

# 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, Fargo, ND, USA

3 <sup>b</sup> Current address: Movement Ecology Group, Senckenberg Biodiversity and Climate

4 Research Centre (SBIK-F), Frankfurt, Germany

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

6 **ORCID IDs:** [0000-0002-5837-633X](https://orcid.org/0000-0002-5837-633X); [0000-0002-8370-4614](https://orcid.org/0000-0002-8370-4614)

7

8

9 Running Head: Comparing variation within datasets

10 **ABSTRACT**

11 Many key questions in evolutionary ecology require the use of variance ratios such as  
12 heritability, repeatability, and individual resource specialization. These ratios allow to  
13 understand how phenotypic variation is structured into genetic and non-genetic  
14 components, to identify how much organisms vary in the resources they use or how  
15 functional traits structure species communities. Understanding how evolutionary and  
16 ecological processes differs among populations and environments therefore often requires  
17 the comparison of these ratios across groups (i.e. populations, sexes, species). Inference  
18 based on comparisons of ratios can be limited, however. Variance ratios can remain the  
19 same across group despite very different values in the numerator and denominator  
20 variances. Moreover, evolutionary ecologists are most often interested in differences in  
21 specific variance component among groups rather than in differences in variance ratios *per*  
22 *se*. Recommendations for how to infer whether groups differ in variance are not clear in the  
23 literature. Using simulations, we show how questions regarding the estimation of variance  
24 components and their differences among groups can be answered with Linear Mixed  
25 Models (LMMs). Frequentist and Bayesian frameworks have similar abilities to identify  
26 differences in variance components. However, variance differences at higher levels of  
27 organization can be difficult to detect with low sample sizes. We provide tools to conduct  
28 power analyses to determine the appropriate sample sizes necessary to detect differences  
29 in variance of a given magnitude. We conclude by supplying guidelines for how to report  
30 and draw inferences based on the comparisons of variance components and variance ratios

31 **SIGNIFICANCE STATEMENT**

32 Many critical questions in ecology and evolution use variance ratios, such as repeatability,  
33 heritability, or individual resource specialization, to make inferences about ecological and  
34 evolutionary processes. In many cases these inferences rely on the comparison of variance  
35 ratios among datasets (populations, sexes, or environments). In this article, we show that  
36 current approaches of drawing inferences about group differences from comparisons of  
37 ratios are inappropriate because ratios can differ due to differences in the numerator,  
38 denominator, or both. We investigated how questions regarding differences in variance  
39 ratios and constituent variance components can be evaluated using Linear Mixed Model  
40 approaches (LMMs) and provide guidance for appropriate sampling schemes under  
41 different scenarios and discuss common pitfalls associated with estimation of differences in  
42 variance component among datasets.

43

44 Running Head: Comparing variation within datasets

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

47 **Declarations**

48 **Funding**

49 This study was funded by NSF IOS-1557951 (to NAD) and the Department of Biological  
50 Sciences at North Dakota State University.

51 **Conflicts of interest/Competing interests**

52 The authors declare no conflict of interest

53 **Availability of data and material**

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

56 **Code availability**

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

59 **Author contribution**

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

61 **Ethics approval**

62 Not applicable

63 **Consent to participate**

64 Not applicable

65 **Consent for publication**

66 Not applicable

67

## 68 INTRODUCTION

69 Our understanding of many evolutionary and ecological processes is underpinned  
70 by an estimation of variance ratios (Table 1). For example, the reporting of repeatability  
71 has become pervasive in behavioral studies as it summarizes the amount of variation in  
72 behavior attributable to differences among individuals. Informally these differences among  
73 individuals can be thought of as differences in their average behaviors. Repeatability then  
74 can be interpreted as how much of the overall variation is attributable to individual  
75 differences

76 Use of variance ratios like repeatability spans a broad swath of evolutionary ecology  
77 (Table 1). This includes the most well-known variance standardized ratio: heritability, and  
78 extends to interest in community ecology regarding the distribution of functional trait  
79 variation expressed within versus among populations or species (Violle et al. 2012).

80 While useful for understanding the relative magnitude of variation, variance ratios  
81 can be highly misleading when compared between groups (Houle 1992; Wilson 2018).  
82 Comparisons of variance ratios are only narrowly interpretable because these ratios can  
83 differ when numerators differ, when denominators differ, or when both differ. Indeed,  
84 variance ratios can be equal despite having different numerators and denominators values.  
85 Put another way, differences between groups in ratios like repeatability are not  
86 informative as to absolute differences in the magnitudes of variation observed.

87 **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>Adjusted Repeatability</i>	$R_A = Vi / (Vi + Vw)$	The proportion of variation attributable to among-individual differences ( $Vi$ ) relative to either the total variation ( $Vi + Vf + Vw$ ) or after adjusting for fixed-effects ( $Vi + Vw$ )	Lessels & Boag 1987
	<i>Unadjusted Repeatability</i>	$R_U = Vi / (Vi + Vf + Vw)$		
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 ( $V_{IR}$ )	

88

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

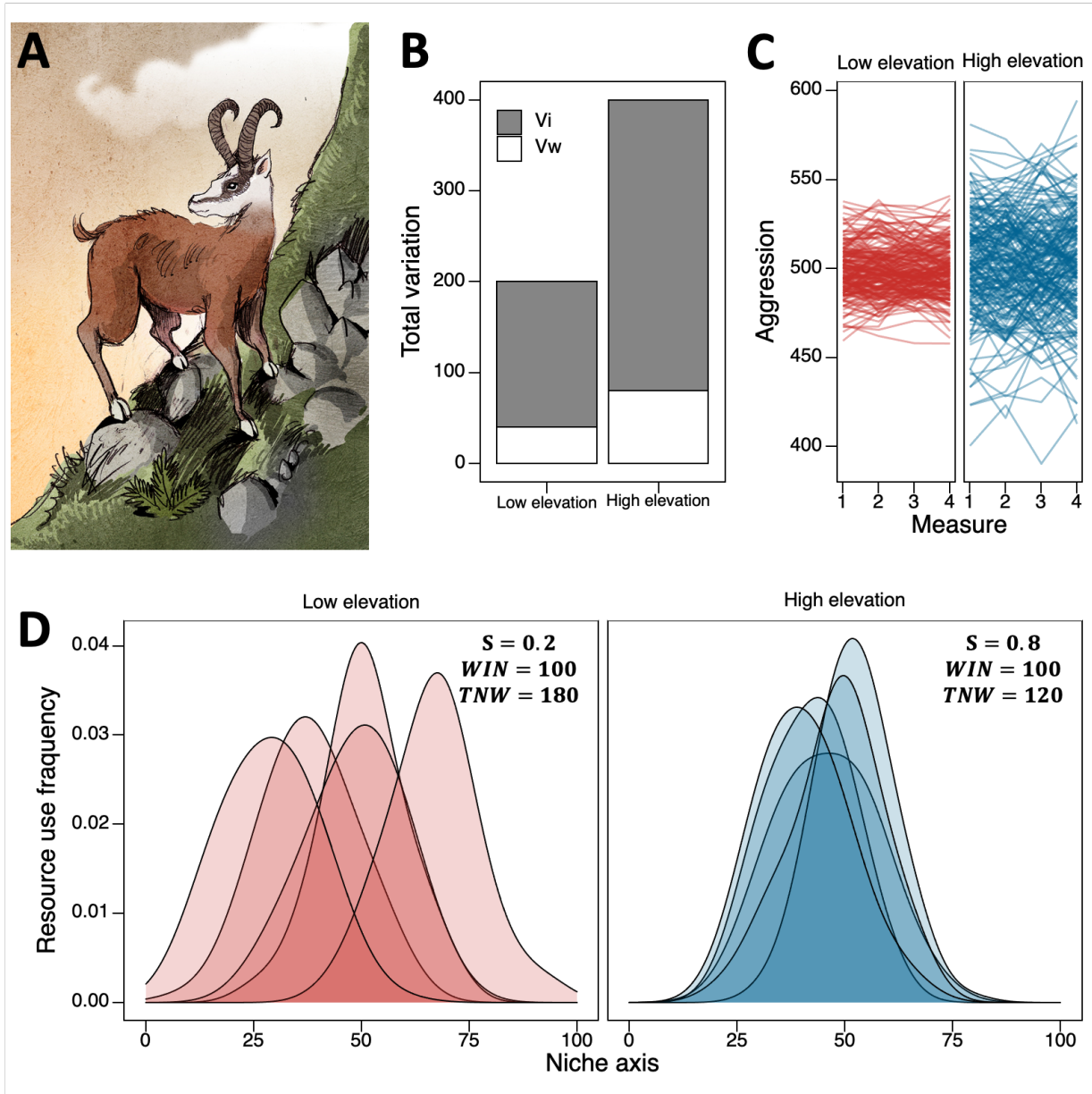
92           To further illustrate the inferential limits of variance ratios, consider the following  
93 scenario: researchers are studying the behaviors and dietary habits of two populations of  
94 the mythical Dahu (*Dahu desterus*; Figure 1A) at different elevations. These elusive  
95 creatures have shorter hind-legs on their left side, thus only allowing for clockwise  
96 movement (Chartois & Claudel 1945; Jacquat 1995). While measuring aggressive  
97 interactions, researchers find no differences in means between populations and similar  
98 behavioral repeatabilities ( $\tau = 0.8$ ; Figure 1B). Put another way, the same relative amount  
99 of variation is attributable to individuals in each population. The researchers notice,  
100 however, that there are large differences in the among-and within-individual variances of  
101 each population. Had researchers only examined repeatabilities and mean differences they  
102 would inappropriately conclude that the populations are behaviorally equivalent. Instead,  
103 the actual variance estimates reveal that individuals from the high-altitude population are  
104 very distinct from one another in their aggressive tendencies while, at low-altitude,  
105 individuals show little departure from the population average (Figure 1B, C).

