

# Comparing ecological and evolutionary variability within datasets

1 Raphaël Royauté <sup>a,b\*</sup> and Ned A. Dochtermann <sup>a</sup>

2 <sup>a</sup> Department of Biological Sciences; North Dakota State University

3 <sup>b</sup>Current address: Behavioural Ecology, Department of Biology, Ludwig-Maximilians

4 University of Munich, Planegg-Martinsried, Germany

5 \* corresponding author: [raphael.royaute@gmail.com](mailto:raphael.royaute@gmail.com)

6

7 Running Head: Comparing variation within datasets

8 **ABSTRACT**

- 9 1. Variance ratios—including heritability, repeatability, and individual resource  
10 specialization—are an integral part of evolutionary ecology. Understanding how  
11 evolutionary and ecological processes differs among populations and environments,  
12 can require the comparison of these ratios across groups.
- 13 2. Inference based on comparisons of ratios is limited because groups can differ due to  
14 differences in the numerator, denominator, or both. Moreover, evolutionary  
15 ecologists are most often interested in differences in specific variance component  
16 among groups rather than in differences in variance ratios *per se*.
- 17 3. Recommendations for how to infer whether groups differ are not clear in the  
18 literature. We show how questions regarding variance components and how they  
19 vary among groups can be answered using Hierarchical Linear Model approaches  
20 (HLMs).
- 21 4. Frequentist and Bayesian frameworks have similar abilities to infer differences in  
22 variance components. However, simulations where differences occur at higher  
23 levels of organization can be difficult to detect at low sample sizes.
- 24 5. We provide guidelines for how to report and draw inferences based on comparisons  
25 of variance components and variance ratios.

26 Running Head: Comparing variation within datasets

27 Keywords: Heritability, repeatability, individual niche specialization, animal personality,  
28 phenotypic variation, functional traits, mixed models, individual variation

## 29 INTRODUCTION

30 Our understanding of many evolutionary and ecological processes is underpinned by an  
31 estimation of variance ratios. For example, evolutionary change is dependent on the ratio  
32 of additive genetic variation ( $V_a$ ) to total phenotypic variation ( $V_p$ ), more commonly known  
33 as narrow-sense heritability ( $\frac{V_a}{V_p}$  or  $h^2$ ):

$$34 \quad \Delta z = h^2 s \quad (\text{equation 1})$$

35 where the change in a population's mean from one generation to the next ( $\Delta z$ ) is based on  
36 the selection differential ( $s$ ) and the trait's heritability ( $h^2$ ) (breeder's equation, Lush 1937).  
37 Considerable effort has been directed toward estimating and comparing heritability  
38 estimates among taxa or among trait types (Mousseau and Roff 1987; Stirling et al. 2002;  
39 Dochtermann et al. 2019), with these comparisons sometimes used to argue that some  
40 traits are under greater selection than others (Mousseau and Roff 1987).

41 Variance ratios are similarly important across ecology. For example, individual  
42 resource specialization can be estimated as the proportion of variation in an individual's  
43 resource use relative to the species' total variation in resource use (Bolnick et al. 2002):

$$44 \quad \textit{specialization} = \frac{WIC}{TNW} \quad (\text{equation 2})$$

45 where TNW is a species' total niche width (total resource variation) and WIC is "the  
46 average variance of resources found within individual's diets".

47 Interest in variance ratios spans a broad swath of evolutionary ecology (Table 1).  
48 This includes interest in repeatability and "animal personality" (Lessells and Boag 1987;  
49 Bell et al. 2009; Dingemanse and Dochtermann 2013; Dochtermann et al. 2015) and

50 interest in community ecology regarding the distribution of functional trait variation  
51 expressed within versus among populations or species (Violle et al. 2012).

52         While the use of variance ratios can facilitate comparison among populations,  
53 inferences based on these ratios can be highly misleading (Houle 1992; Wilson 2018). If a  
54 variance ratio is compared between two groups, this comparison is only narrowly  
55 interpretable. Specifically, such a comparison is not informative regarding the biological  
56 basis of a difference or lack thereof. This is the case because variance ratios can differ when  
57 their numerators differ, their denominators differ, or because both differ. Indeed, variance  
58 ratios can be equal despite having different numerators and denominators values.

59

60 **Table 1.** Examples variance ratios found in the the ecological and evolutionary literature.

Discipline	Variance ratio	Definition	Description	References
Quantitative Genetics	<i>Heritability</i>	$h^2 = Va / Vp$	The proportion of variation attributable to additive genetic variance ( $Va$ )	Mousseau & Roff 1987
Behavioral Ecology	<i>Repeatability</i>	$R = Vi / Vp$	The proportion of variation attributable to among-individual differences ( $Vi$ )	Lessels & Boag 1987
Ecology	<i>Individual Niche Specialization</i>	$S = WIC / TNW$	The proportion of variation attributable to within-individual preference in niche ( $WIC$ ) (usually expressed as standard deviations)	Bolnick et al. 2002
Community Ecology	<i>T-ratios</i>	$T_{IP/IC} = V_{IP} / V_{IC}$	The proportion of variation attributable to within-population variance ( $V_{IP}$ ) relative to the community variance ( $V_{IC}$ )	Violle et al. 2012
		$T_{IC/IR} = V_{IC} / V_{IR}$	The proportion of variation attributable to community variance ( $V_{IC}$ ) relative to the regional pool variance ( $V_{IR}$ )	

