

<sup>1</sup> A Standardized Effect Size for Evaluating the Strength of Phylogenetic Signal, and Why Lambda is not Appropriate

<sup>3</sup>

<sup>4</sup>

<sup>5</sup> **Abstract**

<sup>6</sup> Macroevolutionary studies frequently characterize the phylogenetic signal in phenotypes, and wish to compare  
<sup>7</sup> the strength of that signal across traits. However, analytical tools for such comparisons have largely remained  
<sup>8</sup> underdeveloped. In this study, we evaluated the efficacy of one commonly used parameter (Pagel's  $\lambda$ ) to  
<sup>9</sup> estimate the strength of phylogenetic signal in phenotypic traits, and evaluate the degree to which  $\lambda$  correctly  
<sup>10</sup> identifies known levels of phylogenetic signal. We find that the precision of  $\lambda$  in estimating actual levels of  
<sup>11</sup> phylogenetic signal is often inaccurate, and that biological interpretations of the strength of phylogenetic  
<sup>12</sup> signal based on  $\lambda$  are therefore compromised. We then propose a standardized effect size based on *Kappa*  
<sup>13</sup> ( $Z_K$ ), which measures the strength of phylogenetic signal, and places it on a common scale for statistical  
<sup>14</sup> comparison. Tests based on  $Z_K$  provide a mechanism for formally comparing the strength of phylogenetic  
<sup>15</sup> signal across datasets, in much the same manner as effect sizes may be used to summarize patterns in  
<sup>16</sup> quantitative meta-analysis. Our approach extends the phylogenetic comparative toolkit to address hypotheses  
<sup>17</sup> that compare the strength of phylogenetic signal between various phenotypic traits, even when those traits  
<sup>18</sup> are found in different evolutionary lineages or of in have different units or scale.

19 **Introduction**

20 Investigating macroevolutionary patterns of trait variation requires a phylogenetic perspective, because the  
21 shared ancestry among species generates statistical non-independence violates an assumption of independence  
22 among trait values that is common for statistical tests (Felsenstein 1985; Harvey and Pagel 1991). Accounting  
23 for this evolutionary non-independence is the purview of *phylogenetic comparative methods* (PCMs) ~~±~~: a  
24 suite of analytical tools that condition trends in the data on the phylogeny through the course of statistical  
25 evaluations of phenotypic trends phylogenetic relatedness of observations (e.g., Grafen 1989; Garland and  
26 Ives 2000; Rohlf 2001; Butler and King 2004). The past several decades have witnessed a rapid expansion in  
27 the development of PCMs to address an ever-growing set of macroevolutionary hypotheses (Martins and  
28 Hansen 1997; O'Meara et al. 2006; Revell and Harmon 2008; Beaulieu et al. 2012; Adams 2014b,a; Adams  
29 and Collyer 2018). These methods are predicated on the notion that phylogenetic signal – the tendency  
30 for closely related species to display similar trait values – is present in cross-species datasets (Felsenstein  
31 1985; Pagel 1999; Blomberg et al. 2003). Indeed, under numerous evolutionary models, phylogenetic  
32 signal is to be expected, as stochastic character change along the hierarchical structure of the tree of life  
33 generates trait covariation among related taxa (see Felsenstein 1985; Blomberg et al. 2003; Revell et al. 2008).

34

35 Several analytical tools have been developed to quantify phylogenetic signal in phenotypic datasets, including  
36 measures of serial independence ( $C$ : Abouheif 1999), autocorrelation estimates ( $I$ : Gittleman and Kot 1990),  
37 statistical ratios of trait variation relative to what is expected given the phylogeny ( $Kappa$ : Blomberg et  
38 al. 2003; Adams 2014a), and scaling parameters used in maximum likelihood fitting of the data to the  
39 phylogeny ( $\lambda$ : Pagel 1999), among others (e.g., Klingenberg and Gidaszewski 2010). The statistical properties  
40 of these methods – namely type I error rates and power – have also been investigated to determine when  
41 phylogenetic signal can be detected and under what conditions (e.g., Munkemuller et al. 2012; Pavoine  
42 and Ricotta 2012; Diniz-Filho et al. 2012; Adams 2014a; Molina-Venegas and Rodriguez 2017; see also  
43 Revell et al. 2008; Revell 2010). One of the most widely used methods for characterizing phylogenetic signal  
44 in macroevolutionary studies is Pagel's  $\lambda$  (Pagel 1999). Here, maximum likelihood is used to fit the data  
45 to the phylogeny under a Brownian motion model of evolution. A parameter  $\lambda$  is included, which  
46 transforms parameter  $\lambda$  to the lengths of the internal branches of the phylogeny to improve the fit a  
47 fit of data to the phylogeny via maximum likelihood (Pagel 1999; Freckleton et al. 2002). Pagel's  $\lambda$  ranges  
48 from 0 → 1, with larger values signifying a greater dependence of observed trait variation on the phylogeny.  
49 Pagel's  $\lambda$  also has the appeal that it may be included in phylogenetic generalized least-squares regression

50 (PGLS) to account for the degree of phylogenetic signal in comparative analyses (see Freckleton et al. 2002).

51

52 ~~Evolutionary biologists commonly seek to~~ In addition to functioning as a parameter that is tuned for  
53 appropriate analysis,  $\lambda$  can function as a descriptive statistic to describe the relative strength of phylogenetic  
54 phylogenetic signal in phenotypic traits, to determine the extent to which shared evolutionary history has  
55 influenced trait covariation among taxa. This is often accomplished by interpreting empirical estimates. The  
56 appeal of  $\lambda$ ; with smaller values signifying ‘weak’ phylogenetic signal, while larger values are interpreted  
57 as ‘strong’ phylogenetic signal as a descriptive statistic for evolutionary biologists is a basis for interpreting  
58 “weak” versus “strong” phylogenetic signal; i.e., small versus large values of  $\lambda$ , respectively, in a comparative  
59 sense (e.g., De Meester et al. 2019; Pintanel et al. 2019; Su et al. 2019). Other approaches for interpreting  $\lambda$   
60 are more statistical. For instance, some have evaluated whether the observed  $\lambda$  differs from some expected  
61 value through the use of confidence intervals (Vandeloek et al. 2019) or by performing likelihood ratio tests  
62 that compare the observed model fit to that obtained when  $\lambda = 0$  or  $\lambda = 1$  (Freckleton et al. 2002; Cooper  
63 et al. 2010; Bose et al. 2019). Additionally, qualitative comparisons of  $\lambda$  estimates obtained from multiple  
64 phenotypic traits have been used to infer whether the strength of phylogenetic signal is greater in one trait  
65 as compared to another (e.g., Liu et al. 2019; Bai et al. 2019). Indeed, statements regarding the strength of  
66 phylogenetic signal based on  $\lambda$  are rather common in the evolutionary literature. For instance, of the 204  
67 papers published in 2019 that estimated and reported Pagel’s  $\lambda$  (found from a literature survey we conducted  
68 in Google.scholar), 40% interpreted the strength of phylogenetic signal for at least one phenotypic trait.  
69 Further, because nearly half of the 1,572  $\lambda$  values reported were near 0 or 1 (Figure 1) where the biological  
70 interpretation of  $\lambda$  is known, this percentage is even higher.

71

72

73 [insert Figure 1 here]

74

75 Various other approaches use  $\lambda$  as a parameter that can be varied for inferences akin to sensitivity analysis. For  
76 instance, some have performed likelihood ratio tests that compare observed model fits to those obtained when  
77  $\lambda = 0$  or  $\lambda = 1$  (Freckleton et al. 2002; Cooper et al. 2010; Bose et al. 2019) or evaluated whether observed  
78  $\lambda$  differs from an expected  $\lambda$ , based on confidence intervals generated for the expected value (Vandeloek et  
79 al. 2019). Qualitative comparisons of  $\lambda$  estimates have also been performed for multiple traits on the same  
80 phylogenetic tree to infer whether the strength of phylogenetic signal is greater in one trait as compared to

81 another (e.g., Liu et al. 2019; Bai et al. 2019).