106           These researchers are also curious as to whether the harsher climate at the top of  
107 the mountain range leads to a narrower dietary breadth. Researchers predict that  
108 individual resource specialization will be higher in the low elevation population, as *D.*  
109 *desterus* have more food options to choose from. To the researcher's surprise, they find  
110 much higher individual resource specialization in the high-altitude population:  $S_1 = 0.2$ ,  $S_2 =$   
111  $0.8$ . Upon examining the specific values of among- and within-individual variation in niche,  
112 they find that these differences are a result of the high elevation population having a much  
113 narrower total niche width (Figure 1D) while the within-individual variation in niche  
114 preference is equal between populations. This means that it is the difference in diet

115 preference among individuals that drives the difference between the two populations. With  
116 more varied resources available at low elevation, each individual can specialize along the  
117 total niche axis, yet the breadth of diet preference within-individuals is the same between  
118 populations.

119           For both traits, exclusive reliance on ratios would have led to either inappropriate  
120 or incomplete inferences (i.e. inappropriately concluding behavioral equivalence and  
121 incompletely recognizing the basis of differences in apparent specialization). Due to these  
122 problems with interpretations of variance ratios (Houle 1992; Wilson 2018; Dochtermann  
123 & Royauté 2019), what would be of greater use to researchers is to instead evaluate  
124 differences in specific variance components.





125  
 126 **Figure 1.** Reliance on variance ratios can lead to misleading inferences. (A) The elusive Dahu (*Dahu*  
 127 *dexterus*) in its natural environment. (B) Two populations of Dahus living at different elevations do  
 128 not differ in their repeatability of aggressive interactions ( $\tau$ ). (C) By plotting the individual  
 129 aggression scores over the course of multiple measurements, it is clear that individuals are more  
 130 distinct in their aggressive behavioral strategies at high elevation. This inference cannot be made by  
 131 investigating repeatability alone. (D) The two population have very different resource  
 132 specialization indices ( $S$ ). A more accurate inference is that individuals do not differ in niche width  
 133 ( $WIN$ ), it is instead the total niche width ( $TNW$ ) that is narrower in the high-altitude population.  
 134 Figure code available here: <https://osf.io/5aw42/>

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

137 The statistical procedures necessary for the estimation of variance components and ratios  
138 within a single population have been the subject of much attention (e.g. mixed models for  
139 repeatability: Dingemanse and Dochtermann 2013; animal models for heritability: Wilson  
140 et al. 2010; individual niche specialization: Bolnick et al. 2002; Coblentz et al. 2017;  
141 functional trait variation: Nakagawa and Schielzeth 2012; Violle et al. 2012; Carmona et al.  
142 2016). There is also a long history in quantitative genetics regarding the comparison of  
143 variances and *covariance* structures among groups (Shaw 1991, Arnold & Phillips 1999,  
144 Roff 2002, Roff et al. 2012, Aguirre et al. 2014). Unfortunately, these quantitative genetic  
145 approaches have been poorly disseminated across fields (but see Dochtermann & Roff  
146 2010 and White et al. 2019). Here we describe and investigate methods for detecting  
147 differences in variance components amongst groups. Specifically, we compare the strength  
148 and weaknesses of three statistical approaches: comparison of confidence intervals, model  
149 comparison with AIC, and Bayesian estimation of the difference in variance components.  
150 While this selection of methods encompasses very different philosophical approaches to  
151 data analysis, all three are routinely used in the estimation of repeatability and other ratios.

152 We consider a scenario where a phenotypic attribute,  $y$ , is measured repeatedly for  
153 individual organisms occupying one of two different environments (E1 and E2) and in  
154 which variation occurs among and within-individuals units ( $V_I$  and  $V_W$  respectively). In the  
155 following sections we focus on differences in individual variation and repeatability. Note,  
156 however, that this scenario can also be expanded to the comparison of diet specialization  
157 for individuals occupying different environments or how functional traits vary among and  
158 within species in two different environments.

159           An easy way to compare these variance components and their ratios ( $\tau = V_I/(V_I +$   
160  $V_W)$ ) is to estimate the variance components for each environment in separate statistical  
161 models. We can then test for differences in variances and ratios by environment based on  
162 whether estimate confidence intervals overlap or not. While straightforward, this method  
163 suffers from two key limitations. First, basing inference on the overlap of 95 % confidence  
164 intervals is overly conservative (Barr 1969), especially when sample size is low. It is  
165 instead whether the confidence interval for the *difference* in variances excludes 0 that is  
166 relevant for drawing inferences. This difference cannot be directly estimated from the  
167 approach we have described. However, statistical significance can still be assessed by  
168 comparing the overlap of the 83% confidence intervals for variance components, a  
169 threshold that provides a better approximation for an  $\alpha = 0.05$  for the null hypothesis of no  
170 difference (Austin and Hux 2002; MacGregor-Fors and Payton 2013; Hector 2015; see also  
171 Schenker & Gentleman 2001 for additional caveats). Second, by estimating variance  
172 components in separate statistical models, the hierarchical structure of the data, i.e. the  
173 variance components nested within the environments, has been broken. As a result,  
174 potential average differences in the traits of interest are not appropriately tested.

175           Instead, we suggest that a more appropriate procedure would be the use of a Linear  
176 Mixed Model (LMM) where the among- and within-individual variance is estimated for  
177 each environment within the same statistical model. This statistical model can be described  
178 by the following equation:

179

180  $y_{ij} = \beta_0 + \beta_1 \text{Environment} + ID_{0i} + e_{0ij}$  (equation 3)

181  $ID_{0i} \sim MVN(0, \Omega_{ID}); \quad \Omega_{ID} = \begin{bmatrix} V_{ID0} E_1 & 0 \\ 0 & V_{ID0} E_2 \end{bmatrix}$

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

183 where  $y_{ij}$  describes the phenotypic traits for the  $i$ th individual and  $j$ th observation.  $ID_{0i}$ , is  
 184 the deviation from an overall intercept,  $\beta_0$ , for the  $i$ th individual.  $\beta_1$  represents the  
 185 regression coefficient for the fixed effect of environment (here a contrast coefficient). The  
 186 random intercepts and residual variance ( $e_{0ij}$ ) both follow a multivariate normal  
 187 distribution, and  $\Omega_{ID}$  and  $\Omega_e$ , are the variance-covariance matrices at the among- and  
 188 within-individual levels respectively.

