

<sup>1</sup> A Standardized Effect Size for Evaluating the Strength of Phylo-  
<sup>2</sup> genetic 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  $\kappa$  ( $Z_\kappa$ ),  
<sup>13</sup> which measures the strength of phylogenetic signal, and places it on a common scale for statistical comparison.  
<sup>14</sup> Tests based on  $Z_\kappa$  provide a mechanism for formally comparing the strength of phylogenetic signal across  
<sup>15</sup> datasets, in much the same manner as effect sizes may be used to summarize patterns in quantitative meta-  
<sup>16</sup> analysis. Our approach extends the phylogenetic comparative toolkit to address hypotheses that compare  
<sup>17</sup> the strength of phylogenetic signal between various phenotypic traits, even when those traits are found in  
<sup>18</sup> different evolutionary lineages or have different units or scales.

<sup>19</sup> **Introduction**

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

<sup>33</sup>

<sup>34</sup> Several analytical tools have been developed to quantify phylogenetic signal in phenotypic datasets, including  
<sup>35</sup> measures of serial independence ( $C$ : Abouheif 1999), autocorrelation estimates ( $I$ : Gittleman and Kot 1990),  
<sup>36</sup> statistical ratios of trait variation relative to what is expected given the phylogeny ( $\kappa$ : Blomberg et al. 2003;  
<sup>37</sup> Adams 2014a), and scaling parameters used in maximum likelihood fitting of the data to the phylogeny ( $\lambda$ :  
<sup>38</sup> Pagel 1999), among others (e.g., Klingenberg and Gidaszewski 2010). The statistical properties of these  
<sup>39</sup> methods – namely type I error rates and power – have also been investigated to determine when phylogenetic  
<sup>40</sup> signal can be detected and under what conditions (e.g., Münkemüller et al. 2012; Pavoine and Ricotta 2012;  
<sup>41</sup> Diniz-Filho et al. 2012; Adams 2014a; Molina-Vegas and Rodríguez 2017; see also Revell et al. 2008; Revell  
<sup>42</sup> 2010). One of the most widely used methods for characterizing phylogenetic signal in macroevolutionary  
<sup>43</sup> studies is Pagel's  $\lambda$  (Pagel 1999). The parameter ( $\lambda$ ) transforms the lengths of the internal branches of the  
<sup>44</sup> phylogeny to improve the fit of data to the phylogeny via maximum likelihood (Pagel 1999; Freckleton et al.  
<sup>45</sup> 2002). Pagel's  $\lambda$  ranges from  $0 \rightarrow 1$ , with larger values signifying a greater dependence of observed trait  
<sup>46</sup> variation on the phylogeny. Pagel's  $\lambda$  also has the appeal that it may be included in phylogenetic generalized  
<sup>47</sup> least-squares regression (PGLS) to account for the degree of phylogenetic signal in comparative analyses (see  
<sup>48</sup> Freckleton et al. 2002).

<sup>49</sup>

50 In addition to functioning as a parameter that is tuned for appropriate analysis,  $\lambda$  can function as a descriptive  
51 statistic characterizing the relative strength of phylogenetic signal in phenotypic traits to determine the  
52 extent to which shared evolutionary history has influenced trait covariation among taxa. The appeal of  $\lambda$  as  
53 a descriptive statistic is that it serves as the basis for interpreting “weak” versus “strong” phylogenetic signal;  
54 i.e., small versus large values of  $\lambda$ , respectively, in a comparative sense (e.g., De Meester et al. 2019; Pintanel  
55 et al. 2019; Su et al. 2019). Indeed, statements regarding the strength of phylogenetic signal based on  $\lambda$   
56 are rather common in the evolutionary literature. For instance, of the 204 papers published in 2019 that  
57 estimated and reported Pagel’s  $\lambda$  (found from a literature survey we conducted in Google.scholar), 40%  
58 explicitly interpreted the strength of phylogenetic signal for at least one phenotypic trait. Further, because  
59 nearly half of the 1,572  $\lambda$  values reported were near 0 or 1 (Figure 1) where the biological interpretation of  $\lambda$   
60 is known ( $\lambda = 0$  represents no phylogenetic signal, while  $\lambda = 1$  is phylogenetic signal as expected under  
61 Brownian motion), this percentage is even higher.

62

63 [insert Figure 1 here]

64

65 Various other approaches use  $\lambda$  as a parameter that can be adjusted for inferences that are akin to sensitivity  
66 analysis. For instance, some have performed likelihood ratio tests that compare observed model fits to those  
67 obtained when  $\lambda = 0$  or  $\lambda = 1$  (Freckleton et al. 2002; Cooper et al. 2010; Bose et al. 2019) or evaluated  
68 whether observed  $\lambda$  differs from an expected  $\lambda$ , based on confidence intervals generated for the expected  
69 value (Vandekook et al. 2019). Qualitative comparisons of  $\lambda$  estimates have also been performed for multiple  
70 traits on the same phylogenetic tree to infer whether the strength of phylogenetic signal is greater in one  
71 trait as compared to another (e.g., Liu et al. 2019; Bai et al. 2019).