82 It seems intuitive to interpret the strength of phylogenetic signal based on the value of  $\lambda$ , as  $\lambda$  is a parameter  
83 on a bounded scale ( $0 \rightarrow 1$ ) for which interpretation of its extremal points are understood. Specifically,  
84  $\lambda = 0$  represents no phylogenetic signal, while  $\lambda = 1$  is phylogenetic signal as expected under Brownian  
85 motion. However, equating values of  $\lambda$  directly to the strength of phylogenetic signal presumes two important  
86 statistical properties that have not been fully explored. First, it presumes that values of  $\lambda$  can be precisely  
87 estimated, as biological inferences regarding the strength of phylogenetic signal depend on high accuracy in  
88 its estimation. Therefore, understanding the precision in estimating  $\lambda$  is paramount. One study (Boettiger et  
89 al. 2012) found that estimates of Pagel's  $\lambda$  displayed less variation (i.e., greater precision) when data were  
90 simulated on a large phylogeny ( $N = 281$ ) as compared to a small one ( $N = 13$ ). From this observation it  
91 was concluded that insufficient data (i.e., the number of species) was the underlying cause of the increased  
92 variation across parameter estimates (Boettiger et al. 2012). Indeed, such a pattern is common with  
93 statistical estimators, as summary statistics and parameters are often more precise at greater sample sizes  
94 (Cohen 1988). However, this conclusion also implies that the precision of  $\lambda$  remains constant across its range  
95 ( $\lambda = 0 \rightarrow 1$ ); an assumption that to date, has not been verified. Thus, despite widespread use of Pagel's  
96 (1999)  $\lambda$  in macroevolutionary studies, at present, we lack a general understanding of the precision with  
97 which  $\lambda$  can estimate levels of phylogenetic signal in phenotypic datasets.

98

99 Second, while estimates of  $\lambda$  are within a bounded scale ( $0 \rightarrow 1$ ), this does not *de-facto* imply that the  
100 estimated values of this parameter correspond to the actual strength of the underlying input signal in  
101 the data. For this to be the case,  $\lambda$  must be a statistical effect size. Effect sizes are a measure of the  
102 magnitude of a statistical effect in data, represented on a common scale (Glass 1976; Cohen 1988). Effect  
103 sizes have widespread use in many areas of the quantitative sciences, as they represent measures that may  
104 be readily summarized across datasets as in meta-analyses (Glass 1976; Hedges and Olkin 1985; Arnqvist  
105 and Wooster 1995), or compared among datasets (e.g., Adams and Collyer 2016, 2019a). Unfortunatley,  
106 not all model parameters and ~~test descriptive~~ statistics are effect sizes, and thus many summary measures  
107 must first be converted to statistics with standardized units (i.e., conversion to an effect size) for meaningful  
108 comparison (see Rosenthal 1994). As a consequence, it follows that only if  $\lambda$  is a statistical effect size  
109 can comparisons of estimates across datasets be interpretable. For the case of  $\lambda$ , this has not yet been explored.

110

111 In this study, we evaluate the precision of Pagel's  $\lambda$  for estimating known levels of phylogenetic signal  
112 in phenotypic data. We use computer simulations with differing numbers of species, differently shaped

113 phylogenies, and differing input levels of phylogenetic signal, to explore the degree to which  $\lambda$  correctly  
114 identifies known levels of phylogenetic signal, and under what circumstances. We find that estimates of  
115  $\lambda$  vary widely for a given input value of phylogenetic signal, and that the precision in estimating  $\lambda$  is not  
116 constant across its range. Rather, there is decreased precision when input levels of phylogenetic signal are of  
117 intermediate strength. Additionally, the same estimated values of  $\lambda$  may be obtained from datasets containing  
118 vastly different input levels of phylogenetic signal. Thus,  $\lambda$  is not a reliable indicator of the strength of  
119 phylogenetic signal in phenotypic data. We then describe a standardized effect size for measuring the strength  
120 of phylogenetic signal in phenotypic datasets, and apply the concept to two common measures of phylogenetic  
121 signal:  $\lambda$  and *Kappa*. Through simulations we find that the precision of effect sizes based on  $\lambda$  ( $Z_\lambda$ ) are less  
122 reliable than those based on *Kappa* ( $Z_K$ ), implying that  $Z_K$  is a more robust effect size measure. We  
123 also propose a two-sample test statistic that may be used to compare the strength of phylogenetic signal  
124 among datasets, and provide an empirical example to demonstrate its use. We conclude that estimates of  
125 phylogenetic signal using Pagel's  $\lambda$  are often inaccurate, and thus interpreting strength of phylogenetic signal  
126 in phenotypic datasets based on this measure is compromised. By contrast, effect sizes obtained from *Kappa*  
127 hold promise for characterizing phylogenetic signal, and for comparing the strength of phylogenetic signal  
128 across datasets.

129 [note from Mike](#) Throughout the ms thus far *Kappa* is used. Why not  $\kappa$ , the actual greek symbol for  
130 [Kappa?](#)

## 131 Methods and Results

### 132 *The Precision of $\lambda$ is Variable*

133 We conducted a series of computer simulations to evaluate the precision of Pagel's  $\lambda$ . Our primary  
134 simulations were based on pure-birth phylogenies; however, we also evaluated patterns on both balanced  
135 and pectinate trees to determine whether tree shape affected our findings (see Supporting Information).  
136 First we generated 50 pure-birth phylogenies at each of six different tree sizes, ranging from 32 to 1024  
137 taxa ( $n = 2^5 - 2^{10}$ ). Next, we rescaled the simulated phylogenies by multiplying the internal branches by  
138  $\lambda_{in}$ , using 21 intervals of 0.05 units across its range ( $\lambda_{in} = 0.0 \rightarrow 1.0$ ), resulting in 1050 scaled phylogenies  
139 at each level of species richness ( $n$ ). Continuous traits were then simulated on each phylogeny under a  
140 Brownian motion model of evolution to obtain datasets with differing levels of phylogenetic signal, that  
141 ranged from no phylogenetic signal (when  $\lambda_{in} = 0$ ), to phylogenetic signal reflecting Brownian motion (when

<sup>142</sup>  $\lambda_{in} = 1$ ). For each dataset we then estimated phylogenetic signal ( $\lambda_{est}$ ), and calculated the variance of  $\lambda$   
<sup>143</sup> ( $\sigma_\lambda^2$ ) across datasets at each input level of phylogenetic signal and level of species richness as an estimate  
<sup>144</sup> of precision. We verified that the variance of traits simulated had no effect on phylogenetic signal estimation.

<sup>145</sup>

<sup>146</sup> We also evaluated the precision of  $\lambda$  when estimated in PGLS regression and ANOVA (i.e.,  $Y \sim X$ ). Here,  
<sup>147</sup> an independent variable  $X$  was simulated on each rescaled phylogeny under a Brownian motion model of  
<sup>148</sup> evolution (for PGLS regression). For phylogenetic ANOVA, random groups ( $X$ ) were obtained by simulating  
<sup>149</sup> a discrete (binary, 0 or 1) character on each phylogeny. Next, the dependent variable was simulated in  
<sup>150</sup> such a manner as to contain a known relationship with  $X$  plus random error containing phylogenetic  
<sup>151</sup> signal. This was accomplished as:  $Y = \beta X + \epsilon$ . Here, the The association between  $Y$  and  $X$  was modeled  
<sup>152</sup> using a range of values:  $\beta = (0.0, 0.25, 0.5, 0.75, 1.0)$ , and the residual error ( $\epsilon$ ) was modeled to contain  
<sup>153</sup> phylogenetic signal simulated under a Brownian motion model of evolution on each rescaled phylogeny:  
<sup>154</sup>  $\epsilon = \mathcal{N}(\mu = 0, \sigma = \sigma^2 \mathbf{C})$ : (see Revell 2010 for a similar simulation design). The fit of the phylogenetic  
<sup>155</sup> regression was estimated using maximum likelihood, and parameter estimates ( $\beta_{est}$  and  $\lambda_{est}$ ) were obtained.  
<sup>156</sup> We then calculated precision estimates ( $\sigma_\lambda^2$ ) at each input level of phylogenetic signal and level of species  
<sup>157</sup> richness. We verified that the amount of residual variance simulated had no effect on  $\sigma^2$  but did influence  
<sup>158</sup> the precision of coefficients estimated from the linear model (precision increased with smaller  $\epsilon$ , as expected).