189 The diagonal elements of these matrices represent the among- and within-  
 190 individual variances in each environment:  $E_1$  and  $E_2$ . The off-diagonal elements represent  
 191 the cross-environment correlation (set to 0 if individuals are only ever evaluated in one of  
 192 the two environments). This formulation has the advantage of allowing considerable  
 193 flexibility in the specification of the statistical models considered (Dingemans and  
 194 Dochtermann 2013). LMMs are now available for most statistical software and their  
 195 generalized extensions can accommodate non-normal error distributions (Table 2).

196 Upon fitting LMMs, several methods are then available to determine whether a  
 197 variance ratio or components of the ratio differ by environment. Specific hypotheses of  
 198 which variance component differs across environment can be easily tested via model  
 199 comparison. For example, a model where only the among-individual variance differs by  
 200 environment can be compared to a null model where the among and within- individual  
 201 variance are kept constant across developmental environments (Royauté et al. 2019).

202 These models can be estimated within a frequentist framework via restricted maximum  
203 likelihood or a Bayesian framework and suitable decision criteria can be used to determine  
204 which model best fits the data. In the case of restricted maximum likelihood estimation, it is  
205 also possible to use likelihood ratio tests to compare these models. Note however that the  
206 proper degrees of freedom to apply to each model is unclear and additional care should be  
207 taken when using this method (Pinheiro and Bates 2000; see Santostefano et al. 2016 for a  
208 recent example). We recommend calculating these degrees of freedom by considering each  
209 variance component as a full parameter for more conservative testing (see also the tutorial  
210 in SI3).

211 In many cases, researchers are also interested in whether the difference in variance  
212 components have a biologically meaningful effect. In other words, when asking questions  
213 about whether variance components vary between environments, we are mostly interested  
214 in the *magnitude of the difference* in these components across environments. While model  
215 comparison of LMMs can help us understand whether a statistically detectable difference is  
216 observable across environments, the magnitude of the difference can only be determined  
217 by examining the difference in variance components among environment:  $\Delta V$  estimated as  
218  $V_{E2} - V_{E1}$  in our case. When the trait of interest is expressed as standard deviation units (i.e.  
219 mean centered and scaled to the standard deviation of the dataset across all populations  
220 and environments), this difference can be considered an effect size for the magnitude of the  
221 difference among variance components, thus making comparisons across studies possible  
222 (Royauté et al. 2015; Hamilton et al. 2017; Royauté and Dochtermann 2017). Note that  $\Delta V$   
223 could also be expressed on a ratio scale ( $V_{E2}/V_{E1}$ ) or on a log-additive scale ( $\log(V_{E2}) - \log$   
224 ( $V_{E1}$ )). We will return to the topic of statistical significance vs. appropriate effect sizes later

225 in the paper. For now, we simply consider  $\Delta V$  on an additive scale with data expressed in  
226 standard unit deviations because it allows the most straightforward interpretation and  
227 functions in cases where a variance component is zero or approaching zero.

228 **Table 2.** Packages and softwares allowing to test for differences in variance components using Linear Mixed Models (LMM) along with  
 229 parameter estimation method (maximum likelihood (ML), restricted maximum likelihood (REML), hierachical likelihood (H-ML) or  
 230 Bayesian framework) and inference method (Likelihood Ratio tests (LRT), AIC, bootstrapping or credible interval for  $\Delta V$ ). This list is not  
 231 comprehensive and is instead based on widely-used commercial softwares and R packages.

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, bootstrapping	Yes	Yes	Gaussian		Gilmour et al. (2015)
SAS	Commercial	ML/REML	LRT, AIC, bootstrapping	Yes	Yes	Gaussian, Poisson, Binomial		SAS Institute Inc.
nlme	Free	ML/REML	LRT, AIC, bootstrapping	Yes	Yes	...		Pinheiro and Bates (2000)
lme4	Free	ML/REML	LRT, AIC, bootstrapping	Yes	No	Gaussian, Poisson, Binomial		Bates et al. (2015)
glmmTMB	Free	ML/REML	LRT, AIC, bootstrapping	Yes	Yes	...		Brooks et al. 2017
hglm	Free	H-ML	LRT, AIC, bootstrapping	Yes	Yes	...	Within-unit variance modelled as Gamma distribution	Rönnegård et al. (2010)
R-INLA	Free	Approximate Bayesian	credible intervals for $\Delta V$	Yes	Yes	...		Lindgren, and Rue (2015)
MCMCglmm	Free	Bayesian	DIC, credible intervals for $\Delta V$	Yes	Yes	Gaussian, Poisson, Binomial		Hadfield (2010)
brms	Free	Bayesian	WAIC, LOO, credible intervals for $\Delta V$	Yes	Yes	...	Within-unit variance modelled	Bürkner (2017)





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

240 In summary, approaches based on LMM and their generalized extensions allow  
241 great flexibility and are well suited to study questions related to how variation in  
242 phenotypic traits varies at multiple levels of organization. In the next section, we describe  
243 the performance of LMMs to detect differences in variance components.

244

## 245 **METHODS**

246 The simulations described below focus on interpretation in the context of behavioral  
247 repeatability. However, it is worth noting again that inferences about the ability to estimate  
248 and detect differences in variances generalizes to the components of the ratios described in  
249 Table 1.

### 250 *Data simulations*

251 To compare the performance of statistical procedures for detecting differences in variance  
252 components and variance ratios, we performed a series of simulations based on the  
253 scenarios illustrated in Figure 2. In these scenarios a phenotypic attribute  $y$  is measured in