72

73 It seems intuitive to interpret the strength of phylogenetic signal based on the value of  $\lambda$ , as  $\lambda$  is a parameter  
74 on a bounded scale ( $0 \rightarrow 1$ ) for which interpretation of its extremal points are understood. However, equating  
75 values of  $\lambda$  directly to the strength of phylogenetic signal presumes two important statistical properties that  
76 have not been fully explored. First, it presumes that values of  $\lambda$  can be precisely estimated, as biological  
77 inferences regarding the strength of phylogenetic signal depend on high accuracy in its estimation. Therefore,  
78 understanding the precision in estimating  $\lambda$  is paramount. One study (Boettiger et al. 2012) found that  
79 estimates of Pagel’s  $\lambda$  displayed less variation (i.e., greater precision) when data were simulated on a large  
80 phylogeny ( $N = 281$ ) as compared to a small one ( $N = 13$ ). From this observation it was concluded that  
81 insufficient data (i.e., low number of species) was the underlying cause of the increased variation across

82 parameter estimates (Boettiger et al. 2012). Indeed, such a pattern is common with statistical estimators, as  
83 summary statistics and parameters are often more precise at greater sample sizes (Cohen 1988). However, this  
84 conclusion also implies that the precision of  $\lambda$  remains constant across its range ( $\lambda = 0 \rightarrow 1$ ); an assumption  
85 that to date, has not been verified. Thus, despite widespread use of Pagel's (1999)  $\lambda$  in macroevolutionary  
86 studies, at present, we lack a general understanding of the precision with which  $\lambda$  can estimate levels of  
87 phylogenetic signal in phenotypic datasets.

88

89 Second, while estimates of  $\lambda$  are within a bounded scale ( $0 \rightarrow 1$ ), this does not *de-facto* imply that the  
90 estimated values of this parameter correspond to the actual strength of the underlying input signal in the  
91 data. For this to be the case,  $\lambda$  must be a statistical effect size. Effect sizes are a measure of the magnitude  
92 of a statistical effect in data, represented on a common scale (Glass 1976; Cohen 1988). Effect sizes have  
93 widespread use in many areas of the quantitative sciences, as they represent measures that may be readily  
94 summarized across datasets as in meta-analyses (Glass 1976; Hedges and Olkin 1985; Arnqvist and Wooster  
95 1995), or compared among datasets (e.g., Adams and Collyer 2016, 2019a). Unfortunately, not all model  
96 parameters and descriptive statistics are effect sizes, and thus many summary measures must first be  
97 converted to statistics with standardized units (i.e., conversion to an effect size) for meaningful comparison  
98 (see Rosenthal 1994). As a consequence, it follows that only if  $\lambda$  is a statistical effect size can comparisons of  
99 estimates across datasets be interpretable. However, the calculation and statistical behavior of  $\lambda$  as an effect  
100 size has not yet been explored.

101

102 In this study, we evaluate the precision of Pagel's  $\lambda$  for estimating known levels of phylogenetic signal  
103 in phenotypic data. We use computer simulations with differing numbers of species, differently shaped  
104 phylogenies, and differing input levels of phylogenetic signal, to explore the degree to which  $\lambda$  correctly  
105 identifies known levels of phylogenetic signal, and under what circumstances. We find that estimates of  
106  $\lambda$  vary widely for a given input value of phylogenetic signal, and that the precision in estimating  $\lambda$  is not  
107 constant across its range. Rather, there is decreased precision when input levels of phylogenetic signal are of  
108 intermediate strength. Additionally, the same estimated values of  $\lambda$  may be obtained from datasets containing  
109 vastly different input levels of phylogenetic signal. Thus,  $\lambda$  is not a reliable indicator of the strength of  
110 phylogenetic signal in phenotypic data. We then describe a standardized effect size for measuring the strength  
111 of phylogenetic signal in phenotypic datasets and apply the concept to two common measures of phylogenetic  
112 signal:  $\lambda$  and  $\kappa$ . Through simulations we find that the precision of effect sizes based on  $\lambda$  ( $Z_\lambda$ ) are less reliable  
113 than that those based on  $\kappa$  ( $Z_\kappa$ ), implying that  $Z_\kappa$  is a more robust effect size measure. We also propose a

114 two-sample test statistic that may be used to compare the strength of phylogenetic signal among datasets,  
115 and provide an empirical example to demonstrate its use. We conclude that estimates of phylogenetic signal  
116 using Pagel's  $\lambda$  are often inaccurate, and thus interpreting strength of phylogenetic signal in phenotypic  
117 datasets based on this measure is compromised. By contrast, effect sizes obtained from  $\kappa$  hold promise for  
118 characterizing phylogenetic signal, and for comparing the strength of phylogenetic signal across datasets.

## 119 Methods and Results

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

121 We conducted a series of computer simulations to evaluate the precision of Pagel's  $\lambda$ . Our primary simulations  
122 were based on pure-birth phylogenies; however, we also evaluated patterns on both balanced and pectinate  
123 trees to determine whether tree shape affected our findings (see Supporting Information). First we generated  
124 50 pure-birth phylogenies at each of six different tree sizes, ranging from 32 to 1024 taxa ( $n = 2^5 - 2^{10}$ ).  
125 Next, we rescaled the simulated phylogenies by multiplying the internal branches by  $\lambda_{in}$ , using 21 intervals of  
126 0.05 units across its range ( $\lambda_{in} = 0.0 \rightarrow 1.0$ ), resulting in 1050 scaled phylogenies at each level of species  
127 richness ( $n$ ). Continuous traits were then simulated on each phylogeny under a Brownian motion model of  
128 evolution to obtain datasets with differing levels of phylogenetic signal, that ranged from no phylogenetic  
129 signal (when  $\lambda_{in} = 0$ ), to phylogenetic signal reflecting Brownian motion (when  $\lambda_{in} = 1$ ). For each dataset  
130 we then estimated phylogenetic signal ( $\lambda_{est}$ ), and calculated the variance of  $\lambda$  ( $\sigma_\lambda^2$ ) across datasets at each  
131 input level of phylogenetic signal and level of species richness as an estimate of precision. We verified that  
132 the variance of traits simulated had no effect on phylogenetic signal estimation.