<sup>159</sup>

<sup>160</sup> All analyses were performed in R v3.6.0 (R Core Team 2019) using the packages **geiger** (Harmon et al.  
<sup>161</sup> 2008), **caper** (Orme et al. 2013), **phytools** (Revell 2012), and **geomorph** 3.2.1 (Adams and Otárola-Castillo  
<sup>162</sup> 2013; Adams et al. 2020). R-scripts are found in the Supporting Information.

<sup>163</sup>

<sup>164</sup> *Results.* We found that the precision of  $\lambda_{est}$  varied widely across simulation conditions. Predictably, precision  
<sup>165</sup> improved as the number of species increased (Figure 2). This confirmed earlier findings of Boettiger et al.  
<sup>166</sup> (2012), and adhered to parametric statistical theory. However, in many cases the set of  $\lambda_{est}$  spanned nearly  
<sup>167</sup> the entire range of possible values (e.g.,  $n = 32$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.0 \rightarrow 0.985$ ), revealing that estimates  
<sup>168</sup> of  $\lambda$  were not a reliable indicator of input phylogenetic signal. Importantly, the precision of  $\lambda_{est}$  was not  
<sup>169</sup> uniform across all levels of phylogenetic signal. The worst precision was observed at intermediate levels  
<sup>170</sup> of phylogenetic signal ( $\lambda_{in} \approx 0.5$ ), while precision improved as input levels approached the extremes of  
<sup>171</sup>  $\lambda$ 's range (i.e.,  $\lambda_{in} \rightarrow 0$  &  $\lambda_{in} \rightarrow 1$ ). Thus, estimates of  $\lambda$  were least reflective of the true input signal at  
<sup>172</sup> intermediate values. Additionally, even at large levels of species richness, we found that the range of  $\lambda_{est}$  still  
<sup>173</sup> encompassed a substantial portion of possible values (e.g.,  $n = 512$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.32 \rightarrow 0.68$ ). Likewise,

174 the same  $\lambda_{est}$  could be obtained from datasets containing vastly different input levels of phylogenetic  
175 signal (e.g.,  $n = 512$ ;  $\lambda_{est} = 0.5$ ;  $\lambda_{in} = 0.25 \rightarrow 0.65$ ). These findings were particularly unsettling when  
176 considered in light of our literature survey. Over one quarter of the  $\lambda$  estimates published in empirical  
177 studies (421 of 1,572) were between  $\lambda = 0.25$  and  $\lambda = 0.75$  (Figure 1). This range reflected the region  
178 that our simulations identified as being the least reliable in terms of accurately characterizing levels of  
179 phylogenetic signal, yet 30% of these mid-range empirical estimates were explicitly interpreted in terms of  
180 the strength of phylogenetic signal that they represented (i.e., weak, intermediate, strong phylogenetic signal).

181

182 Finally, when  $\lambda$  was co-estimated with regression parameters in PGLS regression and ANOVA, the results of our  
183 simulations were quite similar. [Here, regression](#) [Regression](#) parameters ( $\beta$ ) were accurately estimated, confirm-  
184 ing earlier findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal  
185 ( $\lambda$ ) were less precise (Figure 3; see also Supporting Information), and the spread of  $\lambda_{est}$  was similar to that  
186 observed when  $\lambda$  was estimated for only the dependent variable, as in Figure 2. Taken together, these findings  
187 reveal that  $\lambda_{est}$  does not precisely characterize observed levels of phylogenetic signal in phenotypic datasets,  
188 and that biological interpretations of the strength of phylogenetic signal based on  $\lambda$  may be highly inaccurate.

189

190 [insert Figure 2 here]

191

192 [insert Figure 3 here]

193

#### 194 **A Standardized Effect Size for Phylogenetic Signal**

195 The results above demonstrate that  $\lambda$  is not a reliable estimate of the phylogenetic signal in phenotypic data.  
196 As such, biological interpretations of the strength of phylogenetic signal, and comparisons of the magnitude  
197 of such effects across datasets are severely compromised when based on this parameter. As an alternative, we  
198 propose that summary estimates of phylogenetic signal be converted to effect sizes for interpretation and  
199 comparison of the relative strength of phylogenetic signal in phenotypic datasets. Statistically, a standardized  
200 effect size may be found as:

$$Z_\theta = \frac{\theta_{obs} - E(\theta)}{\sigma_\theta} \quad (1)$$

201 where  $\theta_{obs}$  is the observed test statistic,  $E(\theta)$  is its expected value under the null hypothesis, and  $\sigma_\theta$  is its  
 202 standard error (Glass 1976; Cohen 1988; Rosenthal 1994).  $Z_\theta$  expresses the magnitude of the effect in  $\theta_{obs}$  by  
 203 transforming the original test statistic to its standard normal deviate (Glass 1976; Kelley and Preacher 2012).  
 204 Typically,  $\theta_{obs}$  and  $\sigma_\theta$  are estimated from the data, while  $E(\theta)$  is obtained from the distribution of  $\theta$  derived  
 205 from parametric theory. However, recent advances in resampling theory (Collyer et al. 2015; Adams and  
 206 Collyer 2016, 2019a; Collyer and Adams 2018) have shown that  $E(\theta)$  and  $\sigma_\theta$  may also be obtained from an  
 207 empirical sampling distribution of  $\theta$  obtained from permutation procedures.

208

209 Adams and Collyer (2019b) suggested that the strength of phylogenetic signal could be represented as an  
 210 effect size based on the *Kappa* statistic and its empirical sampling distribution from permutation. Here we  
 211 formalize that suggestion, resulting in an effect size of:

$$Z_K = \frac{\log(K_{obs}) - \hat{\mu}_{\log(K)}}{\hat{\sigma}_{\log(K)}} \quad (2)$$

212 where  $K_{obs}$  is the observed phylogenetic signal, and  $\hat{\mu}_K$  and  $\hat{\sigma}_K$  are the mean and standard deviation of the  
 213 empirical sampling distribution of  $\log(Kappa)$  obtained via permutation. Note that the logarithm was used be-  
 214 cause *Kappa* takes only positive values ( $0 \rightarrow \infty$ ) and its sampling distribution is log-normally distributed (for a  
 215 similar transformation when calculating multivariate effect sizes see: Appendix 1 of Collyer and Adams 2018).

216

217 An effect size based on  $\lambda$  could be envisioned, which is found as:

$$Z_\lambda = \frac{\lambda_{obs} - 0}{\hat{\sigma}_\lambda}. \quad (3)$$

218 In this case,  $\lambda_{obs}$  and  $\hat{\sigma}_\lambda$  are empirically derived using maximum likelihood, as permutation approaches have  
 219 not been developed for evaluating  $\lambda$ . Note also that under the null hypothesis, no phylogenetic signal is

220 expected (Freckleton et al. 2002), and thus  $E(\lambda) = 0$  under this condition.

221

222 To evaluate the utility of  $Z_K$  and  $Z_\lambda$  we calculated both effect sizes for the simulated datasets generated  
223 above, and summarized the precision of each using its variance ( $\sigma_{Z_K}^2$  and  $\sigma_{Z_\lambda}^2$ , Figure 4: additional results in  
224 the Supporting Information). Here two things are evident. First, estimates of  $Z_K$  linearly track the input  
225 phylogenetic signal whereas estimates of  $Z_\lambda$  do not (Figure 4A, B). Thus, actual changes in the strength  
226 of phylogenetic signal are reflected more evenly in the corresponding values of the effect size  $Z_K$ . Second,  
227 the precision of  $Z_K$  is considerably more stable as compared with  $Z_\lambda$ . This may be seen by calculating  
228 the coefficients of variation for the set of precision estimates (i.e.,  $\sigma_{Z_K}^2$  and  $\sigma_{Z_\lambda}^2$ ) across input levels of  
229 phylogenetic signal. ~~Here coefficients~~ Coefficients of variation in the precision of  $Z_K$  were up to an order of  
230 magnitude smaller for than for  $Z_\lambda$  (Figure 4C), implying that estimates of the strength of phylogenetic signal  
231 were more reliable and robust when using  $Z_K$ .