61

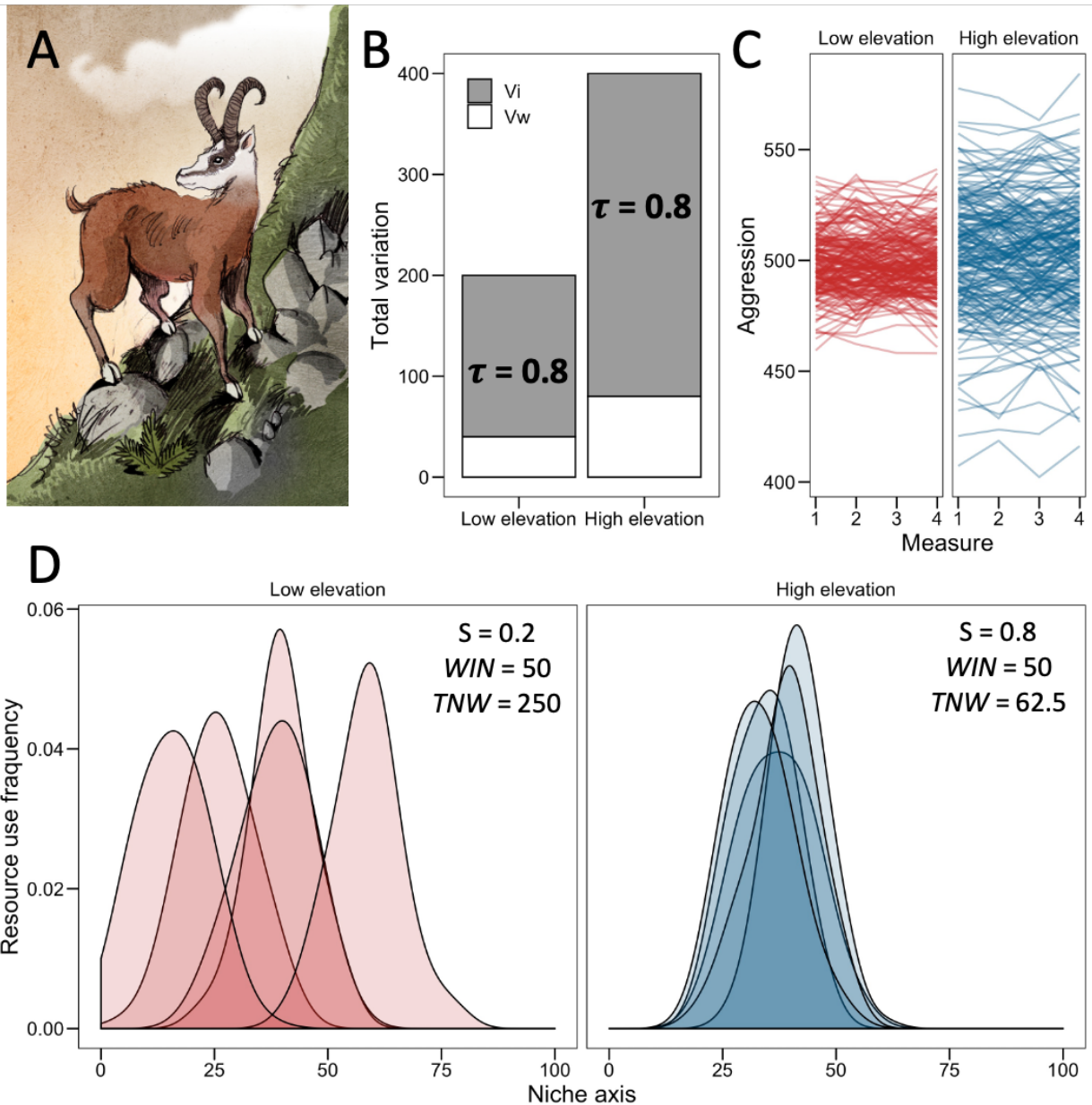
62 Legend:  $Va$ : additive genetic variance in trait,  $Vi$ : among-individual variance in trait,  $Vp$ : total (i.e. phenotypic) variance in trait,  
 63  $WIC$ : within-individual variance in niche preference,  $TNW$ : Total niche width,  $T_{IP}$ : total amount of trait variation in a  
 64 community,  $V_{IP}$ : within-population variance in trait,  $V_{IC}$ : community variance in trait,  $V_{IR}$ : regional pool variance.

65           To illustrate that point further, let us consider the following scenario: researchers  
66 are studying the behaviors and dietary habits of two populations of the mythical Dahu  
67 (*Dahu desterus*; Figure 1A) at different elevations. These elusive creatures have shorter  
68 hind-legs on their left side, thus only allowing for clockwise movement (Chartois & Claudel  
69 1945; Jacquat 1995). While measuring aggressive interactions, researchers find no  
70 differences in means between populations and similar behavioral repeatabilities ( $\tau = 0.8$ ;  
71 Figure 1B). The researchers notice, however, that there are large differences in the among-  
72 and within-individual variances of each population. Had researchers only examined  
73 repeatabilities and mean differences they would inappropriately conclude that the  
74 populations are behaviorally equivalent. However, paying attention to the variance  
75 components reveals that individuals from the high-altitude population are much more  
76 distinct from one another in their aggressive tendencies while, at low-altitude, individuals  
77 show little departure from the population average (Figure 1B, C).

78           These researchers are also curious as to whether the harsher climate at the top of  
79 the mountain range leads to a narrower dietary breadth. Researchers predict that  
80 individual resource specialization will be higher in the low elevation population, as *D.*  
81 *desterus* have more food options to choose from. To the researcher's surprise, they find  
82 much higher individual resource specialization in the high-altitude population:  $S_1 = 0.2$ ,  $S_2 =$   
83  $0.8$ . Upon examining the specific values of among- and within-individual variation in niche,  
84 they find that these differences are a result of the high elevation population having a much  
85 narrower total niche width (Figure 1D) while the within-individual variation in niche  
86 preference is equal between populations. This means that it is the difference in diet  
87 preference among individuals that drives the difference between the two populations. With

88 more diverse resources available at low elevation each individual can specialize along the  
89 total niche axis, yet the breadth of diet preference within-individuals is unchanged in both  
90 populations.

91 For both traits, exclusive reliance on ratios would have led to either inappropriate  
92 or incomplete inferences. Due to these problems with interpretations of variance ratios,  
93 what would be of greater use to researchers is to understand differences in the underlying  
94 variance components themselves.



95  
 96 **Figure 1.** Reliance on variance ratios can lead to misleading inferences. (A) The elusive Dahu (*Dahu*  
 97 *dexterus*) in its natural environment. (B) Two populations of Dahus living at different elevations do  
 98 not differ in their repeatability of aggressive interactions ( $\tau$ ). (C) By plotting the individual  
 99 aggression scores over the course of multiple measurements, it is clear that individuals are more  
 100 distinct in their aggressive behavioral strategies at high elevation. This inference cannot be made by  
 101 investigating repeatability alone. (D) The two population have very different resource  
 102 specialization indices ( $S$ ). A more accurate inference is that individuals do not differ in niche width  
 103 ( $WIN$ ), it is instead the total niche width ( $TNW$ ) that is narrower in the high-altitude population.  
 104 Figure code available here: <https://osf.io/5aw42/>

105 Illustration: [Philippe Semeria](#) (CC BY 3.0 license)



106 *A statistical framework for comparing variance components*

107 The statistical procedures necessary for the estimation of variance components and ratios  
108 within a single population have been the subject of much attention ( e.g. mixed models for  
109 repeatability: Dingemanse and Dochtermann 2013; animal models for heritability: Wilson  
110 et al. 2010; individual niche specialization: Bolnick et al. 2002; Coblentz et al. 2017;  
111 functional trait variation: Nakagawa and Schielzeth 2012; Violle et al. 2012; Carmona et al.  
112 2016). There is also a long history in quantitative genetics regarding the comparison of  
113 variances and covariance structures among groups (Shaw 1991, Arnold & Phillips 1999,  
114 Roff 2002, Roff et al. 2012, Aguirre et al. 2014). Unfortunately, these quantitative genetic  
115 approaches have been poorly disseminated across fields (but see Dochtermann & Roff  
116 2010 and White et al. 2019). Here we describe and investigate methods for detecting  
117 differences in variance components amongst groups. Specifically, we compare the strength  
118 and weaknesses of three statistical approaches: comparison of confidence intervals, model  
119 comparison with AIC, and Bayesian estimation of the difference in variance components.

120 We consider a scenario where a phenotypic attribute,  $y$ , is measured repeatedly for  
121 individual organisms occupying one of two different environments (E1 and E2) and in  
122 which variation occurs among and within experimental units ( $V_H$  and  $V_W$  respectively). We  
123 use the subscripts  $H$  and  $W$  to denote that the among-unit variance ( $V_H$ ) represents the  
124 “higher-level” variance used for comparing differences between the two environments,  
125 while the within-unit variance ( $V_W$ ) indicates differences in trait value occurring within  
126 environments during the course of the experiment. This is a broadly applicable scenario  
127 that can correspond to the comparison of the repeatability of a phenotype between  
128 environments, the comparison of diet specialization for individuals occupying different

129 environments, or how functional traits vary among and within species in two different  
130 environments.