133

134 We also evaluated the precision of  $\lambda$  when estimated in PGLS regression and ANOVA (i.e.,  $Y \sim X$ ). Here,  
135 an independent variable  $X$  was simulated on each rescaled phylogeny under a Brownian motion model of  
136 evolution (for PGLS regression). For phylogenetic ANOVA, random groups ( $X$ ) were obtained by simulating  
137 a discrete (binary, 0 or 1) character on each phylogeny. Next, the dependent variable was simulated in such a  
138 manner as to contain a known relationship with  $X$  plus random error containing phylogenetic signal. This  
139 was accomplished as:  $Y = \beta X + \epsilon$ . The association between  $Y$  and  $X$  was modeled using a range of values:  
140  $\beta = (0.0, 0.25, 0.5, 0.75, 1.0)$ , and the residual error ( $\epsilon$ ) was modeled to contain phylogenetic signal simulated  
141 under a Brownian motion model of evolution on each rescaled phylogeny:  $\epsilon = \mathcal{N}(\mu = 0, \sigma = \sigma^2 \mathbf{C})$ : (see  
142 Revell 2010 for a similar simulation design). The fit of the phylogenetic regression was estimated using

143 maximum likelihood, and parameter estimates ( $\beta_{est}$  and  $\lambda_{est}$ ) were obtained. We then calculated precision  
144 estimates ( $\sigma_\lambda^2$ ) at each input level of phylogenetic signal and level of species richness. We verified that the  
145 amount of residual variance simulated had no effect on  $\sigma_\lambda^2$  but did influence the precision of coefficients  
146 estimated from the linear model (precision increased with smaller  $\epsilon$ , as expected).

147

148 All analyses were performed in R v3.6.0 (R Core Team 2019) using the packages **geiger** (Harmon et al.  
149 2008), **caper** (Orme et al. 2013), **phytools** (Revell 2012), and **geomorph** 3.2.1 (Adams and Otárola-Castillo  
150 2013; Adams et al. 2020). R-scripts are found in the Supporting Information.

151

152 *Results.* We found that the precision of  $\lambda_{est}$  varied widely across simulation conditions. Predictably, precision  
153 improved as the number of species increased (Figure 2). This confirmed earlier findings of Boettiger et  
154 al. (2012), and adhered to parametric statistical theory. However, in many cases the set of  $\lambda_{est}$  spanned  
155 nearly the entire range of possible values (e.g.,  $n = 32$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.0 \rightarrow 0.985$ ), revealing that  
156 estimates of  $\lambda$  were not a reliable indicator of input phylogenetic signal. Importantly, the precision of  $\lambda_{est}$   
157 was not uniform across all levels of phylogenetic signal. The worst precision was observed at intermediate  
158 levels of phylogenetic signal ( $\lambda_{in} \approx 0.5$ ), while precision improved as input levels approached the extremes of  
159  $\lambda$ '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  
160 intermediate values. Additionally, even at large levels of species richness, we found that the range of  $\lambda_{est}$  still  
161 encompassed a substantial portion of possible values (e.g.,  $n = 512$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.32 \rightarrow 0.68$ ). Likewise,  
162 the same  $\lambda_{est}$  could be obtained from datasets containing vastly different input levels of phylogenetic  
163 signal (e.g.,  $n = 512$ ;  $\lambda_{est} = 0.5$ ;  $\lambda_{in} = 0.25 \rightarrow 0.65$ ). These findings were particularly unsettling when  
164 considered in light of our literature survey. Over one quarter of the  $\lambda$  estimates published in empirical  
165 studies (421 of 1,572) were between  $\lambda = 0.25$  and  $\lambda = 0.75$  (Figure 1). This range reflected the region  
166 that our simulations identified as being the least reliable in terms of accurately characterizing levels of  
167 phylogenetic signal, yet 30% of these mid-range empirical estimates were explicitly interpreted in terms of  
168 the strength of phylogenetic signal that they represented (i.e., weak, intermediate, strong phylogenetic signal).

169

170 When phylogenetic signal was investigated on balanced or pectinate trees, patterns in the precision of  $\lambda$  were  
171 largely the same, with decreased precision at intermediate levels of phylogenetic signal (Supporting Informa-  
172 tion). Likewise, when  $\lambda$  was co-estimated with regression parameters in PGLS regression and ANOVA, the  
173 results of our simulations were quite similar. Regression parameters ( $\beta$ ) were accurately estimated, confirming  
174 earlier findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal ( $\lambda$ )

175 were less precise (Figure 3; see also Supporting Information), and the spread of  $\lambda_{est}$  was similar to that  
176 observed when  $\lambda$  was estimated for only the dependent variable, as in Figure 2. Taken together, these findings  
177 reveal that  $\lambda_{est}$  does not precisely characterize known levels of phylogenetic signal in phenotypic datasets,  
178 and that biological interpretations of the strength of phylogenetic signal based on  $\lambda$  may be highly inaccurate.