232

233 [insert Figure 4 here]

234 ***Statistical Comparisons of Phylogenetic Signal***

235 Once the magnitude of phylogenetic signal is characterized using  $Z_K$ , one may wish to compare such measures  
236 across datasets, to determine whether the strength of phylogenetic signal is significantly greater in one  
237 phenotypic trait than another. As with other effect sizes derived from permutation distributions (e.g., Adams  
238 and Collyer 2016, 2019a), a two-sample test statistic may be calculated as:

$$\hat{Z}_{12} = \frac{|(K_1 - \hat{\mu}_{K_1}) - (K_2 - \hat{\mu}_{K_2})|}{\sqrt{\hat{\sigma}_{K_1}^2 + \hat{\sigma}_{K_2}^2}} = \frac{|Z_{K_1} - Z_{K_2}|}{\sqrt{2}} \quad (4)$$

239 where  $K_1$ ,  $K_2$ ,  $\hat{\mu}_{K_1}$ ,  $\hat{\mu}_{K_2}$ ,  $\hat{\sigma}_{K_1}$ , and  $\hat{\sigma}_{K_2}$  are as defined above for equation 2. The right side of the equation  
240 illustrates that if  $Z_K$  has already been calculated for two sampling distributions as in equation 2, the  
241 sampling distributions have unit variance for each of the  $Z_K$  statistics. Estimates of significance of  $\hat{Z}_{12}$  may  
242 be obtained from a standard normal distribution. Typically,  $\hat{Z}_{12}$  is considered a two-tailed test, however  
243 directional (one-tailed) tests may be specified should the empirical situation require it (see Adams and  
244 Collyer 2016, 2019a).

245

246 **Empirical Example**

247 To demonstrate the utility of  $\hat{Z}_{12}$  we quantified and compared the strength of phylogenetic signal of two  
248 phenotypic traits across species of plethodontid salamander. The data were part of a series of studies  
249 examining macroevolutionary trends in phenotypic diversification in this group (Baken and Adams 2019;  
250 Baken et al. 2020). Our dataset contained surface area to volume ratios (SA:V) and relative body width  
251 ( $\frac{BW}{SVL}$ ) for 305 species (Figure 5A). For SA:V, 11 linear body measurements from 2,781 individuals were  
252 taken, from which estimates of the surface area and volume of the head, body, and tail were calculated and  
253 subsequently combined to arrive at the SA:V for each individual (for mathematical details see Baken et al.  
254 2020). Species means were then obtained. Likewise, body size (SVL) and body width (BW) measurements  
255 were taken from 3,371 individuals, and species means of relative body width ( $\frac{BW}{SVL}$ ) were calculated (data  
256 from Baken and Adams 2019). A time-dated molecular phylogeny for the group (Bonett and Blair 2017) was  
257 then pruned to match the species in the dataset, resulting in a phylogeny and corresponding phenotypic  
258 dataset containing 305 species. The phylogenetic signal in each trait was then characterized using *Kappa*,  
259 which was converted to its effect size ( $Z_K$ ) using geomorph 3.2.1 (Adams and Otárola-Castillo 2013; Adams  
260 et al. 2020). Finally, the strength of phylogenetic signal was compared across traits using  $\hat{Z}_{12}$  as described  
261 above (to be incorporated in geomorph upon manuscript acceptance).

262

263 *Results.* Both SA:V and relative body width displayed significant phylogenetic signal ( $Kappa_{SA:V} = 0.7608$ ;  
264  $P = 0.001$ ;  $Kappa_{BW/SVL} = 0.2515$ ;  $P = 0.001$ ). For both phenotypic traits,  $K_{obs}$  differed markedly from  
265 their corresponding permutation distributions, which were found to overlap almost perfectly (Figure 5B).  
266 However, while both traits displayed significant phylogenetic signal, there was nearly a four-fold difference  
267 in the magnitude of their effect sizes, with SA:V displaying the greater phylogenetic signal (Figure 5C).  
268 Using the two-sample test statistic above, this difference was found to be highly significant ( $\hat{Z}_{12} = 4.13$ ;  
269  $P = 0.000036$ ). Thus it may be concluded that SA:V displays significantly stronger phylogenetic signal  
270 than does relative body width, and that shared evolutionary history has strongly influenced trait covariation  
271 among taxa for SA:V. Biologically, this observation corresponds with the fact that tropical species – which  
272 form a monophyletic group within plethodontids – display greater variation in SA:V which covaries with  
273 disparity in their climatic niches (Baken et al. 2020). We hypothesize that because of this macroevolutionary  
274 association, strong phylogenetic signal in SA:V is observed.

275 **Discussion**

276 It is common in comparative evolutionary studies to characterize the phylogenetic signal in phenotypic traits  
277 to determine the extent to which shared evolutionary history has generated trait covariation among taxa.  
278 However, while numerous analytical approaches may be used to quantify phylogenetic signal (e.g., Abouheif  
279 1999; Gittleman and Kot 1990; Pagel 1999; Blomberg et al. 2003; Adams 2014a), methods that explicitly  
280 measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained  
281 underdeveloped. In this study, we evaluated the precision of one common measure, Pagel's  $\lambda$ , and explored its  
282 efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations,  
283 we found that the precision of  $\lambda$  increased with increasing sample sizes; a pattern noted previously (Boettiger  
284 et al. 2012), and one that conformed with parametric statistical theory (Cohen 1988). However, we also found  
285 that vastly different  $\lambda$  estimates could be obtained from data containing the same level of phylogenetic signal,  
286 and that similar  $\lambda$  estimates may be obtained from data containing differing levels of phylogenetic signal.  
287 Further, the precision of  $\lambda$  varied with the strength of phylogenetic signal, where lower precision was observed  
288 when in data whose phylogenetic signal was of intermediate strength. From these findings we conclude that  $\lambda$  is  
289 not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and that biological  
290 interpretations of the strength of signal based on this parameter may inaccurately characterize such effects.

291

292 As an alternative, we described a standardized effect size ( $Z$ ) for assessing the strength of phylogenetic signal.  
293  $Z$  expresses the magnitude of phylogenetic signal as a standard normal deviate, which is easily interpretable  
294 as the strength of phylogenetic signal relative to the mean. We applied this concept to both  $\lambda$  and  $Kappa$ ,  
295 and found that  $Z_K$  was a better estimate of the strength of phylogenetic signal in phenotypic data. First,  $Z_K$   
296 was more precise than  $Z_\lambda$ , and variation in this precision was more consistent across the range of input levels  
297 of phylogenetic signal. Additionally, values of  $Z_K$  more accurately tracked known levels of phylogenetic signal,  
298 with changes in the actual strength of phylogenetic signal reflected in a more linear fashion by concomitant  
299 changes in the values of  $Z_K$ . Thus,  $Z_K$  holds promise as a measure of the relative strength of phylogenetic  
300 signal that reflects the magnitude of this effect in phenotypic data. We therefore recommend that future  
301 studies interested in the strength of phylogenetic signal incorporate  $Z_K$  as a statistical measure of this effect.

302

303 Based on the effect size  $Z_K$ , we then proposed a two-sample test, which provides a statistical means of  
304 determining whether the strength of phylogenetic signal is greater in one phenotypic trait as compared  
305 to another, via a hypothesis test. Prior studies have summarized patterns of variation in phylogenetic

306 signal across datasets using summary test values, such as *Kappa* (e.g., Blomberg et al. 2003). However,  
307 *Kappa* does not scale linearly with input levels of phylogenetic signal, and its variance increases (i.e.,  
308 precision decreases) with increasing strength of phylogenetic signal (Munkemuller et al. 2012; Diniz-Filho  
309 et al. 2012: see also Supporting Information). Thus, *Kappa* should not be considered ~~a standardized an~~  
310 effect size that measures the strength of phylogenetic signal on a common scale. By contrast, ~~converting~~  
311 standardizing *Kappa* ~~to  $Z_K$~~  ( $Z_K$ , via equation 2) alleviates these concerns, and facilitates formal statistical  
312 comparisons of the strength of signal across datasets. Thus when viewed from this perspective, the approach  
313 developed here aligns well with other statistical approaches such as meta-analysis (sensu Hedges and Olkin  
314 1985; Glass 1976; Arnqvist and Wooster 1995), where summary statistics across datasets are converted to  
315 standardized effect sizes for subsequent ~~'higher-order'~~ "higher order" statistical summaries or comparisons.  
316 As such our approach enables evolutionary biologists to quantitatively examine the relative strength  
317 of phylogenetic signal across a wide range of phenotypic traits, and thus opens the door for future dis-  
318 coveries that inform on how phenotypic diversity accumulates in macroevolutionary time across the tree of life.