131 An easy way to compare these variance components and their ratios ( $\tau = V_H/(V_H +$   
132  $V_W)$ ) is to estimate the variance components for each environment in separate statistical  
133 models. We can then test for differences in variance components and ratio by  
134 environments based on whether their confidence intervals overlap or not. While  
135 straightforward, this method suffers from several limitations. First, basing inference on the  
136 overlap of 95 % confidence intervals is overly conservative (Barr 1969), especially when  
137 sample size is low. It is instead whether the confidence interval for the *difference* in  
138 variances excludes 0 that is relevant for drawing inferences. This difference cannot be  
139 directly estimated from the approach we have described. However, statistical significance  
140 can still be assessed by comparing the overlap of the 83% confidence intervals for variance  
141 components, a threshold that provides a better approximation for an  $\alpha = 0.05$  for the null  
142 hypothesis of no difference (Austin and Hux 2002; MacGregor-Fors and Payton 2013;  
143 Hector 2015). Second, by estimating variance components in separate statistical models,  
144 the hierarchical structure of the data, i.e. the variance components nested within the  
145 environments, has been broken. As a result, potential average differences in the traits of  
146 interest are not appropriately tested.

147 Instead, we suggest that a more appropriate procedure would be the use of a  
148 Hierarchical Linear Model (HLM) where the among- and within-unit variance is estimated  
149 for each environment within the same statistical model. This statistical model can be  
150 described by the following equation:

$$151 \quad y_{ij} = \beta_0 + \beta_1 \textit{Environment} + \textit{unit}_{0j} + e_{0ij} \quad (\text{equation 3})$$

152  $unit_{0j} \sim MVN(0, \Omega_{unit}); \quad \Omega_{unit} = \begin{bmatrix} V_{unit0} E_1 & 0 \\ 0 & V_{unit0} E_2 \end{bmatrix}$

153  $e_{0ij} \sim MVN(0, \Omega_e); \quad \Omega_e = \begin{bmatrix} V_{e0} E_1 & 0 \\ 0 & V_{e0} E_2 \end{bmatrix}$

154 where  $y_{ij}$  describes the phenotypic traits for the  $i$ th experimental unit and  $j$ th observation.

155  $unit_{0j}$ , is the deviation from an overall intercept,  $\beta_0$ , for the  $j$ th experimental unit.  $\beta_1$

156 represents the regression coefficient for the fixed effect of environment (here a contrast

157 coefficient). The random intercepts and residual variance ( $e_{0ij}$ ) both follow a multivariate

158 normal distribution, and  $\Omega_{unit}$  and  $\Omega_e$ , are the variance-covariance matrices at the among-

159 and within-unit levels respectively.

160 The diagonal elements of these matrices represent the among- ( $H$ ) and within-unit

161 ( $W$ ) variances by environment and the off-diagonal elements represent the cross-

162 environment correlation (set to 0 if units are only ever evaluated in one of the two

163 environments). This formulation has the advantage of allowing considerable flexibility in

164 the specification of the statistical models considered (Dingemans and Dochtermann

165 2013). HLMs are now available for most statistical software and their generalized

166 extensions can accommodate non-normal error distributions (Table 2).

167 Upon fitting HLMs, several methods are then available to determine whether a

168 variance ratio or components of the ratio differ by environment. Specific hypotheses of

169 which variance component differs across environment can be easily tested via model

170 comparison. For example, a model where only the among-unit variance differs by

171 environment can be compared to a null model where the among and within-unit variance

172 are kept constant across environments (Royauté et al. 2019). These models can be

173 estimated within a frequentist framework via restricted maximum likelihood or a Bayesian

174 framework and suitable decision criteria can be used to determine which model best fits  
175 the data. In the case of restricted maximum likelihood estimation, it is also possible to use  
176 likelihood ratio tests to compare these models. Note however that the proper degrees of  
177 freedom to apply to each model is unclear and additional care should be taken when using  
178 this method (Pinheiro and Bates 2000; see Santostefano et al. 2016 for a recent example).

179         In many cases, researchers are also interested in whether the difference in variance  
180 components have a biologically meaningful effect. In other words, when asking questions  
181 about whether variance components vary between environments, we are mostly interested  
182 in the *magnitude of the difference* in these components across environments. While model  
183 comparison of HLMs can help us understand whether a statistically detectable difference is  
184 observable across environments, the magnitude of the difference can only be determined  
185 by examining the difference in variance components among environment:  $\Delta V$  estimated as  
186  $V_{E2} - V_{E1}$  in our case. When the trait of interest is expressed as standard deviation units (i.e.  
187 mean centered and scaled to the standard deviation of the dataset), this difference can be  
188 considered an effect size for the magnitude of the difference among variance components,  
189 thus making comparisons across studies possible (Royauté et al. 2015; Hamilton et al.  
190 2017; Royauté and Dochtermann 2017). Note that  $\Delta V$  could also be expressed on a ratio  
191 scale ( $V_{E2}/V_{E1}$ ) or on a log-additive scale ( $\log(V_{E2}) - \log(V_{E1})$ ). We used  $\Delta V$  on an additive  
192 scale because it allows the most straightforward interpretation and functions in cases  
193 where a variance component is zero or approaching zero.

194 **Table 2.** Packages and softwares allowing to test for differences in variance components using Hierarchical Linear Models (HLM) along  
 195 with parameter estimation method (maximum likelihood (ML), restricted maximum likelihood (REML) or Bayesian framework) and  
 196 inference method (Likelihood Ratio tests (LRT), AIC or credible interval overlap). Tthis list is non-representative of the diversity of option  
 197 available and is based on widely-used commercial softwares and R packages.  
 198

Package or software	Free or commercial	Estimation	Testing method	Among-unit variance by group	Within-unit variance by group	Distributions handled	Comments	Reference
ASREmL	Commercial	ML/REML	LRT, AIC	Yes	Yes	Gaussian		Gilmour et al. (2015)
SAS	Commercial	ML/REML	LRT, AIC	Yes	Yes	Gaussian, Poisson, Binomial		SAS Institute Inc.
nlme	Free	ML/REML	LRT, AIC	Yes	Yes	... Gaussian		Pinheiro and Bates (2000)
lme4	Free	ML/REML	LRT, AIC	Yes	No	Gaussian, Poisson, Binomial		Bates et al. (2015)
R-INLA	Free	ML/REML	LRT, AIC	Yes	Yes	... Gaussian		Lindgren, and Rue (2015)
glmmTMB	Free	ML/REML	LRT, AIC	Yes	Yes	Gaussian, Poisson, Binomial		Brooks et al. 2017
hglm	Free	ML/REML	LRT, AIC	Yes	Yes	... Gaussian, Poisson, Binomial	Within-unit variance modelled as Gamma distribution	Rönnegård et al. (2010)
MCMCglmm	Free	Bayesian	DIC, overlap of credible intervals	Yes	Yes	... Gaussian, Poisson, Binomial		Hadfield (2010)
brms	Free	Bayesian	WAIC, LOO, overlap of credible intervals	Yes	Yes	... Gaussian, Poisson, Binomial	Within-unit variance modelled as log-normal distribution	Bürkner (2017)

199  $\Delta V$  can be calculated from the maximum likelihood estimates in a frequentist  
200 framework but calculation of the uncertainty around this estimate is not straightforward  
201 and can require additional steps such as bootstrapping. In a Bayesian framework, the  
202 calculations are much simpler given that the distribution of  $\Delta V$  can be directly estimated by  
203 taking the difference in the posterior distribution of  $V_{E2} - V_{E1}$ . The posterior mode of  $\Delta V$  can  
204 then be interpreted as the estimated strength of  $\Delta V$ , with credible intervals representing  
205 the precision around this estimate.