179

180 [insert Figure 2 here]

181

182 [insert Figure 3 here]

183

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

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

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

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

198

199 Adams and Collyer (2019b) suggested that the strength of phylogenetic signal could be represented as an effect  
200 size based on the  $\kappa$  statistic and its empirical sampling distribution from permutation. Here we formalize

201 that suggestion, resulting in an effect size of:

$$Z_\kappa = \frac{\log(\kappa_{obs}) - \hat{\mu}_{\log(\kappa)}}{\hat{\sigma}_{\log(\kappa)}} \quad (2)$$

202 where  $\kappa_{obs}$  is the observed phylogenetic signal, and  $\hat{\mu}_\kappa$  and  $\hat{\sigma}_\kappa$  are the mean and standard deviation of the  
203 empirical sampling distribution of  $\log(\kappa)$  obtained via permutation. Note that the logarithm was used  
204 because  $\kappa$  takes only positive values ( $0 \rightarrow \infty$ ) and its sampling distribution is log-normally distributed (for a  
205 similar transformation when calculating multivariate effect sizes see: Appendix 1 of Collyer and Adams 2018).

206

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

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

208 In this case,  $\lambda_{obs}$  and  $\hat{\sigma}_\lambda$  are empirically derived using maximum likelihood, as permutation approaches have  
209 not been developed for evaluating  $\lambda$ . Note also that under the null hypothesis, no phylogenetic signal is  
210 expected (Freckleton et al. 2002), and thus  $E(\lambda) = 0$  under this condition.

211

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

222

223 [insert Figure 4 here]

224 ***Statistical Comparisons of Phylogenetic Signal***

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

$$\hat{Z}_{12} = \frac{|(\kappa_1 - \hat{\mu}_{\kappa_1}) - (\kappa_2 - \hat{\mu}_{\kappa_2})|}{\sqrt{\hat{\sigma}_{\kappa_1}^2 + \hat{\sigma}_{\kappa_2}^2}} = \frac{|Z_{\kappa_1} - Z_{\kappa_2}|}{\sqrt{2}} \quad (4)$$

229 where  $\kappa_1$ ,  $\kappa_2$ ,  $\hat{\mu}_{\kappa_1}$ ,  $\hat{\mu}_{\kappa_2}$ ,  $\hat{\sigma}_{\kappa_1}$ , and  $\hat{\sigma}_{\kappa_2}$  are as defined above for equation 2. The right side of the equation  
230 illustrates that if  $Z_\kappa$  has already been calculated for two sampling distributions as in equation 2, the sampling  
231 distributions have unit variance for each of the  $Z_\kappa$  statistics. Estimates of significance of  $\hat{Z}_{12}$  may be obtained  
232 from a standard normal distribution. Typically,  $\hat{Z}_{12}$  is considered a two-tailed test, however directional  
233 (one-tailed) tests may be specified should the empirical situation require it (see Adams and Collyer 2016,  
234 2019a).

235

236 ***Empirical Example***

237 To demonstrate the utility of  $\hat{Z}_{12}$  we quantified and compared the strength of phylogenetic signal of two  
238 phenotypic traits across species of plethodontid salamander. The data were part of a series of studies  
239 examining macroevolutionary trends in phenotypic diversification in this group (Baken and Adams 2019;  
240 Baken et al. 2020). Our dataset contained surface area to volume ratios (SA:V) and relative body width  
241 ( $\frac{BW}{SVL}$ ) for 305 species (Figure 5A). For SA:V, 11 linear body measurements from 2,781 individuals were  
242 taken, from which estimates of the surface area and volume of the head, body, and tail were calculated and  
243 subsequently combined to arrive at the SA:V for each individual (for mathematical details see Baken et al.  
244 2020). Species means were then obtained. Likewise, body size (SVL) and body width (BW) measurements  
245 were taken from 3,371 individuals, and species means of relative body width ( $\frac{BW}{SVL}$ ) were calculated (data  
246 from Baken and Adams 2019). A time-dated molecular phylogeny for the group (Bonett and Blair 2017)  
247 was then pruned to match the species in the phenotypic dataset. The phylogenetic signal in each trait was  
248 then characterized using  $\kappa$ , which was converted to its effect size ( $Z_\kappa$ ) using geomorph 3.2.1 (Adams and  
249 Otárola-Castillo 2013; Adams et al. 2020). Finally, the strength of phylogenetic signal was compared across

250 traits using  $\hat{Z}_{12}$  as described above (**to be incorporated in geomorph upon manuscript acceptance**).

251

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

## 264 Discussion

265 It is common in comparative evolutionary studies to characterize the phylogenetic signal in phenotypic traits  
266 to determine the extent to which shared evolutionary history has generated trait covariation among taxa.  
267 However, while numerous analytical approaches may be used to quantify phylogenetic signal (e.g., Abouheif  
268 1999; Gittleman and Kot 1990; Pagel 1999; Blomberg et al. 2003; Adams 2014a), methods that explicitly  
269 measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained  
270 underdeveloped. In this study, we evaluated the precision of one common measure, Pagel's  $\lambda$ , and explored its  
271 efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations,  
272 we found that the precision of  $\lambda$  increased with increasing sample sizes; a pattern noted previously (Boettiger  
273 et al. 2012), and one that conformed with parametric statistical theory (Cohen 1988). However, we also found  
274 that vastly different  $\lambda$  estimates could be obtained from data containing the same level of phylogenetic signal,  
275 and that similar  $\lambda$  estimates may be obtained from data containing differing levels of phylogenetic signal.  
276 Further, the precision of  $\lambda$  varied with the strength of phylogenetic signal, where lower precision was observed  
277 when in data whose phylogenetic signal was of intermediate strength. From these findings we conclude that  $\lambda$  is  
278 not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and that biological  
279 interpretations of the strength of signal based on this parameter may inaccurately characterize such effects.