319

320 One important advantage of the approach advocated here is that the resulting effect sizes ( $Z_K$ ) are  
321 dimensionless, as the units of measurement cancel out during the calculation of  $Z$  (Sokal and Rohlf 2012).  
322 Thus,  $Z_K$  represents the strength of phylogenetic signal on a common and comparable scale – measured  
323 in standard deviation ~~units~~ – regardless of the initial units and original scale of the phenotypic variables  
324 under investigation. This means that the strength of phylogenetic signal may be compared across datasets  
325 for continuous phenotypic traits measured in different units and scale, because those units have been  
326 standardized through their conversion to  $Z_K$ . For example, our approach could be utilized to determine  
327 whether the strength of phylogenetic signal (say, in response to ecological differentiation) is stronger in  
328 morphological traits (linear traits:  $mm$ ), physiological traits (metabolic rate:  $\frac{O^2}{min}$ ), or behavioral traits  
329 (aggression:  $\frac{\#displays}{second}$ ). In fact, our empirical example provided such a comparison, as SA:V is represented in  
330  $mm^{-1}$  while relative body size is a unitless ratio ( $\frac{BW}{SVL}$ ). Additionally, our method is capable of comparing  
331 the strength of phylogenetic signal in traits of different dimensionality, as estimates of phylogenetic signal  
332 using *Kappa* have been generalized for multivariate data ( $K_{mult}$ : see Adams 2014a). Furthermore, tests  
333 based on  $\hat{Z}_{12}$  may be utilized for comparing the strength of phylogenetic signal among datasets containing a  
334 different number of species, and even for phenotypes obtained from species in different lineages, because  
335 their phylogenetic non-independence and observed variation are taken into account in the generation of the  
336 empirical sampling distribution via permutation.

337

338 **Finally... Need Closing paragraph.** Phylogenetic signal can be thought of as both an attribute to be  
339 measured in the data and a parameter that can be tuned to account for the phylogenetic non-independence  
340 among observations, for analysis of the data. As such,  $\lambda$  is appealing, as a statistic that potentially fulfills  
341 both roles. However, the inability to estimate phylogenetic signal with  $\lambda$  for data simulated with known  
342 phylogenetic signal is troublesome, and we recommend evolutionary biologists refrain from viewing it as a  
343 useful statistic to describe the amount of phylogenetic signal in the data. Interestingly,  $Kappa$  is a better  
344 statistic for measuring the amount of phylogenetic signal in data simulated with respect to known levels of  
345  $\lambda$ . Although  $\lambda$  might be viewed as an important parameter for modifying the conditional estimation of  
346 linear model coefficients with respect to phylogeny, it is neither a statistic that has meaningful comparative  
347 value as a measure of phylogenetic signal nor a statistic that lends itself well to reliable calculation of a test  
348 statistic. By contrast,  $Kappa$  has been shown here to be a reliable statistic, but only when standardized by  
349 the mean and standard deviation of its empirical sampling distribution. Because one has control over the  
350 number of permutations used in analysis, one can be assured with many permutations that the empirical  
351 sampling distribution is representative of true probability distributions (Adams and Collyer 2018). With  
352 low coefficients of variation for  $Z_K$  (Figure 4), it is difficult to imagine that a hypothesis test can improve  
353 equation 4 for efficiently comparing phylogenetic signal for different traits, different trees, or a combination  
354 of both.

355    **References**

- 356    Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.  
357         Evolutionary Ecology Research 1:895–909.
- 358    Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other  
359         high-dimensional multivariate data. Systematic Biology 63:685–697.
- 360    Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other  
361         high-dimensional multivariate data. Evolution 68:2675–2688.
- 362    Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative  
363         modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- 364    Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration  
365         across morphometric datasets. Evolution 70:2623–2631.
- 366    Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,  
367         and a refined permutation procedure. Evolution 72:1204–1215.
- 368    Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate  
369         phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- 370    Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric  
371         analyses. R package version 3.2.1.
- 372    Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of  
373         geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- 374    Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.  
375         Trends in Ecology and Evolution 10:236–240.
- 376    Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with  
377         phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil  
378         434:305–326.
- 379    Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution  
380         9:7005–7016.

- 381 Baken, E. K., L. E. Mellenthin, and D. C. Adams. 2020. Macroevolution of desiccation-related morphology  
382 in plethodontid salamanders as inferred from a novel surface area to volume ratio estimation approach.  
383 *Evolution* 74:476–486.
- 384 Beaulieu, J. M., D. C. Jhwueng, C. Boettiger, and B. C. O'Meara. 2012. Modeling stabilizing selection:  
385 Expanding the ornstein-uhlenbeck model of adaptive evolution. *Evolution* 66:2369–2383.
- 386 Blomberg, S. P., T. Garland, and A. R. Ives. 2003. Testing for phylogenetic signal in comparative data:  
387 Behavioral traits are more labile. *Evolution* 57:717–745.
- 388 Boettiger, C., G. Coop, and P. Ralph. 2012. Is your phylogeny informative? Measuring the power of  
389 comparative methods. *Evolution* 67:2240–2251.
- 390 Bonett, R. M., and A. L. Blair. 2017. Evidence for complex life cycle constraints on salamander body form  
391 diversification. *Proceedings of the National Academy of Sciences, U.S.A.* 114:9936–9941.
- 392 Bose, R., B. R. Ramesh, R. Pélassier, and F. Munoz. 2019. Phylogenetic diversity in the western ghats  
393 biodiversity hotspot reflects environmental filtering and past niche diversification of trees. *Journal of  
394 Biogeography* 46:145–157.
- 395 Butler, M. A., and A. A. King. 2004. Phylogenetic comparative analysis: A modeling approach for adaptive  
396 evolution. *American Naturalist* 164:683–695.
- 397 Cohen, J. 1988. Statistical power analysis for the behavioral sciences. Routledge.
- 398 Collyer, M. L., and D. C. Adams. 2018. RRPP: An r package for fitting linear models to high-dimensional  
399 data using residual randomization. *Methods in Ecology and Evolution* 9:1772–1779.
- 400 Collyer, M. L., D. J. Sekora, and D. C. Adams. 2015. A method for analysis of phenotypic change for  
401 phenotypes described by high-dimensional data. *Heredity* 115:357–365.
- 402 Cooper, N., W. Jetz, and R. P. Freckleton. 2010. Phylogenetic comparative approaches for studying niche  
403 conservatism. *Journal of Evolutionary Biology* 23:2529–2539.
- 404 De Meester, G., K. Huyghe, and R. Van Damme. 2019. Brain size, ecology and sociality: A reptilian  
405 perspective. *Biological Journal of the Linnean Society* 126:381–391.
- 406 Diniz-Filho, J. A. F., T. Santos, T. F. Rangel, and L. M. Bini. 2012. A comparison of metrics for estimating  
407 phylogenetic signal under alternative evolutionary models. *Genetics and Molecular Biology* 35:673–679.

- 408 Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* 125:1–15.
- 409 Freckleton, R. P., P. H. Harvey, and M. Pagel. 2002. Phylogenetic analysis and comparative data: A test and  
410 review of evidence. *American Naturalist* 160:712–726.
- 411 Garland, T. J., and A. R. Ives. 2000. Using the past to predict the present: Confidence intervals for regression  
412 equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
- 413 Gittleman, J. L., and M. Kot. 1990. Adaptation: Statistics and a null model for estimating phylogenetic  
414 effects. *Systematic Zoology* 39:227–241.
- 415 Glass, G. V. 1976. Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
- 416 Grafen, A. 1989. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London B,*  
417 *Biological Sciences* 326:119–157.
- 418 Harmon, L. J., J. T. Weir, C. D. Brock, R. E. Glor, and W. Challenger. 2008. GEIGER: Investigating  
419 evolutionary radiations. *Bioinformatics* 24:129–131.
- 420 Harvey, P. H., and M. D. Pagel. 1991. The comparative method in evolutionary biology. Oxford University  
421 Press, Oxford.
- 422 Hedges, L. V., and I. Olkin. 1985. Statistical methods for meta-analysis. Elsevier.
- 423 Kelley, K., and K. J. Preacher. 2012. On effect size. *Psychological Methods* 17:137–152.
- 424 Klingenberg, C. P., and N. A. Gidaszewski. 2010. Testing and quantifying phylogenetic signals and homoplasy  
425 in morphometric data. *Systematic biology* 59:245–261.
- 426 Liu, H., C. P. Osborne, D. Yin, R. P. Freckleton, G. Jiang, and M. Liu. 2019. Phylogeny and ecological  
427 processes influence grass coexistence at different spatial scales within the steppe biome. *Oecologia*  
428 191:25–38.
- 429 Martins, E. P., and T. F. Hansen. 1997. Phylogenies and the comparative method: A general approach  
430 to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist*  
431 149:646–667.
- 432 Molina-Venegas, R., and M. A. Rodriguez. 2017. Revisiting phylogenetic signal; strong or negligible impacts  
433 of polytomies and branch length information? *BMC evolutionary biology* 17:53.
- 434 Munkemuller, T., S. Lavergne, B. Bzeznik, S. Dray, T. Jombart, K. Schiffers, and W. Thuiller. 2012. How to