206 In summary, approaches based on HLM and their generalized extensions allow great  
207 flexibility and are well suited to study questions related to how variation in phenotypic  
208 traits varies at multiple levels of organization. In the next section, we describe the  
209 performance of HLMs to detect differences in variance components.

## 210 **METHODS**

### 211 *Data simulations*

212 To compare the performance of statistical procedures for the detection of differences in  
213 variance components and variance ratios, we performed a series of simulations based on  
214 the scenarios illustrated in Figure 2. In these scenarios a phenotypic attribute  $y$  is  
215 measured in two different environments (E1 and E2) and variation occurs among and  
216 within experimental units ( $V_H$  and  $V_W$  respectively). In scenarios A through C the variance  
217 ratio differs by an equal amount between the two environments ( $\Delta\tau = 0.3$ ), but the  
218 underlying driver of this difference is either due to a difference in the among-unit variance  
219 (A), in the within-unit variance (B) or in both the among and within-unit variance (C). Note  
220 that for scenario C, the total variance remains the same between environments. In

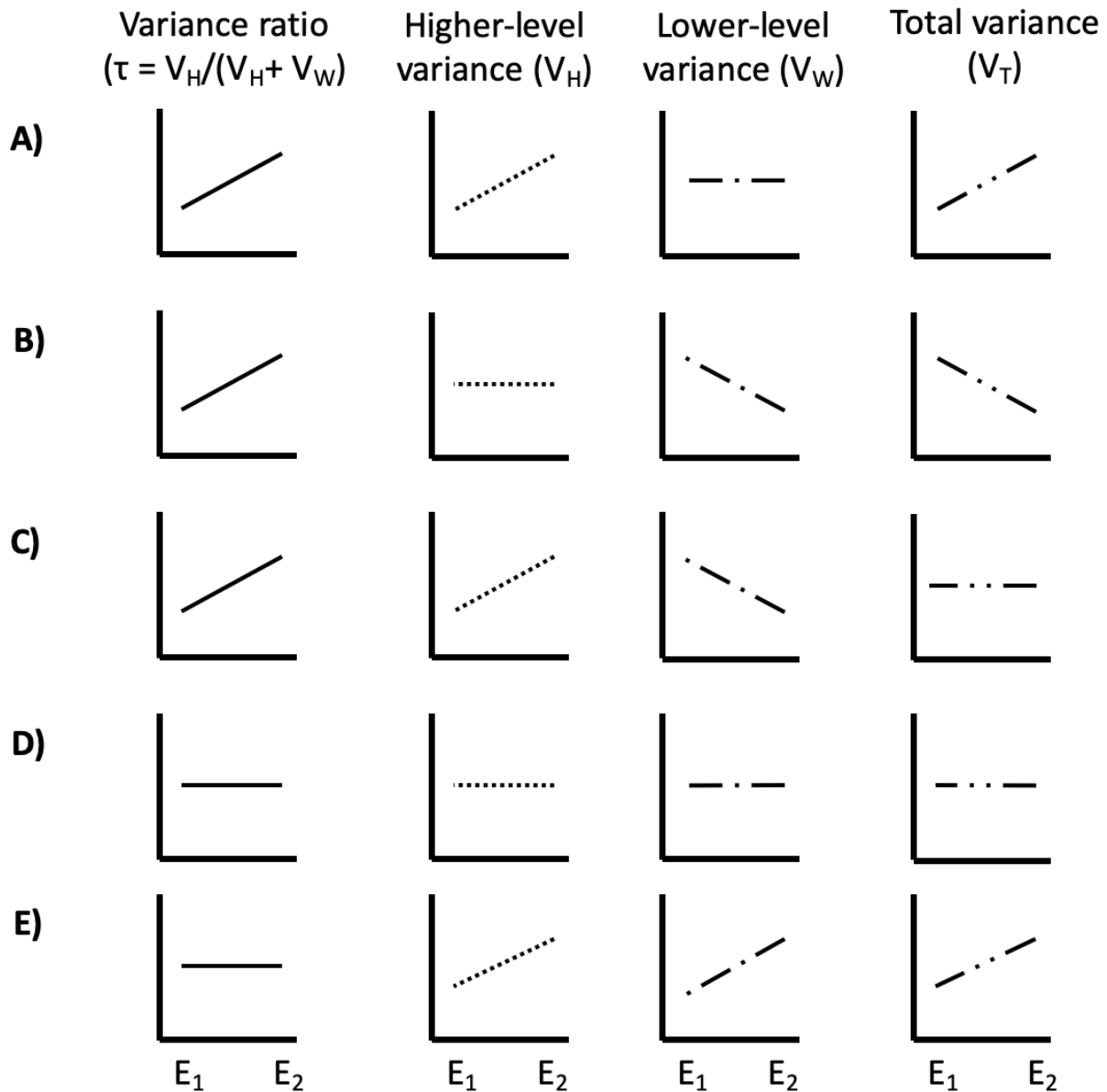
221 scenarios D and E, we explore cases where the variance ratios are equal among  
222 environment, either because all variance components are equal as well (D) or in spite of  
223 differences in all other variance components (E) (see Table S1 for exact values for all  
224 parameters).

225 Using the R statistical environment (R Core Team 2017), we generated 500 datasets for  
226 each of the following combinations:

- 227 • Sample size varying from 20 to 200 units by increments of 20 for each environment  
228 (sample size was equal between the two environments)
- 229 • Number of repeated measures taken on each unit varying from 2 to 6 repeated  
230 measures by increments of 1
- 231 • Five different scenarios of known difference in variance ratios as described in  
232 Figure 1 and Table S1.

233 Each dataset was simulated by sampling from a Gaussian distribution for the random  
234 (among-unit values) and the error (within-unit) terms. This resulted in a total of 125,000  
235 datasets on which we tested three different statistical procedures to detect differences in  
236 variance components and variance ratios. We provide all R code for data generation and  
237 analysis in Supporting Information 1.

238



239

240 **Figure 2.** Scenarios used in simulations detailing how differences or lack of difference in  
 241 variance ratios can arise from different patterns in the underlying variance components  
 242 (Exact values can be found in Table S1). Panels A-C indicate scenarios where the total  
 243 variation differs between two environments (E1 and E2) due to differences in the higher  
 244 group level variance (A), the lower level variance (B) or both (C). Panels D-E indicate  
 245 scenarios where the ratios remains constant across environments, because all variance  
 246 components are identical (D) or in spite of variance component being different among  
 247 environments (E).

248



249 *Comparison of confidence interval overlap from separate mixed models*

250 We first compared the overlap of 83 % confidence intervals for variance component when  
251 estimated from separate linear mixed models. We specified one mixed model for  
252 environment 1 and one for environment 2. These models are a simplified version of the one  
253 presented in equation (3):

$$254 \quad y_{ij} = \beta_0 + unit_{0j} + e_{0ij} \quad (\text{equation 4})$$

$$255 \quad unit_{0j} \sim \mathcal{N}(0, V_{unit});$$

$$256 \quad e_{0ij} \sim \mathcal{N}(0, V_e)$$

257 The experimental units in the environment of interest are included as random effects and  
258 no additional fixed effect are needed. Upon fitting these models, we computed 83 %  
259 confidence intervals for the among and within-unit variance. Datasets where these  
260 intervals did not overlap were considered as statistically different.