280

281 As an alternative, we described a standardized effect size ( $Z$ ) for assessing the strength of phylogenetic signal.  
282  $Z$  expresses the magnitude of phylogenetic signal as a standard normal deviate, which is easily interpretable  
283 as the strength of phylogenetic signal relative to the mean. We applied this concept to both  $\lambda$  and  $\kappa$ , and  
284 found that  $Z_\kappa$  was a better estimate of the strength of phylogenetic signal in phenotypic data. First,  $Z_\kappa$  was  
285 more precise than  $Z_\lambda$ , and precision was more consistent across the range of input levels of phylogenetic  
286 signal. Additionally, values of  $Z_\kappa$  more accurately tracked known changes in the magnitude of phylogenetic  
287 signal, as demonstrated by the linear relationship between  $Z_\kappa$  and  $\lambda_{in}$ . Thus,  $Z_\kappa$  holds promise as a measure  
288 of the relative strength of phylogenetic signal that reflects the magnitude of this effect in phenotypic data.  
289 We therefore recommend that future studies interested in the strength of phylogenetic signal incorporate  $Z_\kappa$   
290 as a statistical measure of this effect.

291

292 Based on the effect size  $Z_\kappa$ , we then proposed a two-sample test, which provides means of determining  
293 whether the strength of phylogenetic signal is greater in one phenotypic trait as compared to another, via a  
294 hypothesis test. Prior studies have summarized patterns of variation in phylogenetic signal across datasets  
295 using summary test values, such as  $\kappa$  (e.g., Blomberg et al. 2003). However,  $\kappa$  does not scale linearly with  
296 input levels of phylogenetic signal, and its variance increases (i.e., precision decreases) with increasing  
297 strength of phylogenetic signal (Münkemüller et al. 2012; Diniz-Filho et al. 2012: see also Supporting  
298 Information). Thus,  $\kappa$  should not be considered an effect size that measures the strength of phylogenetic  
299 signal on a common scale. By contrast, standardizing  $\kappa$  ( $Z_\kappa$ , via equation 2) alleviates these concerns, and  
300 facilitates formal statistical comparisons of the strength of signal across datasets. Thus when viewed from this  
301 perspective, the approach developed here aligns well with other statistical approaches such as meta-analysis  
302 (sensu Hedges and Olkin 1985; Glass 1976; Arnqvist and Wooster 1995), where summary statistics across  
303 datasets are converted to standardized effect sizes for subsequent “higher order” statistical summaries or  
304 comparisons. As such, our approach enables evolutionary biologists to quantitatively examine the relative  
305 strength of phylogenetic signal across a wide range of phenotypic traits, and thus opens the door for future dis-  
306 coveries that inform on how phenotypic diversity accumulates in macroevolutionary time across the tree of life.

307

308 One important advantage of the approach advocated here is that the resulting effect sizes ( $Z_\kappa$ ) are  
309 dimensionless, as the units of measurement cancel out during the calculation of  $Z$  (Sokal and Rohlf 2012).  
310 Thus,  $Z_\kappa$  represents the strength of phylogenetic signal on a common and comparable scale – measured  
311 in standard deviations – regardless of the initial units and original scale of the phenotypic variables under

312 investigation. This means that the strength of phylogenetic signal may be compared across datasets for  
313 continuous phenotypic traits measured in different units and scale, because those units have been standardized  
314 through their conversion to  $Z_\kappa$ . For example, our approach could be utilized to determine whether the  
315 strength of phylogenetic signal (say, in response to ecological differentiation) is stronger in morphological  
316 traits (linear traits:  $mm$ ), physiological traits (metabolic rate:  $\frac{O^2}{min}$ ), or behavioral traits (aggression:  
317  $\frac{\#displays}{second}$ ). In fact, our empirical example provided such a comparison, as SA:V is represented in  $mm^{-1}$   
318 while relative body size is a unitless ratio ( $\frac{BW}{SVL}$ ). Additionally, our method is capable of comparing the  
319 strength of phylogenetic signal in traits of different dimensionality, as estimates of phylogenetic signal using  $\kappa$   
320 have been generalized for multivariate data ( $\kappa_{mult}$ : see Adams 2014a). Furthermore, tests based on  $\hat{Z}_{12}$  may  
321 be utilized for comparing the strength of phylogenetic signal among datasets containing a different number  
322 of species, and even for phenotypes obtained from species in different lineages, because their phylogenetic  
323 non-independence and observed variation are taken into account in the generation of the empirical sampling  
324 distribution via permutation.

325

326 This study is not the first to compare  $\lambda$  and  $\kappa$  for their ability as statistics to measure phylogenetic signal.  
327 Our results for  $\lambda$  and  $\kappa$  values are consistent with those found in the simulations performed by Münkemüller  
328 et al. (2012), but that study investigated type I error rates and statistical power, finding that  $\lambda$  performed  
329 better in both regards, irrespective of species number in trees. Although not the central focus of their study,  
330 the same tendency for variable  $\lambda$  and consistent  $\kappa$  at intermediate phylogenetic signal strengths was observed  
331 (see Fig. 2, Münkemüller et al. 2012). Recent work by Molina-Venegas and Rodríguez (2017) found that  
332  $\kappa$  but not  $\lambda$  tended to inflate the estimate of phylogenetic signal, leading to moderate type I and type II  
333 biases, if polytomic chronograms were used. Their work more thoroughly addressed previous observations of  
334 inflated  $\kappa$  for incompletely resolved phylogenetic trees (Davies et al. 2012; Münkemüller et al. 2012). An  
335 interesting question is whether an inflated  $\kappa$  value leads to an inflated  $Z_\kappa$  or does a tendency of a particular  
336 tree to inflate estimates of  $\kappa$  also inflate the values in random permutations of a test, in which case  $Z_\kappa$  is  
337 robust to polytomies? We repeated the analyses in Figure 4, adjusting trees to have 50% collapsed nodes, per  
338 the technique of Molina-Venegas and Rodríguez (2017), and found results were consistent (see Supporting  
339 Information). This confirms that any tendency of incompletely resolved trees to inflate  $\kappa$  as a descriptive  
340 statistic does not inflate  $Z_\kappa$  as an effect size. Furthermore, because comparison of effect sizes in a test is a  
341 comparison of locations of observed values in their sampling distributions, which would shift concomitantly  
342 because of this tendency, the  $Z_{12}$  test statistic in equation 4 appears to be robust in spite of unresolved trees.