- 435 measure and test phylogenetic signal. *Methods in Ecology and Evolution* 3:743–756.
- 436 O'Meara, B. C., C. Ane, M. J. Sanderson, and P. C. Wainwright. 2006. Testing for different rates of  
437 continuous trait evolution using likelihood. *Evolution* 60:922–933.
- 438 Orme, D., R. Freckleton, G. Thomas, T. Petzoldt, S. Fritz, and N. Isaac. 2013. CAPER: Comparative  
439 analyses of phylogenetics and evolution in r. *Methods in Ecology and Evolution* 3:145–151.
- 440 Pagel, M. D. 1999. Inferring the historical patterns of biological evolution. *Nature* 401:877–884.
- 441 Pavoine, S., and C. Ricotta. 2012. Testing for phylogenetic signal in biological traits: The ubiquity of  
442 cross-product statistics. *Evolution: International Journal of Organic Evolution* 67:828–840.
- 443 Pintanel, P., M. Tejedo, S. R. Ron, G. A. Llorente, and A. Merino-Viteri. 2019. Elevational and microclimatic  
444 drivers of thermal tolerance in andean pristimantis frogs. *Journal of Biogeography* 46:1664–1675.
- 445 R Core Team. 2019. R: A language and environment for statistical computing. R Foundation for Statistical  
446 Computing, Vienna, Austria.
- 447 Revell, L. J. 2010. Phylogenetic signal and linear regression on species data. *Methods in Ecology and  
448 Evolution* 1:319–329.
- 449 Revell, L. J. 2012. Phytools: An r package for phylogenetic comparative biology (and other things). *Methods  
450 in Ecology and Evolution* 3:217–223.
- 451 Revell, L. J., and L. J. Harmon. 2008. Testing quantitative genetic hypotheses about the evolutionary rate  
452 matrix for continuous characters. *Evolutionary Ecology Research* 10:311–331.
- 453 Revell, L. J., L. J. Harmon, and D. C. Collar. 2008. Phylogenetic signal, evolutionary process, and rate.  
454 *Systematic Biology* 57:591–601.
- 455 Rohlf, F. J. 2001. Comparative methods for the analysis of continuous variables: Geometric interpretations.  
456 *Evolution* 55:2143–2160.
- 457 Rosenthal, R. 1994. The handbook of research synthesis. Pp. 231–244 in L. V. Cooper H Hedges, ed. Russell  
458 Sage Foundation.
- 459 Sokal, R. R., and F. J. Rohlf. 2012. Biometry. 4th ed. W.H. Freeman & Co., San Francisco.
- 460 Su, G., S. Villéger, and S. Brosse. 2019. Morphological diversity of freshwater fishes differs between realms,  
461 but morphologically extreme species are widespread. *Global ecology and biogeography* 28:211–221.

- <sup>462</sup> Vandelook, F., S. Janssens, P. Gijbels, E. Fischer, W. Van den Ende, O. Honnay, and S. Abrahamczyk. 2019.  
<sup>463</sup> Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124:269–279.

464

## Figure Legends

465     **Figure 1.** Frequency distribution of  $\lambda$  estimates published in 2019. The majority of these values were  
466     close to 0 or 1, and from phylogenies with fewer than 200 taxa.

467

468     **Figure 2.** Precision of Pagel's  $\lambda$  across known levels of input phylogenetic signal ( $\lambda_{in}$ ) on phylogenies  
469     of various sizes. As phylogenies increase in size, variation in  $\lambda_{in}$  decreases; however the precision is  
470     not constant across the range of input levels ( $\lambda_{in} : 0 \rightarrow 1$ ), and is highest at intermediate levels of  
471     phylogenetic signal.

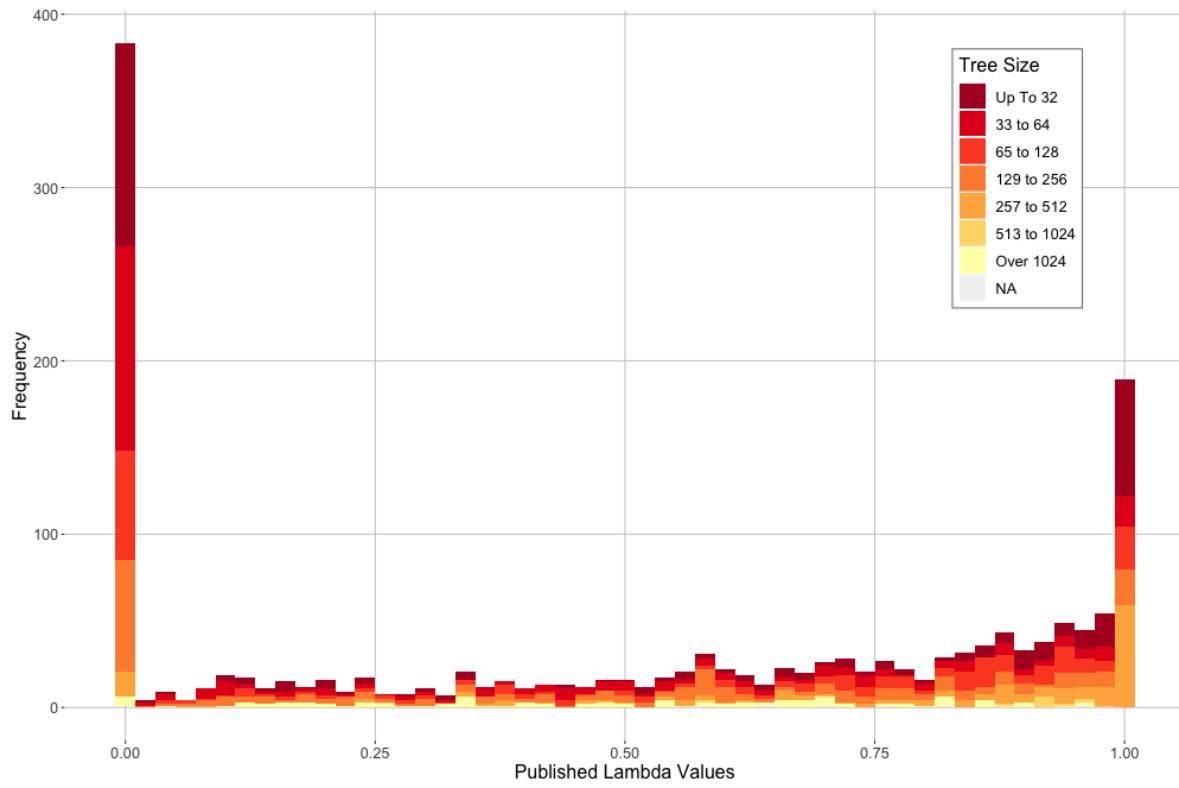
472

473     **Figure 3.** Precision of Pagel's  $\lambda$  when incorporated in phylogenetic regression ( $Y \sim X$ ), across known  
474     levels of input phylogenetic signal ( $\lambda_{in}$ ) on phylogenies of various sizes. As phylogenies increase in  
475     size, variation in  $\lambda_{in}$  decreases; however the precision is not constant across the range of input levels  
476     ( $\lambda_{in} : 0 \rightarrow 1$ ), and is highest at intermediate levels of phylogenetic signal.