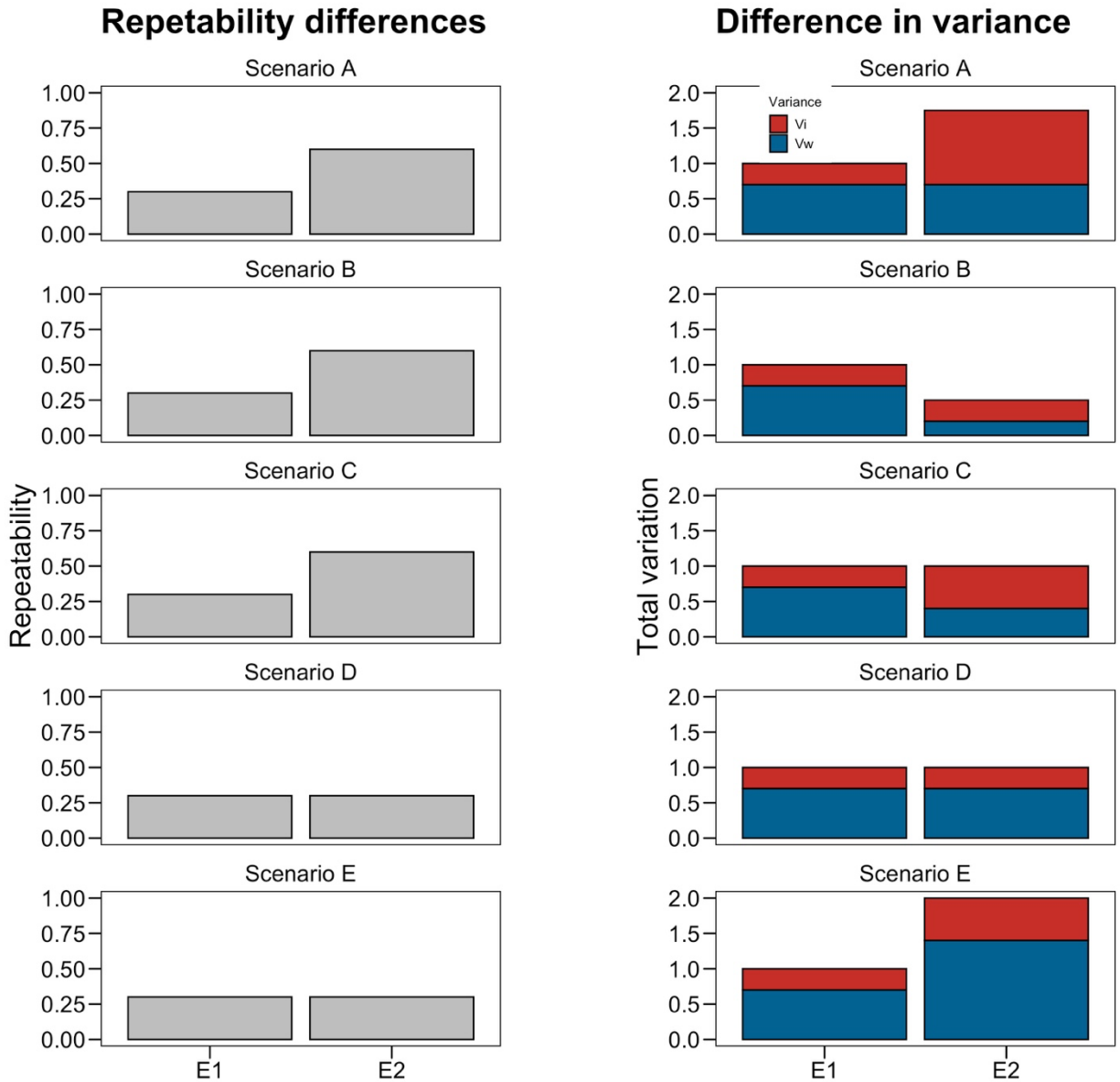
254 two different environments (E1 and E2) and variation occurs among and within individuals  
255 ( $V_I$  and  $V_W$  respectively). In scenarios A through C the repeatability ( $\tau$ ) differs by an equal  
256 amount between the two environments ( $\Delta\tau = 0.3$ ), but the underlying driver of this  
257 difference is either due to a difference in the among-individual variance (A), in the within-  
258 individual variance (B) or in both the among and within-individual variance (C). Note that  
259 for scenario C, the total variance remains the same between environments. In scenarios D  
260 and E, we explore cases where the variance ratios are equal among environment, either  
261 because all variance components are equal as well (D) or in spite of differences in all  
262 variance components (E) (see Table S1 for exact values for all parameters).

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

- 265 • Sample size varying from 20 to 200 individuals by increments of 20 for each  
266 environment (sample size was equal between the two environments)
- 267 • Number of repeated measures taken on each individual varying from 2 to 6  
268 repeated measures by increments of 1
- 269 • Five different scenarios of known difference in variance ratios as described in  
270 Figure 1 and Table S1.

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

276



277

278 **Figure 2.** Scenarios used in simulations detailing how differences or lack of difference in  
 279 repeatability (right-side column) can arise from different patterns in the underlying  
 280 variance components (left-side column; exact values can be found in Table S1). Scenarios  
 281 A-C correspond to cases where the total variation differs between two environments (E1  
 282 and E2) due to differences in the higher group level variance ( $V_i$ , A), the lower level  
 283 variance ( $V_w$ , B) or both (C). Scenarios D-E indicate cases where the ratios remains  
 284 constant across environments, because all variance components are identical (D) or in  
 285 spite of variance component being different among environments (E).

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

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

291  $y_{ij} = \beta_0 + ind_{0j} + e_{0ij}$  (equation 4)

292  $ind_{0j} \sim \mathcal{N}(0, V_{ind});$

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

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

298 *Frequentist LMM with AIC model comparison*

299 Our second approach was to fit the LMM approach described above and test for the for the  
300 significance of the difference in among- and within-individual variance using likelihood  
301 ratio tests. Specifically, we specified four different mixed models corresponding to the four  
302 different possibilities by which variance components may differ (see also Royauté et al.  
303 2019 ; Buckleaw and Dochtermann 2021):

- 304 • Model 1: a null model where the among ( $V_I$ ) and within-individual variance ( $V_W$ )  
305 was kept constant among environments.

- 306 • Model 2: a model where only the among-individual variance differs among  
307 environments, while the within-individual variance is kept constant ( $V_I \neq$  &  $V_W =$ )
- 308 • Model 3: a model where only the within-individual variance differs among  
309 environments while the among-individual variance is kept constant ( $V_I =$  &  $V_W \neq$ )
- 310 • Model 4: a model where both the among and within-individual variance were  
311 allowed to vary among environments ( $V_I \neq$  &  $V_W \neq$ )

312 For each dataset combination, we then compared each model's Aikake's Information  
313 Criterion value (AIC). AIC allows to compare the relative fit of statistical models and models  
314 with lower AIC values indicate better support relative to competing models. These  
315 simulations and this analytical framework are similar to previously used approaches (e.g.  
316 Jenkins 2011; Shaw 1991; Tüzün et al. 2017). These models were specified using the *nlme*  
317 package for mixed models (Pinheiro et al. 2000) using Restricted Maximum Likelihood  
318 (REML).

### 319 *Bayesian LMM and difference in variance components*

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

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

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

332 We calculated the probability of detecting the model with the correct difference in variance  
333 components (hereafter “abridged” to probability of correct model identification), precision,  
334 relative bias and accuracy under each scenario and sampling design to compare the  
335 performance of maximum likelihood and Bayesian mixed models. For Method 1 (overlap of  
336 83 % intervals), we assigned values of 1 when significant differences in variance  
337 components were detected in directions predicted by the data generating process, and 0  
338 otherwise. For Method 2, we calculated the probability of correct model identification as  
339 the proportion of times the model with the lowest AIC matched the generating model. For  
340 Method 3, we calculated whether a given model detected a difference in variance  
341 components based on the overlap of the 95 % credible intervals of the  $\Delta V$  posterior  
342 distribution with 0. As in Method 1, we then assigned values of 0 or 1 based on whether the  
343 detected difference matched with the data generation process of the corresponding  
344 scenario. We calculated the probability of correct model identification as the proportion of  
345 analyzed datasets in which we detected differences in the direction predicted by each  
346 scenario and statistical method. Precision, indicating the similarity of the results produced  
347 by simulations with a given scenario, was calculated as the difference between 25 % and 75  
348 % quantiles of estimates (van de Pol 2012). To calculate the relative bias (in %) for each  
349 statistical approach by scenario, we calculated the mean difference between the expected