343

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

361    **References**

- 362    Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.  
363                      Evolutionary Ecology Research 1:895–909.
- 364    Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other  
365                      high-dimensional multivariate data. Systematic Biology 63:685–697.
- 366    Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other  
367                      high-dimensional multivariate data. Evolution 68:2675–2688.
- 368    Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative  
369                      modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- 370    Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration  
371                      across morphometric datasets. Evolution 70:2623–2631.
- 372    Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,  
373                      and a refined permutation procedure. Evolution 72:1204–1215.
- 374    Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate  
375                      phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- 376    Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric  
377                      analyses. R package version 3.2.1.
- 378    Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of  
379                      geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- 380    Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.  
381                      Trends in Ecology and Evolution 10:236–240.
- 382    Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with  
383                      phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil  
384                      434:305–326.
- 385    Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution  
386                      9:7005–7016.

- 387 Baken, E. K., L. E. Mellenthin, and D. C. Adams. 2020. Macroevolution of desiccation-related morphology  
388 in plethodontid salamanders as inferred from a novel surface area to volume ratio estimation approach.  
389 *Evolution* 74:476–486.
- 390 Beaulieu, J. M., D. C. Jhwueng, C. Boettiger, and B. C. O'Meara. 2012. Modeling stabilizing selection:  
391 Expanding the ornstein-uhlenbeck model of adaptive evolution. *Evolution* 66:2369–2383.
- 392 Blomberg, S. P., T. Garland, and A. R. Ives. 2003. Testing for phylogenetic signal in comparative data:  
393 Behavioral traits are more labile. *Evolution* 57:717–745.
- 394 Boettiger, C., G. Coop, and P. Ralph. 2012. Is your phylogeny informative? Measuring the power of  
395 comparative methods. *Evolution* 67:2240–2251.
- 396 Bonett, R. M., and A. L. Blair. 2017. Evidence for complex life cycle constraints on salamander body form  
397 diversification. *Proceedings of the National Academy of Sciences, U.S.A.* 114:9936–9941.
- 398 Bose, R., B. R. Ramesh, R. Pélassier, and F. Munoz. 2019. Phylogenetic diversity in the western ghats  
399 biodiversity hotspot reflects environmental filtering and past niche diversification of trees. *Journal of*  
400 *Biogeography* 46:145–157.
- 401 Butler, M. A., and A. A. King. 2004. Phylogenetic comparative analysis: A modeling approach for adaptive  
402 evolution. *American Naturalist* 164:683–695.
- 403 Cohen, J. 1988. Statistical power analysis for the behavioral sciences. Routledge.
- 404 Collyer, M. L., and D. C. Adams. 2018. RRPP: An r package for fitting linear models to high-dimensional  
405 data using residual randomization. *Methods in Ecology and Evolution* 9:1772–1779.
- 406 Collyer, M. L., D. J. Sekora, and D. C. Adams. 2015. A method for analysis of phenotypic change for  
407 phenotypes described by high-dimensional data. *Heredity* 115:357–365.
- 408 Cooper, N., W. Jetz, and R. P. Freckleton. 2010. Phylogenetic comparative approaches for studying niche  
409 conservatism. *Journal of Evolutionary Biology* 23:2529–2539.
- 410 Davies, T. J., N. J. Kraft, N. Salamin, and E. M. Wolkovich. 2012. Incompletely resolved phylogenetic trees  
411 inflate estimates of phylogenetic conservatism. *Ecology* 93:242–247. Wiley Online Library.
- 412 De Meester, G., K. Huyghe, and R. Van Damme. 2019. Brain size, ecology and sociality: A reptilian  
413 perspective. *Biological Journal of the Linnean Society* 126:381–391.