261 *Frequentist HLM with AIC model comparison*

262 Our second approach was to fit the HLM approach described above and test for the for the  
263 significance of the difference in among- and within-unit variance using likelihood ratio  
264 tests. Specifically, we compared the following models:

265 We specified four different mixed models corresponding to the four different possibilities  
266 by which variance components may differ (see also Royauté et al. 2019):

- 267 • Model 1: a null model where the among ( $V_H$ ) and within-unit variance ( $V_W$ ) was kept  
268 constant among environments.

- 269 • Model 2: a model where only the among-unit variance differs among environments,  
270 while the within-unit variance is kept constant ( $V_H \neq$  &  $V_W =$ )
- 271 • Model 3: a model where only the within-unit variance differs among environments  
272 while the among-unit variance is kept constant ( $V_H =$  &  $V_W \neq$ )
- 273 • Model 4: a model where both the among and within-unit variance were allowed to  
274 vary among environments ( $V_H \neq$  &  $V_W \neq$ )

275 For each dataset combination, we then compared each model's Akaike's Information  
276 Criterion value (AIC). AIC allows to compare the relative fit of statistical models and models  
277 with lower AIC values indicate better support relative to competing models. These  
278 simulations and this analytical framework is similar to previously used approaches (e.g.  
279 Jenkins 2011; Shaw 1991; Tüzün et al. 2017). These models were specified using the *nlme*  
280 package for mixed models (Pinheiro et al. 2000) using Restricted Maximum Likelihood  
281 (REML).

### 282 *Bayesian HLM and difference in variance components*

283 We next fit a mixed model where variances among and within units were allowed to vary  
284 between environments (as in model 4 described above) to each randomly generated  
285 dataset. We calculated the posterior mode for the difference in variance components  
286 (calculated as  $\Delta V = V_{E2} - V_{E1}$ ) and estimated the 95 % credible intervals based on the  
287 Highest Posterior Density of this distribution. 95 % credible intervals excluding 0 were  
288 taken to indicate statistically detectable differences in variance components among  
289 environments. All models were run with the *MCMCglmm* package (Hadfield 2010) using  
290 default iteration settings to shorten computing time (13000 iterations, 3000 burn-in

291 iterations and thinning interval of 10 iterations). We used priors that were minimally  
292 informative for the variance components (See SI1 and SI3 for prior specification and a  
293 discussion on priors).

294 *Probability of correct model identification, precision, bias and accuracy estimations*

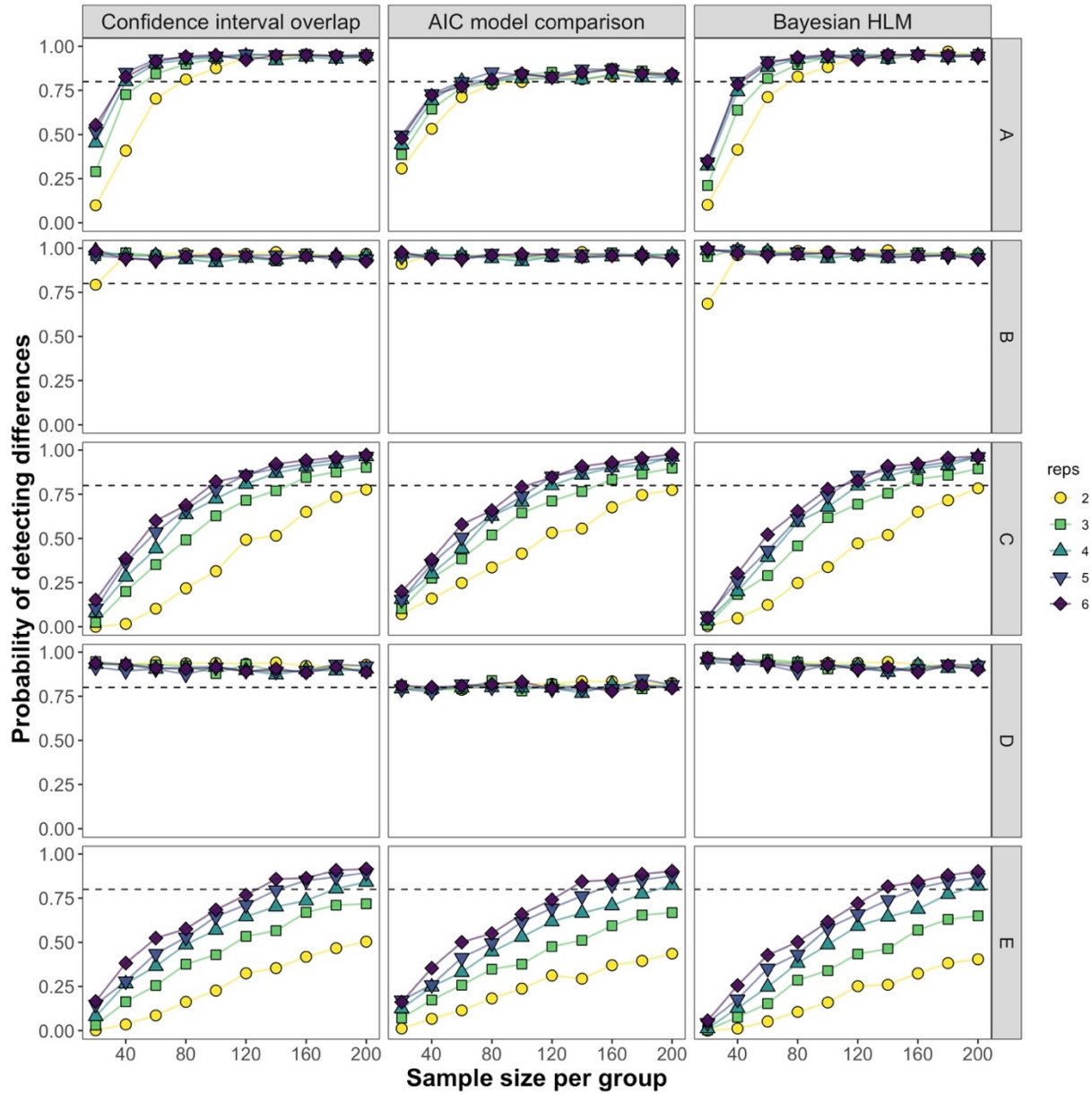
295 We calculated the probability of detecting the model with the correct difference in variance  
296 components (hereafter abridged to probability of detecting differences), precision, relative  
297 bias and accuracy under each scenario and sampling design to compare the performance of  
298 maximum likelihood and Bayesian mixed models. For Method 1 (overlap of 83 % intervals),  
299 we assigned values of 1 when significant differences in variance components were detected  
300 in directions predicted by the data generating process, and 0 otherwise. For Method 2, we  
301 calculated the probability of detecting differences as the proportion of times the model  
302 with the lowest AIC matched the generating model. For Method 3, we calculated whether a  
303 given model detected a difference in variance components based on the overlap of the 95 %  
304 credible intervals of the  $\Delta V$  posterior distribution with 0. As in Method 1, we then assigned  
305 values of 0 or 1 based on whether the detected difference matched with the data  
306 generation process of the corresponding scenario. We calculated the probability of  
307 detecting differences as the proportion of analyzed datasets in which we detected  
308 differences in the direction predicted by each scenario and statistical method. Precision,  
309 indicating the similarity of the results produced by simulations with a given scenario, was  
310 calculated as the difference between 25 % and 75 % quantiles of estimates (van de Pol  
311 2012). To calculate the relative bias (in %) for each statistical approach by scenario, we  
312 calculated the mean difference between the expected value and the value observed in each

313 of the 500 simulations. Finally, we report the root mean square of error (RMSE) for each  
314 scenario and sample sizes. This metric calculates how close estimates are to the expected  
315 values and serves as an estimate of the accuracy of each statistical approach by scenario.