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

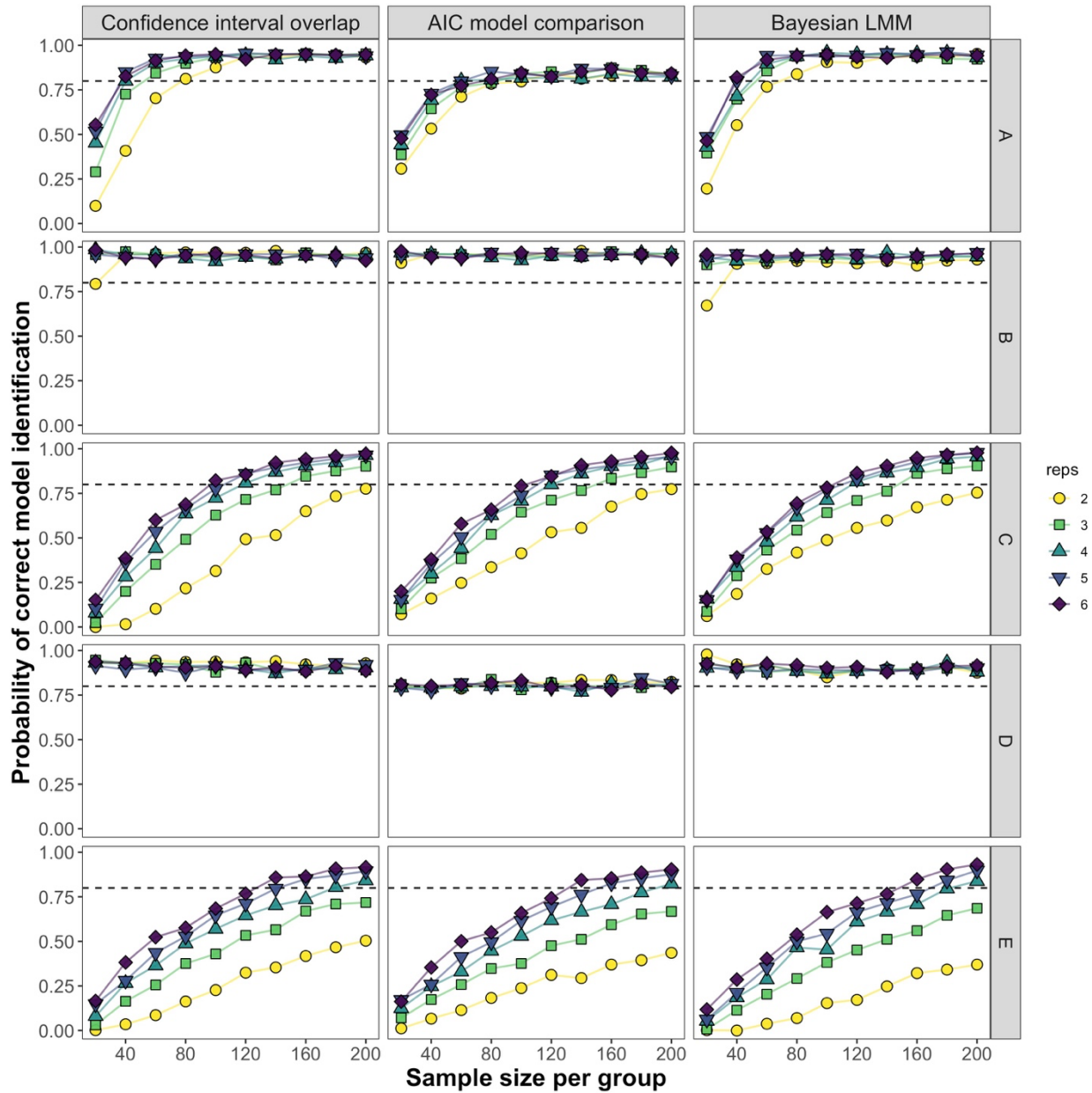
354

## 355 **RESULTS**

356 The probability of correctly detecting differences in variance components did not differ  
357 substantially between frequentist and Bayesian methods of estimation (Figure 3). The  
358 highest probability of correct model identification was observed for in cases where the  
359 variance ratio differs as a result of changes to the within-individual variance (scenario B)  
360 or when variation remained equal between environments (scenario D). The statistical  
361 power to differentiate between alternative scenarios (i.e. scenarios A, C and E) was lower,  
362 especially with small sample sizes and low number of repeated measures (Figure 3).

363 Importantly, no statistical method seemed to outperform all others across scenarios. Our  
364 results are consistent with previous simulations showing that the among-individual  
365 variance component is particularly difficult to estimate at small sample sizes (Dingemanse  
366 & Dochtermann 2013).

367



368

369 **Figure 3.** Effect of sampling design on the probability of correct model identification by  
 370 scenario type and statistical modeling approach. Each point represents the probability of  
 371 detecting the correct differences in variance averaged over 500 simulated datasets. A  
 372 represents a scenario where only the among-individual variance ( $V_I$ ) varies between  
 373 environments, B represents a case where the within-individual variance ( $V_W$ ) varies  
 374 between environments, and both among and within- individual variance vary between  
 375 environments in scenario C. In scenario D, all variance components are equal while in  
 376 scenario E, variance components are different but variance ratios are equal across  
 377 environments. Dashed lines correspond to 80 % threshold similar to recommendations for  
 378 power analyses.



379 In scenarios B and D, the correct differences among variance components was  
380 identified > 80 % of the time, even at low sample sizes (Figure 3). In all other scenarios this  
381 threshold was only reached with high sample sizes and a high number of repeated  
382 measures. For scenarios C and E—which correspond to cases where the variance ratio  
383 differs as a result of among-individual variance (C) or when the variance ratio remains the  
384 same despite changes to both among- and within-individual variance (E)—datasets with  
385 only 2 repeated measures per individual never achieved a probability of identifying the  
386 generating model above 0.8, even with sample sizes above 200 units per environment (i.e. a  
387 minimum of 800 total measurements, Figure 3). Increasing the number of repeated  
388 measures only marginally alleviated the problem. For example, in scenario C, only datasets  
389 with 4 or more repeated measures per individual reached statistical power above 0.8 with  
390 sample sizes above 120 units per environments, which is higher than many ecological or  
391 evolutionary studies can provide under realistic scenarios.

392 Note that for AIC model comparison, we calculated power as the number of times  
393 the best model corresponded to the generating model. A more conservative approach is to  
394 calculate the proportion of times the best model is at least 2 AIC units lower than the  
395 second model. This method corresponds to a common threshold to detect statistically  
396 distinct models (Burnham and Anderson 1998). When using this more conservative  
397 threshold (Figure S1), datasets generated according to scenarios A and D were never  
398 statistically distinguishable from non-generating models, although the correct model was  
399 consistently ranked as the best model. This discrepancy is likely because when the  
400 generating model does not include differences in the within-individual variability  
401 (scenarios A and D), sampling error is erroneously identified as heterogeneity. At smaller

402 sample sizes this error is greater on average, and thus detectable. At larger sample sizes  
403 this sampling error is smaller but more easily detected and therefore manifests as a  
404 difference between groups. To address this, in addition to measures of variance differences  
405 like the described  $\Delta V$  statistic, researchers should also compare mean-standardized  
406 variance estimates like the coefficient of variation or Houle's evolvability between groups  
407 (Houle 1992; Hansen et al. 2011; Dochtermann and Royauté 2019).

408         The comparison of relative bias, precision, and accuracy among statistical methods  
409 produced mixed results. On average, Bayesian LMMs consistently underestimated the  
410 among-individual variance for scenarios in which the among-individual variance differed  
411 between environments (scenarios A, C, and E) resulting in a bias at small sample sizes  
412 (Figure S2). However, Bayesian LMMs also had higher precision and accuracy compared to  
413 maximum likelihood (Figure S3, S4). This means that Bayesian estimates tend to be  
414 consistently more conservative than maximum likelihood regarding the magnitude of the  
415 among-individual variance but that these estimates nonetheless more closely matched  
416 simulation conditions.

## 417 **DISCUSSION**

418 Comparing variability across datasets is important for many questions in evolutionary  
419 ecology (e.g. Table 1). However, variance ratios are not sufficient to address questions  
420 about how variation is expressed across environments, populations, or sexes. The inability  
421 to determine why groups differ based on ratios is in addition to the numerous conceptual  
422 and theoretical problems inherent to the estimation of ratios (Houle 1992; Hansen et al.  
423 2011). Instead, many questions require the direct comparison of variances.