- 414 Diniz-Filho, J. A. F., T. Santos, T. F. Rangel, and L. M. Bini. 2012. A comparison of metrics for estimating  
415 phylogenetic signal under alternative evolutionary models. *Genetics and Molecular Biology* 35:673–679.
- 416 Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* 125:1–15.
- 417 Freckleton, R. P., P. H. Harvey, and M. Pagel. 2002. Phylogenetic analysis and comparative data: A test and  
418 review of evidence. *American Naturalist* 160:712–726.
- 419 Garland, T. J., and A. R. Ives. 2000. Using the past to predict the present: Confidence intervals for regression  
420 equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
- 421 Gittleman, J. L., and M. Kot. 1990. Adaptation: Statistics and a null model for estimating phylogenetic  
422 effects. *Systematic Zoology* 39:227–241.
- 423 Glass, G. V. 1976. Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
- 424 Grafen, A. 1989. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London B,*  
425 *Biological Sciences* 326:119–157.
- 426 Harmon, L. J., J. T. Weir, C. D. Brock, R. E. Glor, and W. Challenger. 2008. GEIGER: Investigating  
427 evolutionary radiations. *Bioinformatics* 24:129–131.
- 428 Harvey, P. H., and M. D. Pagel. 1991. The comparative method in evolutionary biology. Oxford University  
429 Press, Oxford.
- 430 Hedges, L. V., and I. Olkin. 1985. Statistical methods for meta-analysis. Elsevier.
- 431 Kelley, K., and K. J. Preacher. 2012. On effect size. *Psychological Methods* 17:137–152.
- 432 Klingenberg, C. P., and N. A. Gidaszewski. 2010. Testing and quantifying phylogenetic signals and homoplasy  
433 in morphometric data. *Systematic biology* 59:245–261.
- 434 Liu, H., C. P. Osborne, D. Yin, R. P. Freckleton, G. Jiang, and M. Liu. 2019. Phylogeny and ecological  
435 processes influence grass coexistence at different spatial scales within the steppe biome. *Oecologia*  
436 191:25–38.
- 437 Martins, E. P., and T. F. Hansen. 1997. Phylogenies and the comparative method: A general approach  
438 to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist*  
439 149:646–667.
- 440 Molina-Venegas, R., and M. A. Rodríguez. 2017. Revisiting phylogenetic signal; strong or negligible impacts

- 441 of polytomies and branch length information? BMC evolutionary biology 17:53.
- 442 Münkemüller, T., S. Lavergne, B. Bzeznik, S. Dray, T. Jombart, K. Schiffers, and W. Thuiller. 2012. How to  
443 measure and test phylogenetic signal. Methods in Ecology and Evolution 3:743–756.
- 444 O'Meara, B. C., C. Ane, M. J. Sanderson, and P. C. Wainwright. 2006. Testing for different rates of  
445 continuous trait evolution using likelihood. Evolution 60:922–933.
- 446 Orme, D., R. P. Freckleton, G. H. Thomas, T. Petzoldt, S. A. Fritz, and N. Isaac. 2013. CAPER: Comparative  
447 analyses of phylogenetics and evolution in r. Methods in Ecology and Evolution 3:145–151.
- 448 Pagel, M. D. 1999. Inferring the historical patterns of biological evolution. Nature 401:877–884.
- 449 Pavoine, S., and C. Ricotta. 2012. Testing for phylogenetic signal in biological traits: The ubiquity of  
450 cross-product statistics. Evolution: International Journal of Organic Evolution 67:828–840.
- 451 Pintanel, P., M. Tejedo, S. R. Ron, G. A. Llorente, and A. Merino-Viteri. 2019. Elevational and microclimatic  
452 drivers of thermal tolerance in andean pristimantis frogs. Journal of Biogeography 46:1664–1675.
- 453 R Core Team. 2019. R: A language and environment for statistical computing. R Foundation for Statistical  
454 Computing, Vienna, Austria.
- 455 Revell, L. J. 2010. Phylogenetic signal and linear regression on species data. Methods in Ecology and  
456 Evolution 1:319–329.
- 457 Revell, L. J. 2012. Phytools: An r package for phylogenetic comparative biology (and other things). Methods  
458 in Ecology and Evolution 3:217–223.
- 459 Revell, L. J., and L. J. Harmon. 2008. Testing quantitative genetic hypotheses about the evolutionary rate  
460 matrix for continuous characters. Evolutionary Ecology Research 10:311–331.
- 461 Revell, L. J., L. J. Harmon, and D. C. Collar. 2008. Phylogenetic signal, evolutionary process, and rate.  
462 Systematic Biology 57:591–601.
- 463 Rohlf, F. J. 2001. Comparative methods for the analysis of continuous variables: Geometric interpretations.  
464 Evolution 55:2143–2160.
- 465 Rosenthal, R. 1994. The handbook of research synthesis. Pp. 231–244 in L. V. Cooper H Hedges, ed. Russell  
466 Sage Foundation.
- 467 Sokal, R. R., and F. J. Rohlf. 2012. Biometry. 4th ed. W.H. Freeman & Co., San Francisco.

<sup>468</sup> Su, G., S. Villéger, and S. Brosse. 2019. Morphological diversity of freshwater fishes differs between realms,  
<sup>469</sup> but morphologically extreme species are widespread. *Global ecology and biogeography* 28:211–221.

<sup>470</sup> Vandelook, F., S. Janssens, P. Gijbels, E. Fischer, W. Van den Ende, O. Honnay, and S. Abrahameczyk. 2019.  
<sup>471</sup> Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124:269–279.