316

## 317 **RESULTS**

318 The probability of correctly detecting differences in variance components did not differ  
319 substantially between frequentist and Bayesian methods of estimation (Figure 3). The  
320 highest probability to detect differences was observed for analyses of scenario B (the  
321 variance ratio differs due to a difference in within-unit variance) and scenario D (no  
322 difference in variance components or ratio). This was the case for all statistical approaches  
323 (Figure 3). Statistical power to differentiate scenarios A, C and E was lower, especially with  
324 small sample sizes and low number of repeated measures (Figure 3). Importantly, no  
325 statistical method seemed to outperform all others in across scenarios. Our results are  
326 consistent with previous simulations showing that the among-unit variance component is  
327 particularly difficult to estimate at small sample sizes (Dingemanse & Dochtermann 2013).



328  
 329 **Figure 3.** Effect of sampling design on the probability to detect differences in variance  
 330 components by scenario type and statistical modeling approach. Each point represents the  
 331 probability of detecting the correct differences in variance averaged over 500 simulated  
 332 datasets. A represents a scenario where only the among-unit variance ( $V_H$ ) varies between  
 333 environments, B represents a case where the within-unit variance ( $V_W$ ) varies between  
 334 environments, and both among and within-unit variance vary between environments in  
 335 scenario C. In scenario D, all variance components are equal while in scenario E, variance  
 336 components are different but variance ratios are equal across environments. Dashed lines  
 337 correspond to 80 % threshold similar to recommendations for power analyses.  
 338

339  
340 In scenarios B and D, the correct differences among variance components was  
341 identified > 80 % of the time, even at low sample sizes (Figure 3). In all other cases this  
342 threshold was only reached with high sample sizes and a high number of repeated  
343 measures. For scenarios C and E, datasets with only 2 repeated measures per unit never  
344 achieved a power above 0.8 even with sample sizes above 200 units per environment (i.e. a  
345 minimum of 800 total measurements, Figure 3). Increasing the number of repeated  
346 measures only marginally alleviated the problem. For example, in scenario C, only datasets  
347 with 4 or more repeated measures per unit reached statistical power above 0.8 with  
348 sample sizes above 120 units per environments, which is higher than many ecological or  
349 evolutionary studies can provide under realistic scenarios.

350 Note that for AIC model comparison, we calculated power as the number of times  
351 the best model corresponded to the generating model. A more conservative approach is to  
352 calculate the proportion of times the best model is at least 2 AIC units lower than the  
353 second model. This method corresponds to a common threshold to detect statistically  
354 distinct models (Burnham and Anderson 1998). When using this more conservative  
355 threshold (Figure S1), datasets generated according to scenarios A and D were never  
356 statistically distinguishable from non-generating models, although the correct model was  
357 consistently ranked as the best model. This is likely because when the generating model  
358 does not include differences in the within-unit variability (scenarios A and D), sampling  
359 error is erroneously identified as heterogeneity. At smaller sample sizes this error is  
360 greater on average, and thus detectable. At larger sample sizes this sampling error is  
361 smaller but more easily detected and therefore manifests as different between groups. To  
362 address this, in addition to measures of variance differences like the described  $\Delta V$  statistic,

363 researchers should also compare mean-standardized variance estimates like the coefficient  
364 of variation or Houle's evolvability between groups (Houle 1992; Hansen et al. 2011;  
365 Dochtermann and Royauté 2019).

366         The comparison of relative bias, precision, and accuracy among statistical methods  
367 produced mixed results. On average, Bayesian HLMs consistently underestimated the  
368 among-unit variance for scenarios in which the among-unit variance differed between  
369 environments (scenarios A, C, and E) resulting in a severe bias at small sample sizes (Figure  
370 S2). However, Bayesian HLMs also had higher precision and accuracy compared to  
371 maximum likelihood (Figure S3, S4). This means that Bayesian estimates tend to be  
372 consistently more conservative than maximum likelihood regarding the magnitude of the  
373 among-unit variance but that these estimates nonetheless more closely matched simulation  
374 conditions.

## 375 **DISCUSSION**

376 Comparing variability across datasets is important for many questions in evolutionary  
377 ecology (e.g. Table 1). However, variance ratios are not sufficient to address questions  
378 about how variance components vary across environments, populations, or sexes. The  
379 inability to determine why groups differ based on ratios is in addition to the numerous  
380 conceptual and theoretical problems inherent to the estimation of ratios (Houle 1992;  
381 Hansen et al. 2011). Instead, many questions require the explicit comparison of variance  
382 components.

383         Our simulations show that regardless of the statistical methods used, comparing  
384 variance components across groups is a "data hungry" question. Scenarios where the

385 among-unit variance differed between environments were particularly hard to detect at  
386 low sample sizes. Our objective was not to provide a full exploration of parameter space in  
387 order to define the proper sample sizes to detect differences of various magnitude for each  
388 variance component. Instead, we focused on a subset of scenarios that are likely to be  
389 common in ecology and evolution.

390         Given the issues discussed above, how should researchers interested in ecological  
391 and evolutionary variation design their studies and report their findings? Based on our  
392 simulations, the probability to detect differences in variance components will depend in  
393 large part on the ability to estimate the among-unit variance component ( $V_H$ ). A simple rule  
394 for sampling can therefore be to estimate the sample size needed to detect the lowest  
395 among-unit variance value (see, for example, Martin et al. 2011; van de Pol 2012;  
396 Dingemanse and Dochtermann 2013) and multiplying that value by the number of  
397 experimental groups involved. We also recommend that power calculations be conducted  
398 prior to the experiment whenever possible (see R code for *a priori* power analyses in SI2  
399 and R Markdown tutorial in SI3).

400         We suggest that researchers report their results in a manner that focuses on the  
401 magnitude of the difference in variability between experimental groups rather than solely  
402 focus on statistical significance. To this effect, we believe that reporting the results of the  
403 full model rather than just the most parsimonious model will be most appropriate in most  
404 cases (i.e. model 4 in our conceptual example). This is because model selection only gives  
405 information on whether differences among groups are statistically detectable. In contrast,  
406 questions regarding the magnitude and precision of the estimated differences are  
407 answerable only with interpretation of the most complete statistical model (see tutorial in



408 SI4). In addition to presenting results of the full model, we suggest that measures of effect  
409 sizes for the differences in variance component also be presented. As reported above,  $\Delta V$   
410 provides a simple metric to estimate the magnitude of these differences, but it is by no  
411 mean the only one. In our theoretical example, the mean trait value did not differ by  
412 environments, but in many cases mean and variance are related. In such cases, using  
413 comparisons based on Houle's (1992)  $I^2$  value or coefficients of variation for each  
414 component as opposed to variance component themselves can be preferable (Hansen et al.  
415 2011; Dochtermann and Royauté 2019). Effect sizes based on the coefficient of variation  
416 can also be calculated within an HLM framework as described by Nakagawa et al. (2015)  
417 (see also Carmona et al. 2016 and Fontana et al. 2018 for approaches relevant to functional  
418 trait diversity).