424 *What are appropriate sample sizes for detecting differences in variance?*

425           Our simulations show that regardless of the statistical methods used, comparing  
426 variance components across groups is a “data hungry” question. Scenarios where the  
427 among-individual variance differed between environments were particularly hard to detect  
428 at low sample sizes. Note that our objective was not to provide a full exploration of  
429 parameter space. Instead, we focused on a subset of scenarios that are likely to be common  
430 in ecology and evolution (Figure 2). Based on our simulations, the probability to detect  
431 differences in variance components will depend in large part on the ability to estimate the  
432 among-individual variance component ( $V_I$ ). In the most complex case where differences  
433 occur among and within-individuals (scenario E), researchers would require a minimum of  
434 1,600 observations to correctly detect differences (i.e, 200 individuals measured 4 times in  
435 each environment). This is far higher than sample sizes needed for single populations,  
436 where moderate repeatabilities only need ~100 observations to be estimated with > 0.8  
437 power (at least 25 individuals measured 4 times to detect a repeatability of 0.3; see  
438 Dingemanse & Dochtermann 2013).

439           Given these challenges, we recommend that researchers conduct power calculations  
440 prior to the experiment whenever possible (see R code for *a priori* power analyses in SI2  
441 and an R Markdown tutorial in SI3). If not, a simple rule for sampling can be to estimate the  
442 sample size needed to detect the lowest among-individual variance value of interest (see,  
443 for example, Martin et al. 2011; van de Pol 2012; Dingemanse and Dochtermann 2013) and  
444 multiplying that sample size by the number of experimental groups involved.

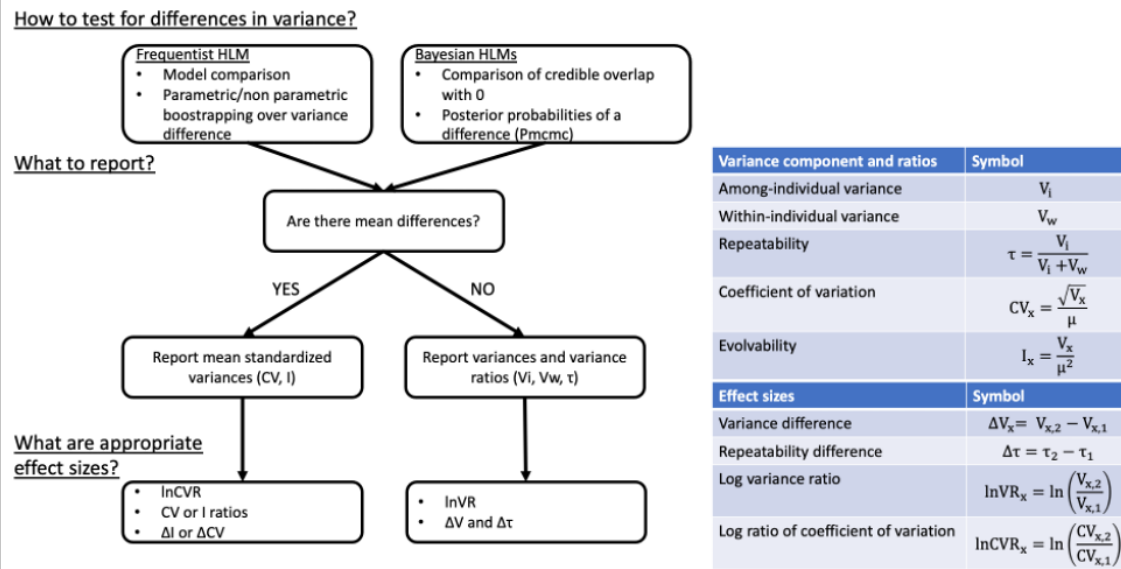
445 *How to report results? Statistical significance vs. effect sizes*

446           Given the issues discussed above, how should researchers interested in ecological  
447 and evolutionary variation design their studies and report their findings? We suggest that  
448 researchers report their results in a manner that focuses on the magnitude of the difference  
449 in variability between experimental groups rather than solely focus on statistical  
450 significance.

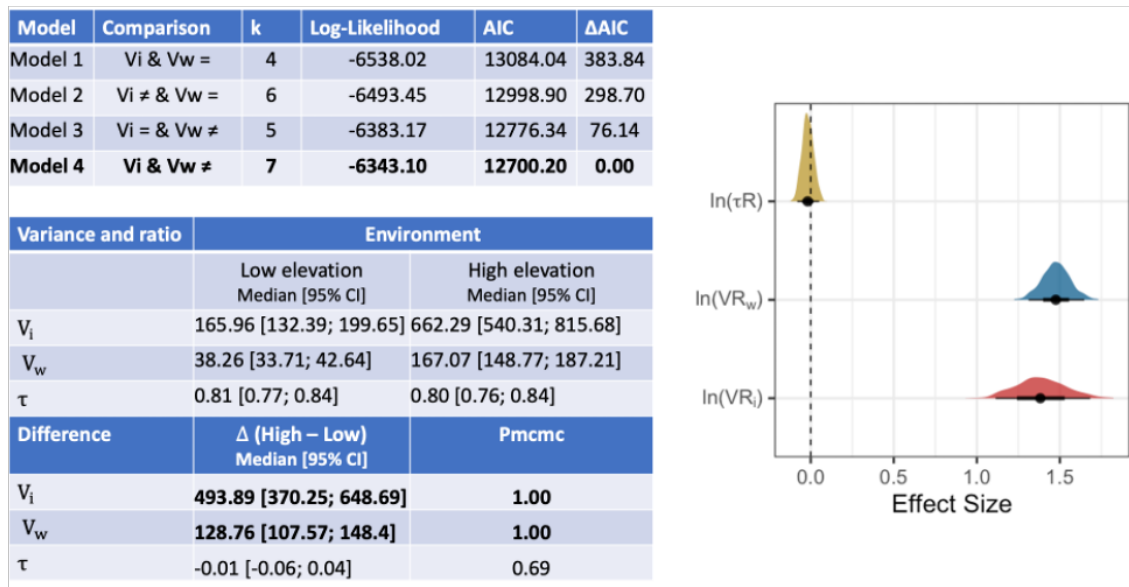
451           To this effect, we believe that reporting the results of the full model rather than just  
452 the most parsimonious model will be most appropriate in most cases (i.e. model 4 in our  
453 conceptual example). This is because model selection only gives information on whether  
454 differences among groups are statistically detectable. In contrast, questions regarding the  
455 magnitude and precision of the estimated differences are answerable only with  
456 interpretation of the most complete statistical model (see tutorial in ESM4).

457           In addition to presenting results of the full model, we suggest that measures of effect  
458 sizes for the differences in variance component also be presented. As reported above,  $\Delta V$   
459 provides a simple metric to estimate the magnitude of these differences, but it is by no  
460 mean the only one. In our theoretical example, the mean trait value did not differ by  
461 environments, but in many cases mean and variance are related. In such cases, using  
462 comparisons based on Houle's (1992)  $I^2$  value or coefficients of variation for each  
463 component as opposed to variance component themselves can be preferable (Hansen et al.  
464 2011; Dochtermann and Royauté 2019). Effect sizes based on the coefficient of variation  
465 can also be calculated within an LMM framework as described by Nakagawa et al. (2015)  
466 (see also Carmona et al. 2016 and Fontana et al. 2018 for approaches relevant to functional  
467 trait diversity).

## A



## B



468

469

470

471

472

473

474

475

476

477

478

479

480

**Figure 4.** A) Flowchart showing decision rules regarding how to test for differences in variance components, which metrics to report and which effect sizes can be calculated, along with their definitions in table format. B) Reporting example based on the simulated case study in Figure 1B-C. The first table used REML model selection with AIC to compare the support for different hypotheses for how variance components of aggression may differ between the low and high elevation populations. The best model is one where among and within-individual variances are higher in the high elevation population. The second table compares all components by environment (posterior medians and 95 % credible intervals estimated from a Bayesian mixed model with model 4, note that frequentist confidence interval can also be reported using non-parametric bootstrapping as shown in SI3). Finally, because aggression does not differ on average between populations, InVR is an appropriate metric to report the effect size for the difference in variance between populations.