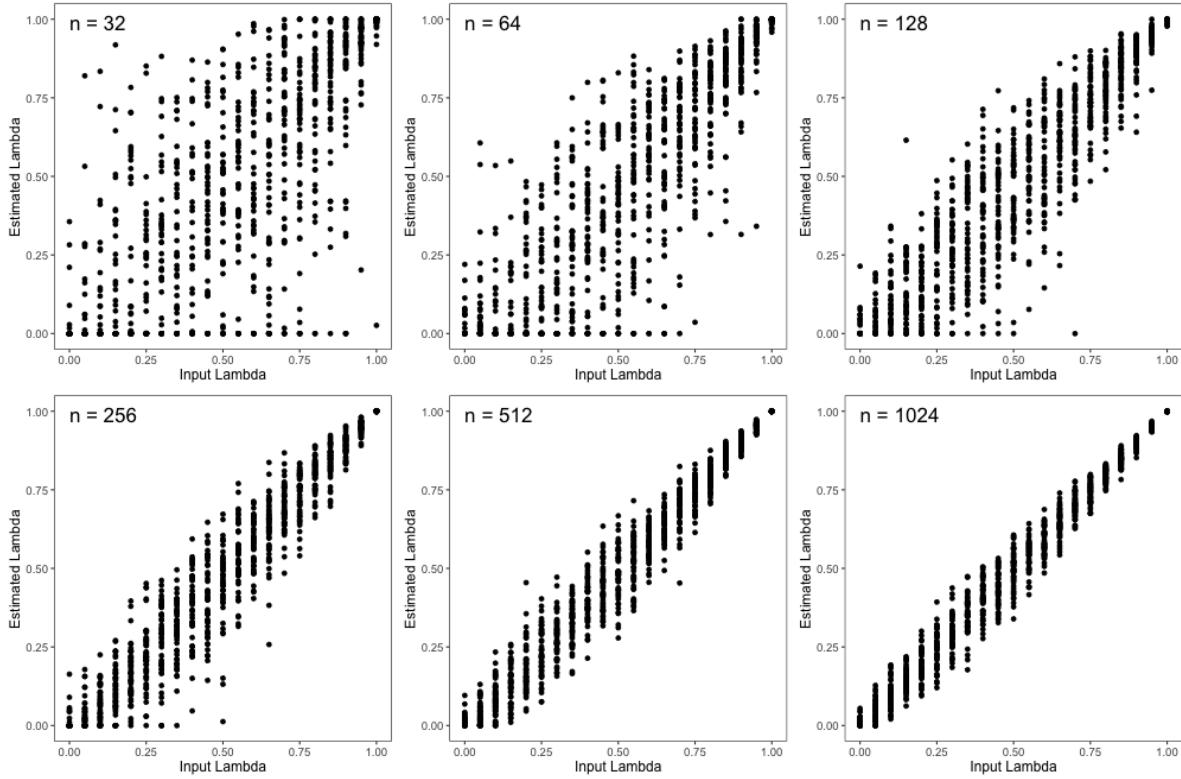
477

478     **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.  
479     (A) Estimates  $Z_\lambda$  for data simulated on phylogenies with 32 taxa ( $n = 32$ ), (B) Estimates of  $Z_K$  for data  
480     simulated on phylogenies with 32 taxa ( $n = 32$ ), (C) Coefficients of variation of precision estimates of  $Z_\lambda$   
481     and  $Z_K$  across input levels of phylogenetic signal, estimated on phylogenies containing differing numbers  
482     of species.

483     **Figure 5.** (A) Linear measures for relative body size, and regions of the body used to estimate surface  
484     area to volume (SA:V) ratios. (B) Permutation distributions of phylogenetic signal for SA:V and  $\frac{BW}{SVL}$ ,  
485     with observed values shown as vertical bars. (C) Effect sizes ( $Z_K$ ) for SA:V and  $\frac{BW}{SVL}$ , with their 95%  
486     confidence intervals (CI not standardized by  $\sqrt{(n)}$ ).

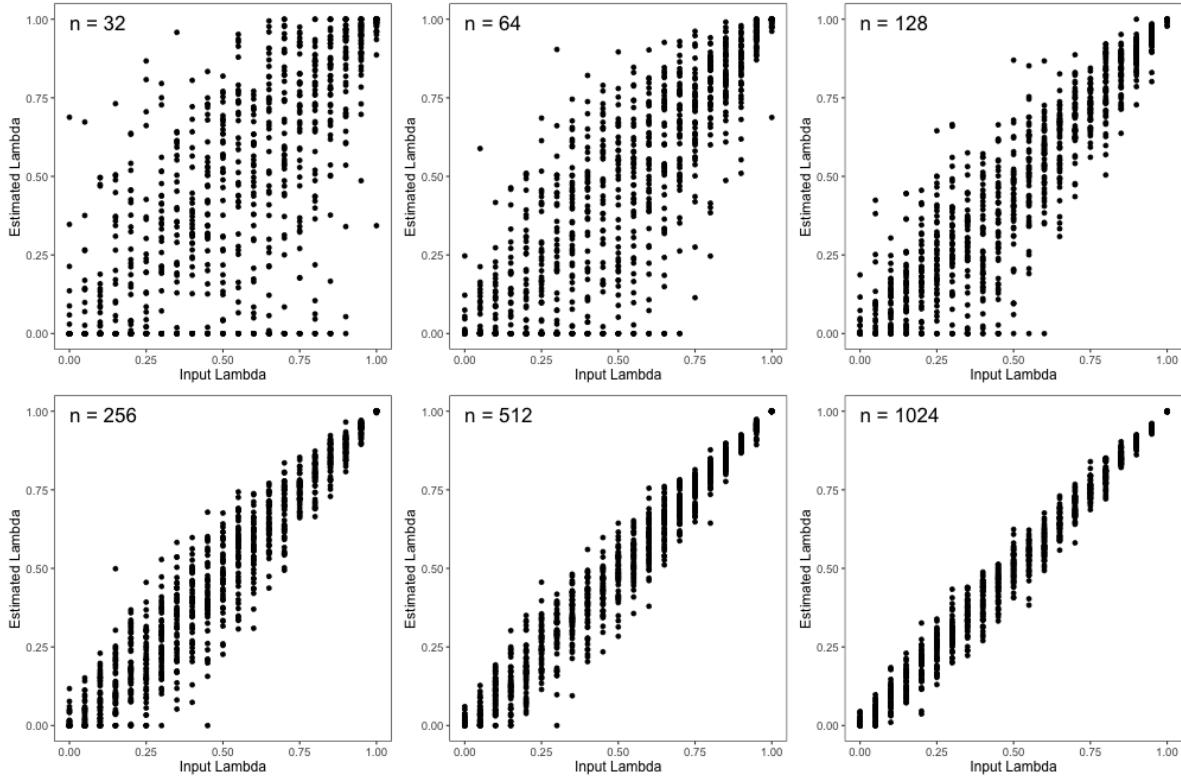


488 **Figure 1.** Frequency distribution of  $\lambda$  estimates published in 2019. The majority of these values were close  
489 to 0 or 1, and from phylogenies with fewer than 200 taxa.