419         While we limited our conceptual example to comparisons between two  
420 environments, the HLM approach we propose is by no mean restricted to two-groups  
421 comparisons. For example, Jenkins (2011) used model comparison to tease apart the  
422 relative influence of sex, species and their interaction on the expression of behavioral  
423 variation in kangaroo rats. Similarly, Coblenz et al. (2017) show how model selection  
424 combined with Bayesian HGLM can allow the comparison of indices of diet specialization  
425 within and among species. In both cases, model section can provide a first pass at whether  
426 differences in variance components are detectable among groups, while specific pairwise  
427 comparisons of effect sizes (using  $\Delta V$  or other metrics) will allow discernment of the most  
428 pronounced differences in variance component. Regardless of the statistical approach used,  
429 we suggest it is important that researchers clearly outline the direction and, when possible,  
430 magnitude of the expected effects in their predictions.

431 Finally, our conceptual examples focus exclusively on the case of “well-behaved”  
432 data with normal error distributions. While these comparisons can be made with  
433 generalized extensions to HLMS (i.e. HGLMs), extra care must be taken to appropriately  
434 estimate and compare the within-unit variance depending on the error distribution  
435 specified (Nakagawa & Schielzeth 2010).

## 436 **CONCLUSIONS**

437 Variance ratios are straightforward metrics to describe how various ecological and  
438 evolutionary processes occur. However, comparing these ratios across studies or group can  
439 be misleading if poor attention is given to the specific variance components making up  
440 those ratios. More importantly, as we have shown, a lack of difference in these ratios does  
441 not mean that variance components are equal among groups. Given these limitations, we  
442 advocate for techniques allowing the estimation of differences in each variance  
443 components rather than focusing solely on variance ratios. The statistical tools allowing  
444 comparison of trait variation have become increasingly sophisticated and now allow asking  
445 very precise questions. Specifically, we can now ask how trait variation is generated and  
446 how variation differs among groups. However, despite the availability of these tools,  
447 researchers interested in ecological and evolutionary variation must remain careful in their  
448 study designs. As our simulations show, scenarios involving differences in among-unit  
449 variance are particularly difficult to detect without substantial sample sizes. Finally, we  
450 hope the statistical approaches and tools for power analysis presented here will allow for  
451 appropriate comparisons of trait variation in ecological and evolutionary studies.

452

453

454 **Acknowledgments**

455 We thank the participants of the Statistical Quantification of Individual Differences (SQuID)  
456 Symposium at the 2016 ISBE Congress for helpful discussions. We also thank Russel  
457 Bonduriansky, Ben Bolker and two anonymous reviewers for helpful comments on a  
458 previous version of this manuscript. This study was funded by NSF IOS-1557951 (to NAD)  
459 and the Department of Biological Sciences at North Dakota State University.

460

461 **Author contribution**

462 Each author contributed equally to the design, analysis and writing of the manuscript.

463

464 **Data availability**

465 All code and data for simulations is available on the Open Science Framework's project for  
466 this article: <https://osf.io/5aw42/>

467

468 **REFERENCES**

- 469 Aguirre, J., E. Hine, K. McGuigan, and M. Blows. 2014. Comparing **G**: multivariate analysis of  
470 genetic variation in multiple populations. *Heredity* 112:21-29.
- 471 Arnold, S. J., and P. C. Phillips. 1999. Hierarchical comparison of genetic variance-  
472 covariance matrices. II. Coastal-inland divergence in the garter snake, *Thamnophis*  
473 *elegans*. *Evolution* 53:1516-1527.
- 474 Austin, P. C., and J. E. Hux. 2002. A Brief Note on Overlapping Confidence Intervals. *Journal*  
475 *of Vascular Surgery* 36:194–195.
- 476 Barr, D. R. 1969. Using confidence intervals to test hypotheses. *Journal of Quality*  
477 *Technology* 1:256–258.
- 478 Bates, D., M. Maechler, B. Bolker, S. Walker, R. H. B. Christensen, H. Singmann, B. Dai, et al.  
479 2015. Package 'lme4.'
- 480 Bell, A. M., S. J. Hankison, and K. L. Laskowski. 2009. The repeatability of behaviour: a meta-  
481 analysis. *Animal behaviour* 77:771–783.
- 482 Bolnick, D. I., R. Svanbäck, J. A. Fordyce, L. H. Yang, J. M. Davis, C. D. Hulsey, and M. L.  
483 Forister. 2002. The ecology of individuals: incidence and implications of individual  
484 specialization. *The American Naturalist* 161:1–28.
- 485 Brooks, M. E., Kristensen, K., van Benthem, K. J., Magnusson, A., Berg, C. W., Nielsen, A.,  
486 Skaug, H. J., Machler, M. and B. M. Bolker 2017. glmmTMB balances speed and  
487 flexibility among packages for zero-inflated generalized linear mixed modeling. *The*  
488 *R Journal* 9:378-400.
- 489 Bürkner, P.-C. 2017. brms: An R package for Bayesian multilevel models using Stan. *Journal*  
490 *of Statistical Software* 80:1–28.
- 491 Burnham, K. P., and D. R. Anderson. 1998. Practical use of the information-theoretic  
492 approach. Pages 75–117 *in* *Model Selection and Inference*. Springer.
- 493 Carmona, C. P., F. de Bello, N. W. Mason, and J. Lepš. 2016. Traits without borders:  
494 integrating functional diversity across scales. *Trends in ecology & evolution* 31:382–  
495 394.
- 496 Chartois, J., & Claudel, C. 1945. Hunting the dahut: a french folk custom. *The Journal of*  
497 *American Folklore* 58:21-24.
- 498 Coblenz, K. E., A. E. Rosenblatt, and M. Novak. 2017. The application of Bayesian  
499 hierarchical models to quantify individual diet specialization. *Ecology* 98:1535–  
500 1547.

- 501 Dingemanse, N. J., and N. A. Dochtermann. 2013. Quantifying individual variation in  
502 behaviour: mixed-effect modelling approaches. *Journal of Animal Ecology* 82:39–54.
- 503 Dochtermann, N. A., and D. A. Roff. 2010. Applying a quantitative genetics framework to  
504 behavioural syndrome research. *Philosophical Transactions of the Royal Society B-  
505 Biological Sciences* 365:4013-4020.
- 506 Dochtermann, N. A., and R. Royauté. 2019. The mean matters: going beyond repeatability to  
507 interpret behavioural variation. *Animal Behaviour* 153:147–150.
- 508 Dochtermann, N. A., T. Schwab, M. Anderson Berdal, J. Dalos, and R. Royauté. 2019. The  
509 Heritability of Behavior: A Meta-analysis. *Journal of Heredity*.
- 510 Dochtermann, N. A., T. Schwab, and A. Sih. 2015. The contribution of additive genetic  
511 variation to personality variation: heritability of personality. *Proceedings of the  
512 Royal Society B: Biological Sciences* 282:20142201.
- 513 Fontana, S., M. K. Thomas, M. Moldoveanu, P. Spaak, and F. Pomati. 2018. Individual-level  
514 trait diversity predicts phytoplankton community properties better than species  
515 richness or evenness. *The ISME journal* 12:356.
- 516 Gilmour, A. R., B. J. Gogel, B. R. Cullis, S. J. Welham, and R. Thompson. 2015. ASReml user  
517 guide release 4.1 structural specification. Hemel Hempstead: VSN international ltd.
- 518 Hadfield, J. D. 2010. MCMC methods for multi-response generalized linear mixed models:  
519 the MCMCglmm R package. *Journal of Statistical Software* 33:1–22.
- 520 Hamilton, J. A., R. Royauté, J. W. Wright, P. Hodgskiss, and F. T. Ledig. 2017. Genetic  
521 conservation and management of the California endemic, Torrey pine (*Pinus  
522 torreyana* Parry): Implications of genetic rescue in a genetically depauperate  
523 species. *Ecology and Evolution* 7:7370–7381.
- 524 Hansen, T. F., C. Pélabon, and D. Houle. 2011. Heritability is not Evolvability. *Evolutionary  
525 Biology* 38:258.
- 526 Hector, A. 2015. *The New Statistics with R: An Introduction for Biologists*. 1<sup>st</sup> edition.  
527 Oxford ; New York, NY: Oxford University Press.
- 528 Houle, D. 1992. Comparing evolvability and variability of quantitative traits. *Genetics*  
529 130:195–204.
- 530 Jacquat, M. S. 1995. *Le dahu: monographie ethno-étho-biologique publiée à l'occasion de  
531 l'exposition inaugurée le 1er avril 1995*. Editions de la Girafe, Musée d'histoire  
532 naturelle.
- 533 Jenkins, S. H. 2011. Sex differences in repeatability of food-hoarding behaviour of kangaroo  
534 rats. *Animal Behaviour* 81:1155–1162.

