogenetic
605 ics and evolutionary analyses in R. *Bioinformatics*, 35:526–528, 2019. doi: 10.1093/
606 bioinformatics/bty633.

607 H. D. Patterson and R. Thompson. Recovery of inter-block information when block sizes are
608 unequal. *Biometrika*, 58(3):545–554, 1971. doi: 10.1093/biomet/58.3.545.

609 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation
610 for Statistical Computing, Vienna, Austria, 2021. URL <https://www.R-project.org/>.

611 Rafael Rios Moura, Marcelo Oliveira Gonzaga, Nelson Silva Pinto, João Vasconcellos-Neto,
612 and Gustavo S Requena. Assortative mating in space and time: Patterns and biases.
613 *Ecology Letters*, 24:1089–1102, 2021. doi: 10.1111/ele.13690.

614 J. Sanchez-Meca and F. Marin-Martinez. Confidence intervals for the overall effect size in
615 random-effects meta-analysis. *Psychological Methods*, 13(1):31–48, 2008. doi: 10.1037/
616 1082-989x.13.1.31.

617 Alistair M Senior, Catherine E Grueber, Tsukushi Kamiya, Malgorzata Lagisz, Katie
618 O’Dwyer, Eduardo SA Santos, and Shinichi Nakagawa. Heterogeneity in ecological and
619 evolutionary meta-analyses: Its magnitude and implications. *Ecology*, 97(12):3293–3299,
620 2016. doi: 10.1002/ecy.1591.

621 C. Song, S. D. Peacor, C. W. Osenberg, and J. R. Bence. An assessment of statistical
622 methods for nonindependent data in ecological meta-analyses. *Ecology*, 101(12):e03184,
623 2020. doi: 10.1002/ecy.3184.

624 Alexander J Sutton and Julian PT Higgins. Recent developments in meta-analysis. *Statistics*
625 *in Medicine*, 27(5):625–650, 2008. doi: 10.1002/sim.2934.

626 Wolfgang Viechtbauer. Conducting meta-analyses in R with the metafor package. *Journal*

627 *of Statistical Software*, 36(3):1–48, 2010. doi: 10.18637/jss.v036.i03. URL <http://www.jstatsoft.org/v36/i03/>.