481 We provide a synthetic guide for which statistical tests and effect sizes are most  
482 appropriate depending on the nature of the dataset in Figure 4A. Returning to our dahu  
483 example, an appropriate analysis of the difference in aggression variance would follow the  
484 tables and figures from Figure 4B. Here the repeatability is unchanged between  
485 environments (posterior median [95 % credible interval];  $\Delta\tau = -0.01 [-0.06; 0.04]$ ,  
486 probability of difference:  $P_{mcmc} = 0.68$ ). However, the high-elevation population shows  
487 significantly higher variation among and within-individuals ( $\Delta VR_I = 493.89 [370.25;$   
488  $648.69]$ ,  $P_{mcmc} = 1.00$ ;  $\Delta VR_W = 128.76 [107.57; 148.40]$ ,  $P_{mcmc} = 1.00$ ). This difference is  
489 also biologically relevant since the effect sizes are also  $> 1$  ( $\ln VR_I = 1.38 [1.10, 1.66]$ ;  $\Delta VR_W$   
490  $= 1.48 [1.31, 1.64]$ ). Biologically, this means that the high elevation population is composed  
491 of individuals that are more distinct in behavior compared to the low elevation population.

492 While we limited our conceptual example to comparisons between two  
493 environments, the LMM approach we propose is by no mean restricted to two-groups  
494 comparisons. For example, Jenkins (2011) used model comparison to tease apart the  
495 relative influence of sex, species and their interaction on the expression of behavioral  
496 variation in kangaroo rats. Similarly, Coblentz et al. (2017) show how model selection  
497 combined with Bayesian HGLM can allow the comparison of indices of diet specialization  
498 within and among species. In both cases, model selection can provide a first pass at  
499 whether differences in variance components are detectable among groups, while specific  
500 pairwise comparisons of effect sizes (using  $\Delta V$  or other metrics) will allow discernment of  
501 the most pronounced differences in variance component. Regardless of the statistical  
502 approach used, we suggest it is important that researchers clearly outline the direction and,  
503 when possible, magnitude of the expected effects in their predictions.

504 Finally, our conceptual examples focus exclusively on the case of “well-behaved”  
505 data with normal error distributions. While these comparisons can be made with  
506 generalized extensions to LMMS (i.e. GLMMs), researchers must take extra precautions  
507 when calculating and comparing the within-individual variances (i.e. the residual variance).  
508 Indeed, in the case of non-Gaussian data, the residual variance depends on both the link  
509 function used and how the software deals with overdispersion (additive vs. multiplicative  
510 overdispersion). Nakagawa & Schielzeth 2010 provides a very useful and extensive guide  
511 explaining how the correct residual variation can be calculated.

## 512 **CONCLUSIONS**

513 Variance ratios are straightforward metrics to describe how various ecological and  
514 evolutionary processes occur. However, comparing these ratios across studies or group can  
515 be misleading if poor attention is given to the specific variance components making up  
516 those ratios. More importantly, a lack of difference in these ratios does not mean that  
517 variation is expressed equally among groups. Given these limitations, we advocate for  
518 techniques allowing the estimation of differences in each variance components rather than  
519 focusing solely on variance ratios. The statistical tools allowing comparison of trait  
520 variation have become increasingly sophisticated and now allow asking very precise  
521 questions. Specifically, we can now ask how trait variation is generated and how variation  
522 differs among groups. However, despite the availability of these tools, researchers  
523 interested in ecological and evolutionary variation must remain careful in their study  
524 designs. As our simulations show, scenarios involving differences in among-individual  
525 variance are particularly difficult to detect without substantial sample sizes. Finally, we

526 hope the statistical approaches and tools for power analysis presented here will allow for  
527 appropriate comparisons of trait variation in ecological and evolutionary studies.

528

### 529 **Acknowledgments**

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

535

### 536 **Availability of data and material**

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

539

### 540 **Code availability**

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

543

### 544 **Author contribution**

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

546



547 **REFERENCES**

- 548 Aguirre, J., E. Hine, K. McGuigan, and M. Blows. 2014. Comparing **G**: multivariate analysis of  
549 genetic variation in multiple populations. *Heredity* 112:21-29.
- 550 Arnold, S. J., and P. C. Phillips. 1999. Hierarchical comparison of genetic variance-  
551 covariance matrices. II. Coastal-inland divergence in the garter snake, *Thamnophis*  
552 *elegans*. *Evolution* 53:1516-1527.
- 553 Austin, P. C., and J. E. Hux. 2002. A Brief Note on Overlapping Confidence Intervals. *Journal*  
554 *of Vascular Surgery* 36:194–195.
- 555 Barr, D. R. 1969. Using confidence intervals to test hypotheses. *Journal of Quality*  
556 *Technology* 1:256–258.
- 557 Bates, D., M. Maechler, B. Bolker, S. Walker, R. H. B. Christensen, H. Singmann, B. Dai, et al.  
558 2015. Package 'lme4.'
- 559 Bell, A. M., S. J. Hankison, and K. L. Laskowski. 2009. The repeatability of behaviour: a meta-  
560 analysis. *Animal behaviour* 77:771–783.
- 561 Bolnick, D. I., R. Svanbäck, J. A. Fordyce, L. H. Yang, J. M. Davis, C. D. Hulsey, and M. L.  
562 Forister. 2002. The ecology of individuals: incidence and implications of individual  
563 specialization. *The American Naturalist* 161:1–28.
- 564 Bucklaew, A. and N.A. Dochtermann. 2021. The effects of exposure to predators on  
565 personality and plasticity. *Ethology* 127:158-165.
- 566 Brooks, M. E., Kristensen, K., van Benthem, K. J., Magnusson, A., Berg, C. W., Nielsen, A.,  
567 Skaug, H. J., Machler, M. and B. M. Bolker 2017. glmmTMB balances speed and  
568 flexibility among packages for zero-inflated generalized linear mixed modeling. *The*  
569 *R Journal* 9:378-400.
- 570 Bürkner, P.-C. 2017. brms: An R package for Bayesian multilevel models using Stan. *Journal*  
571 *of Statistical Software* 80:1–28.
- 572 Burnham, K. P., and D. R. Anderson. 1998. Practical use of the information-theoretic  
573 approach. Pages 75–117 *in* *Model Selection and Inference*. Springer.
- 574 Carmona, C. P., F. de Bello, N. W. Mason, and J. Lepš. 2016. Traits without borders:  
575 integrating functional diversity across scales. *Trends in ecology & evolution* 31:382–  
576 394.
- 577 Chartois, J., & Claudel, C. 1945. Hunting the dahut: a french folk custom. *The Journal of*  
578 *American Folklore* 58:21-24.