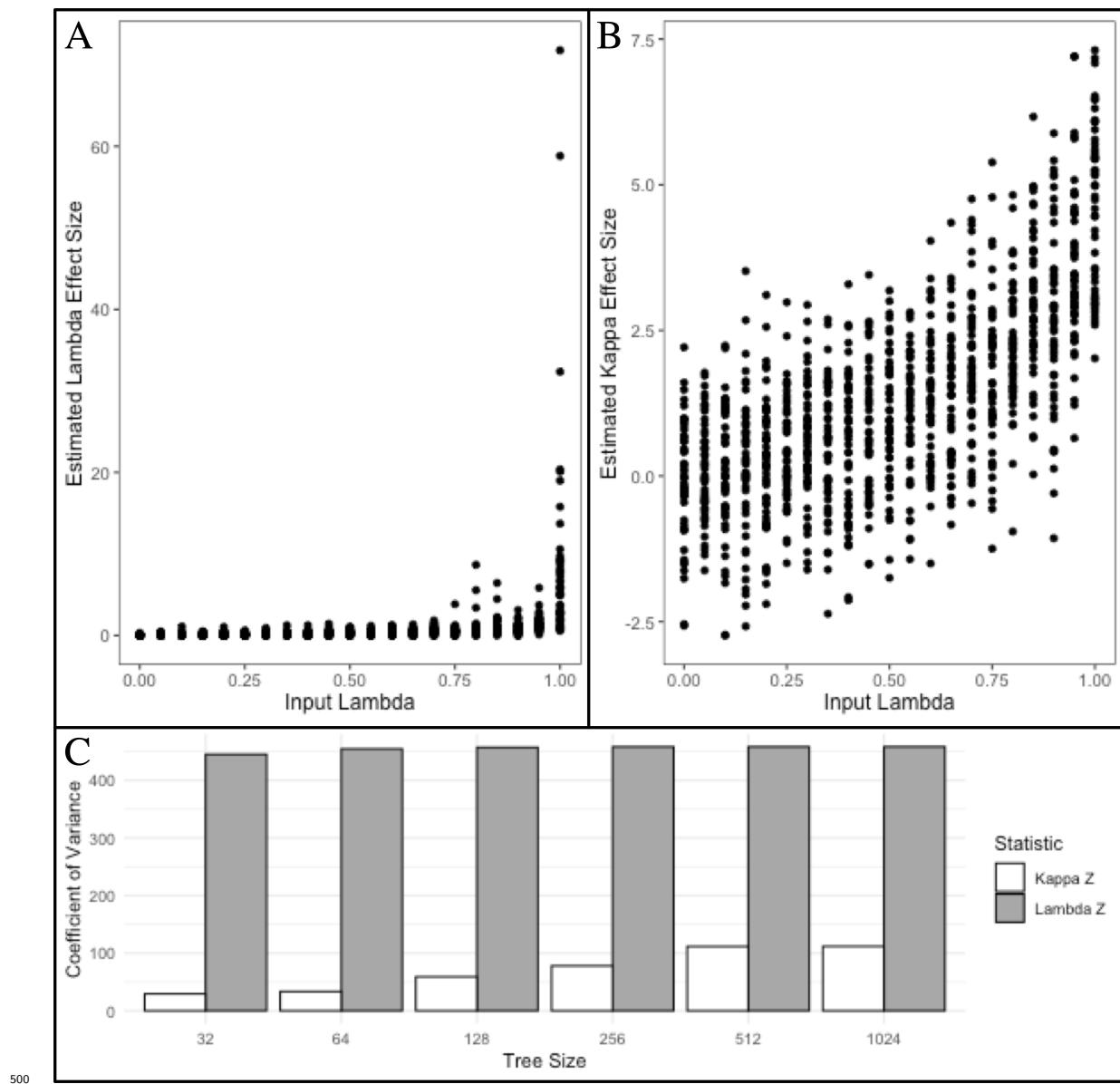
490

491 **Figure 2.** Precision of Pagel's  $\lambda$  across known levels of input phylogenetic signal ( $\lambda_{in}$ ) on phylogenies of  
 492 various sizes. As phylogenies increase in size, variation in  $\lambda_{in}$  decreases; however the precision is not  
 493 constant across the range of input levels ( $\lambda_{in} : 0 \rightarrow 1$ ), and is highest at intermediate levels of phylogenetic  
 494 signal.

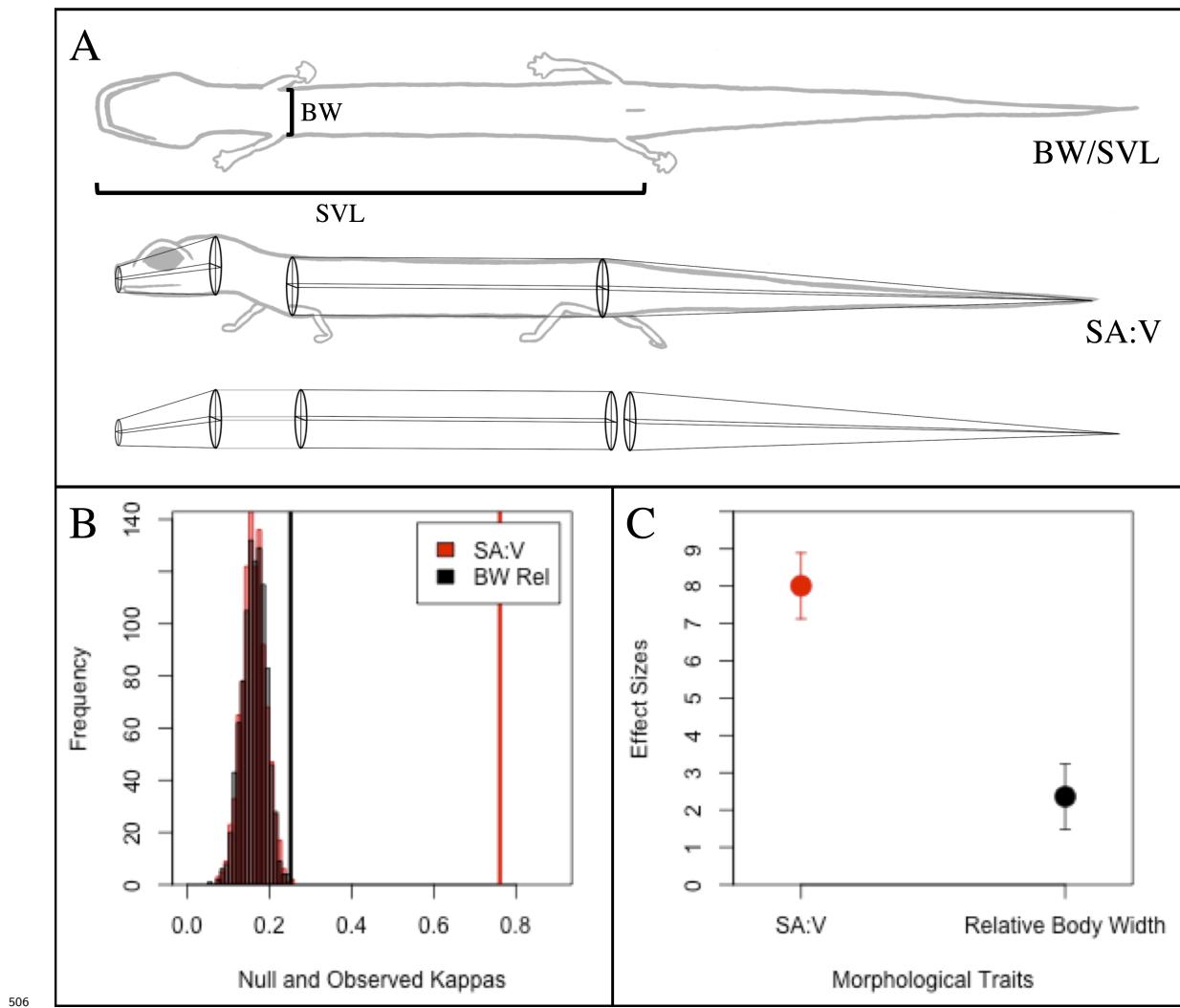


495

496 **Figure 3.** Precision of Pagel's  $\lambda$  when incorporated in phylogenetic regression ( $Y \sim X$ ), across known levels  
 497 of input phylogenetic signal ( $\lambda_{in}$ ) on phylogenies of various sizes. As phylogenies increase in size,  
 498 variation in  $\lambda_{in}$  decreases; however the precision is not constant across the range of input levels ( $\lambda_{in} : 0 \rightarrow 1$ ), and  
 499 is highest at intermediate levels of phylogenetic signal.



501 **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.  
 502 (A) Estimates  $Z_\lambda$  for data simulated on phylogenies with 32 taxa ( $n = 32$ ), (B) Estimates of  $Z_K$  for data  
 503 simulated on phylogenies with 32 taxa ( $n = 32$ ), (C) Coefficients of variation of precision estimates of  $Z_\lambda$   
 504 and  $Z_K$  across input levels of phylogenetic signal, estimated on phylogenies containing differing numbers  
 505 of species.



507 **Figure 5.** (A) Linear measures for relative body size, and regions of the body used to estimate surface  
 508 area to volume (SA:V) ratios. (B) Permutation distributions of phylogenetic signal for SA:V and  $\frac{BW}{SVL}$ ,  
 509 with observed values shown as vertical bars. (C) Effect sizes ( $Z_K$ ) for SA:V and  $\frac{BW}{SVL}$ , with their 95%  
 510 confidence intervals (CI not standardized by  $\sqrt{(n)}$ ).