- 535 Lessells, C. M., and P. T. Boag. 1987. Unrepeatable repeatabilities: a common mistake. *The*  
536 *Auk* 104:116–121.
- 537 Lindgren, F., and H. Rue. 2015. Bayesian spatial modelling with R-INLA. *Journal of*  
538 *Statistical Software* 63:1-25.
- 539 Lush, J. 1937. *Animal Breeding Plans*. Iowa State College Press, Ames, Iowa.
- 540 Martin, J. G., D. H. Nussey, A. J. Wilson, and D. Réale. 2011. Measuring individual differences  
541 in reaction norms in field and experimental studies: a power analysis of random  
542 regression models. *Methods in Ecology and Evolution* 2:362–374.
- 543 MacGregor-Fors, I., and M. E. Payton. 2013. Contrasting Diversity Values: Statistical  
544 Inferences Based on Overlapping Confidence Intervals. *PloS One* 8, no. 2.  
545 <http://dx.plos.org/10.1371/journal.pone.0056794>.
- 546 Mousseau, T. A., and D. A. Roff. 1987. Natural selection and the heritability of fitness  
547 components. *Heredity* 59:181.
- 548 Nakagawa, S., R. Poulin, K. Mengersen, K. Reinhold, L. Engqvist, M. Lagisz, and A. M. Senior.  
549 2015. Meta-analysis of variation: ecological and evolutionary applications and  
550 beyond. *Methods in Ecology and Evolution* 6:143–152.
- 551 Nakagawa, S., and H. Schielzeth. 2010. Repeatability for Gaussian and non-Gaussian data: a  
552 practical guide for biologists. *Biological Reviews* 85:935–956.
- 553 ———. 2012. The mean strikes back: mean–variance relationships and heteroscedasticity.  
554 *Trends in Ecology & Evolution* 27:474–475.
- 555 Pinheiro, J., and D. Bates. 2000. *Mixed-Effects Models in S and S-PLUS*. Springer Science &  
556 Business Media.
- 557 Roff, D. 2002. Comparing **G** matrices: A MANOVA approach. *Evolution* 56:1286-1291.
- 558 Roff, D. A., J. M. Prokkola, I. Krams, and M. J. Rantala. 2012. There is more than one way to  
559 skin a **G** matrix. *Journal of Evolutionary Biology* 25:1113-1126.
- 560 Rönnegård, L., X. Shen, and M. Alam. 2010. hglm: A package for fitting hierarchical  
561 generalized linear models. *The R Journal* 2:20–28.
- 562 Royauté, R., C. M. Buddle, and C. Vincent. 2015. Under the influence: sublethal exposure to  
563 an insecticide affects personality expression in a jumping spider. *Functional Ecology*  
564 29:962–970.
- 565 Royauté, R., and N. A. Dochtermann. 2017. When the mean no longer matters:  
566 Developmental diet affects behavioral variation but not population averages in the  
567 house cricket (*Acheta domesticus*). *Behavioral Ecology* 28:337–345.

568 Royauté, R., C. Garrison, J. Dalos, M. A. Berdal, and N. A. Dochtermann. 2019. Current energy  
569 state interacts with the developmental environment to influence behavioural  
570 plasticity. *Animal Behaviour* 148:39–51.

571 Santostefano, F., A. J. Wilson, Y. G. Araya-Ajoy, and N. J. Dingemanse. 2016. Interacting with  
572 the enemy: indirect effects of personality on conspecific aggression in crickets.  
573 *Behavioral Ecology* 27:1235–1246.

574 Shaw, R. G. 1991. The comparison of quantitative genetic-parameters between populations.  
575 *Evolution* 45:143-151

576 Stirling, D. G., D. Réale, and D. A. Roff. 2002. Selection, structure and the heritability of  
577 behaviour. *Journal of Evolutionary Biology* 15:277–289.

578 Tüzün, N., S. Müller, K. Koch, and R. Stoks. 2017. Pesticide-induced changes in personality  
579 depend on the urbanization level. *Animal behaviour* 134:45–55.

580 van de Pol, M. 2012. Quantifying individual variation in reaction norms: how study design  
581 affects the accuracy, precision and power of random regression models. *Methods in*  
582 *Ecology and Evolution* 3:268–280.

583 Violle, C., B. J. Enquist, B. J. McGill, L. I. N. Jiang, C. H. Albert, C. Hulshof, V. Jung, et al. 2012.  
584 The return of the variance: intraspecific variability in community ecology. *Trends in*  
585 *ecology & evolution* 27:244–252.

586 White, S. J., Pascall, D. J., and A. J. Wilson. 2019. Towards a comparative approach to the  
587 structure of animal personality variation. *Behavioral Ecology*.

588 Wilson, A. J., D. Réale, M. N. Clements, M. M. Morrissey, E. Postma, C. A. Walling, L. E. B.  
589 Kruuk, et al. 2010. An ecologist’s guide to the animal model. *Journal of Animal*  
590 *Ecology* 79:13–26.

591 Wilson, A. J. 2018. How should we interpret estimates of individual repeatability? *Evolution*  
592 *Letters* 2: 4-8.

593

594

595 **Supporting Information**

596 **SI1:** Zip folder containing the raw data from simulations along with R code for data analysis  
597 and figures (<https://osf.io/5aw42/>).

598 **SI2:** R code for conducting *a priori* power analysis (<https://osf.io/5aw42/>).

599 **SI3:** R tutorial for comparing variance components using *nlme*, *MCMCglmm* and *brms*  
600 packages (<https://osf.io/5aw42/>).

601 **Table S1.** Scenarios tested in simulations to estimate the power to detect differences in  
602 variance components of varying magnitude.

603 **Figure S1.** Effect of sampling design on the probability to detect differences in variance  
604 components by scenario type and statistical modeling approach with  $\Delta AIC > 2$  threshold for  
605 model comparison.

606 **Figure S2.** Effect of sampling design on relative bias by scenario type and statistical  
607 modeling approach.

608 **Figure S3.** Effect of sampling design on estimate precision (width of the interquartile  
609 interval) by scenario type and statistical modeling approach.

610 **Figure S4.** Effect of sampling design on model accuracy (estimated as the root mean square  
611 of error, RMSE) by scenario type and statistical modeling approach.