- 579 Coblentz, K. E., A. E. Rosenblatt, and M. Novak. 2017. The application of Bayesian  
580 hierarchical models to quantify individual diet specialization. *Ecology* 98:1535–  
581 1547.
- 582 Dingemanse, N. J., and N. A. Dochtermann. 2013. Quantifying individual variation in  
583 behaviour: mixed-effect modelling approaches. *Journal of Animal Ecology* 82:39–54.
- 584 Dochtermann, N. A., and D. A. Roff. 2010. Applying a quantitative genetics framework to  
585 behavioural syndrome research. *Philosophical Transactions of the Royal Society B-  
586 Biological Sciences* 365:4013-4020.
- 587 Dochtermann, N. A., and R. Royauté. 2019. The mean matters: going beyond repeatability to  
588 interpret behavioural variation. *Animal Behaviour* 153:147–150.
- 589 Dochtermann, N. A., T. Schwab, M. Anderson Berdal, J. Dalos, and R. Royauté. 2019. The  
590 Heritability of Behavior: A Meta-analysis. *Journal of Heredity*.
- 591 Dochtermann, N. A., T. Schwab, and A. Sih. 2015. The contribution of additive genetic  
592 variation to personality variation: heritability of personality. *Proceedings of the  
593 Royal Society B: Biological Sciences* 282:20142201.
- 594 Fontana, S., M. K. Thomas, M. Moldoveanu, P. Spaak, and F. Pomati. 2018. Individual-level  
595 trait diversity predicts phytoplankton community properties better than species  
596 richness or evenness. *The ISME journal* 12:356.
- 597 Gilmour, A. R., B. J. Gogel, B. R. Cullis, S. J. Welham, and R. Thompson. 2015. ASReml user  
598 guide release 4.1 structural specification. Hemel Hempstead: VSN international ltd.
- 599 Hadfield, J. D. 2010. MCMC methods for multi-response generalized linear mixed models:  
600 the MCMCglmm R package. *Journal of Statistical Software* 33:1–22.
- 601 Hamilton, J. A., R. Royauté, J. W. Wright, P. Hodgskiss, and F. T. Ledig. 2017. Genetic  
602 conservation and management of the California endemic, Torrey pine (*Pinus  
603 torreyana* Parry): Implications of genetic rescue in a genetically depauperate  
604 species. *Ecology and Evolution* 7:7370–7381.
- 605 Hansen, T. F., C. Pélabon, and D. Houle. 2011. Heritability is not Evolvability. *Evolutionary  
606 Biology* 38:258.
- 607 Hector, A. 2015. *The New Statistics with R: An Introduction for Biologists*. 1<sup>st</sup> edition.  
608 Oxford ; New York, NY: Oxford University Press.
- 609 Houle, D. 1992. Comparing evolvability and variability of quantitative traits. *Genetics*  
610 130:195–204.

- 611 Jacquat, M. S. 1995. Le dahu: monographie ethno-étho-biologique publiée à l'occasion de  
612 l'exposition inaugurée le 1er avril 1995. Editions de la Girafe, Musée d'histoire  
613 naturelle.
- 614 Jenkins, S. H. 2011. Sex differences in repeatability of food-hoarding behaviour of kangaroo  
615 rats. *Animal Behaviour* 81:1155–1162.
- 616 Lessells, C. M., and P. T. Boag. 1987. Unrepeatable repeatabilities: a common mistake. *The*  
617 *Auk* 104:116–121.
- 618 Lindgren, F., and H. Rue. 2015. Bayesian spatial modelling with R-INLA. *Journal of*  
619 *Statistical Software* 63:1-25.
- 620 Lush, J. 1937. *Animal Breeding Plans*. Iowa State College Press, Ames, Iowa.
- 621 Martin, J. G., D. H. Nussey, A. J. Wilson, and D. Réale. 2011. Measuring individual differences  
622 in reaction norms in field and experimental studies: a power analysis of random  
623 regression models. *Methods in Ecology and Evolution* 2:362–374.
- 624 MacGregor-Fors, I., and M. E. Payton. 2013. Contrasting Diversity Values: Statistical  
625 Inferences Based on Overlapping Confidence Intervals. *PloS One* 8, no. 2.  
626 <http://dx.plos.org/10.1371/journal.pone.0056794>.
- 627 Mousseau, T. A., and D. A. Roff. 1987. Natural selection and the heritability of fitness  
628 components. *Heredity* 59:181.
- 629 Nakagawa, S., R. Poulin, K. Mengersen, K. Reinhold, L. Engqvist, M. Lagisz, and A. M. Senior.  
630 2015. Meta-analysis of variation: ecological and evolutionary applications and  
631 beyond. *Methods in Ecology and Evolution* 6:143–152.
- 632 Nakagawa, S., and H. Schielzeth. 2010. Repeatability for Gaussian and non-Gaussian data: a  
633 practical guide for biologists. *Biological Reviews* 85:935–956.
- 634 Nakagawa, S., and H. Schielzeth. 2012. The mean strikes back: mean–variance relationships  
635 and heteroscedasticity. *Trends in Ecology & Evolution* 27:474–475.
- 636 Pinheiro, J., and D. Bates. 2000. *Mixed-Effects Models in S and S-PLUS*. Springer Science &  
637 Business Media.
- 638 Roff, D. 2002. Comparing **G** matrices: A MANOVA approach. *Evolution* 56:1286-1291.
- 639 Roff, D. A., J. M. Prokkola, I. Krams, and M. J. Rantala. 2012. There is more than one way to  
640 skin a **G** matrix. *Journal of Evolutionary Biology* 25:1113-1126.
- 641 Rönnegård, L., X. Shen, and M. Alam. 2010. hglm: A package for fitting hierarchical  
642 generalized linear models. *The R Journal* 2:20–28.

- 643 Royauté, R., C. M. Buddle, and C. Vincent. 2015. Under the influence: sublethal exposure to  
644 an insecticide affects personality expression in a jumping spider. *Functional Ecology*  
645 29:962–970.
- 646 Royauté, R., and N. A. Dochtermann. 2017. When the mean no longer matters:  
647 Developmental diet affects behavioral variation but not population averages in the  
648 house cricket (*Acheta domesticus*). *Behavioral Ecology* 28:337–345.
- 649 Royauté, R., C. Garrison, J. Dalos, M. A. Berdal, and N. A. Dochtermann. 2019. Current energy  
650 state interacts with the developmental environment to influence behavioural  
651 plasticity. *Animal Behaviour* 148:39–51.
- 652 Santostefano, F., A. J. Wilson, Y. G. Araya-Ajoy, and N. J. Dingemanse. 2016. Interacting with  
653 the enemy: indirect effects of personality on conspecific aggression in crickets.  
654 *Behavioral Ecology* 27:1235–1246.
- 655 Schenker, N., and Gentleman, J. F. (2001). On judging the significance of differences by  
656 examining the overlap between confidence intervals. *The American Statistician*  
657 55:182-186.
- 658 Shaw, R. G. 1991. The comparison of quantitative genetic-parameters between populations.  
659 *Evolution* 45:143-151
- 660 Stirling, D. G., D. Réale, and D. A. Roff. 2002. Selection, structure and the heritability of  
661 behaviour. *Journal of Evolutionary Biology* 15:277–289.
- 662 Tüzün, N., S. Müller, K. Koch, and R. Stoks. 2017. Pesticide-induced changes in personality  
663 depend on the urbanization level. *Animal behaviour* 134:45–55.
- 664 van de Pol, M. 2012. Quantifying individual variation in reaction norms: how study design  
665 affects the accuracy, precision and power of random regression models. *Methods in*  
666 *Ecology and Evolution* 3:268–280.
- 667 Violle, C., B. J. Enquist, B. J. McGill, L. I. N. Jiang, C. H. Albert, C. Hulshof, V. Jung, et al. 2012.  
668 The return of the variance: intraspecific variability in community ecology. *Trends in*  
669 *ecology & evolution* 27:244–252.
- 670 White, S. J., Pascall, D. J., and A. J. Wilson. 2019. Towards a comparative approach to the  
671 structure of animal personality variation. *Behavioral Ecology*.
- 672 Wilson, A. J., D. Réale, M. N. Clements, M. M. Morrissey, E. Postma, C. A. Walling, L. E. B.  
673 Kruuk, et al. 2010. An ecologist's guide to the animal model. *Journal of Animal*  
674 *Ecology* 79:13–26.
- 675 Wilson, A. J. 2018. How should we interpret estimates of individual repeatability? *Evolution*  
676 *Letters* 2: 4-8.

677

678 **Supporting Information (SI) and Electronic Supplementary Materials (ESM)**

679 **ESM 1:** Zip folder containing the raw data from simulations along with R code for data  
680 analysis and figures (<https://osf.io/5aw42/>).

681 **ESM 2:** R code for conducting *a priori* power analysis (<https://osf.io/5aw42/>).

682 **ESM 3:** R tutorial for comparing variance components using *nlme*, *MCMCglmm* and *brms*  
683 packages (<https://osf.io/5aw42/>).

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

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

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

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

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