472      **Figure Legends**

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

475

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

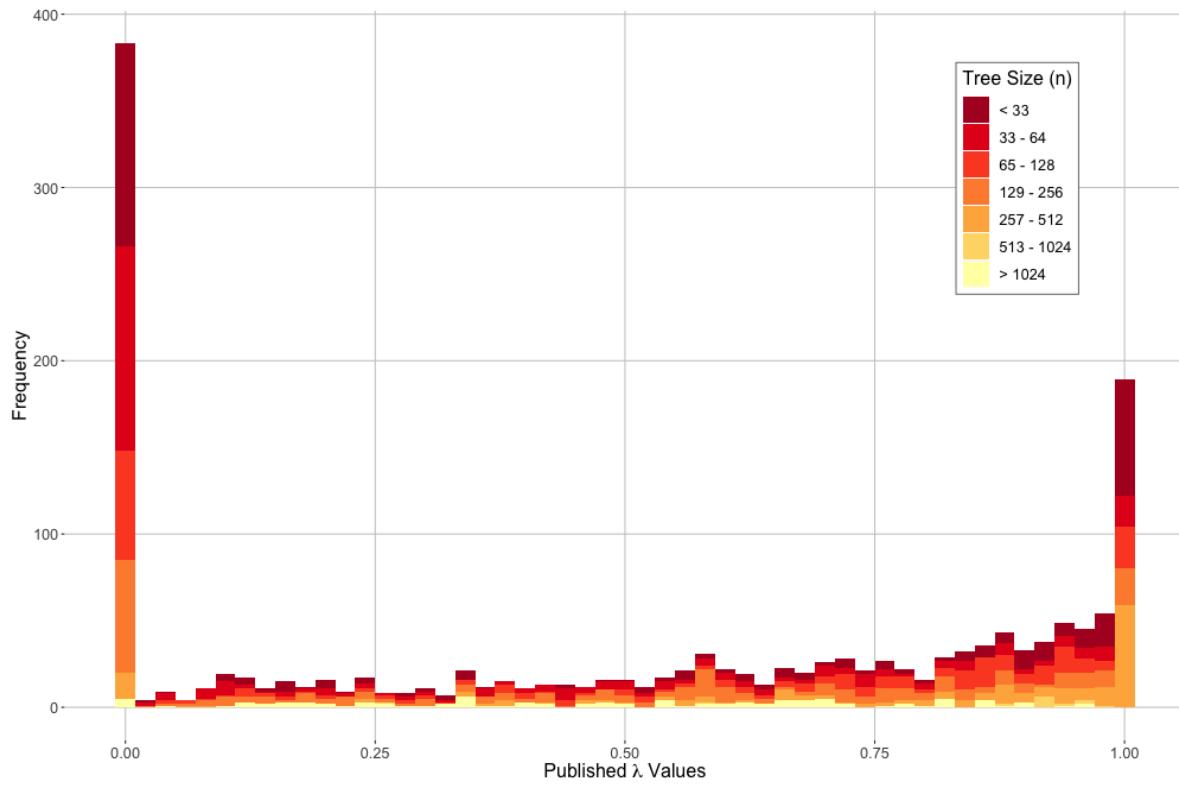
480

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

485

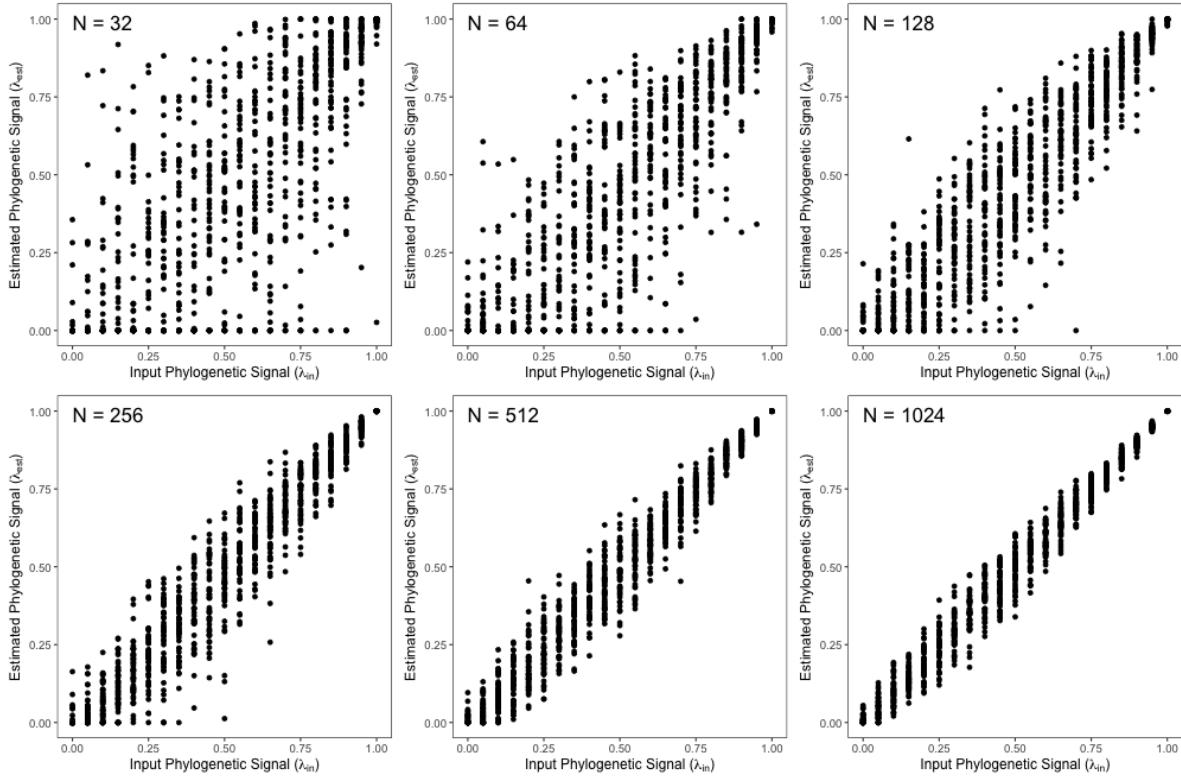
486      **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.  
487      (A) Estimates  $Z_\lambda$  for data simulated on phylogenies with 32 taxa ( $n = 32$ ), (B) Estimates of  $Z_\kappa$  for data  
488      simulated on phylogenies with 32 taxa ( $n = 32$ ), (C) Coefficients of variation of precision estimates of  $Z_\lambda$   
489      and  $Z_\kappa$  across input levels of phylogenetic signal, estimated on phylogenies containing differing numbers  
490      of species.

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



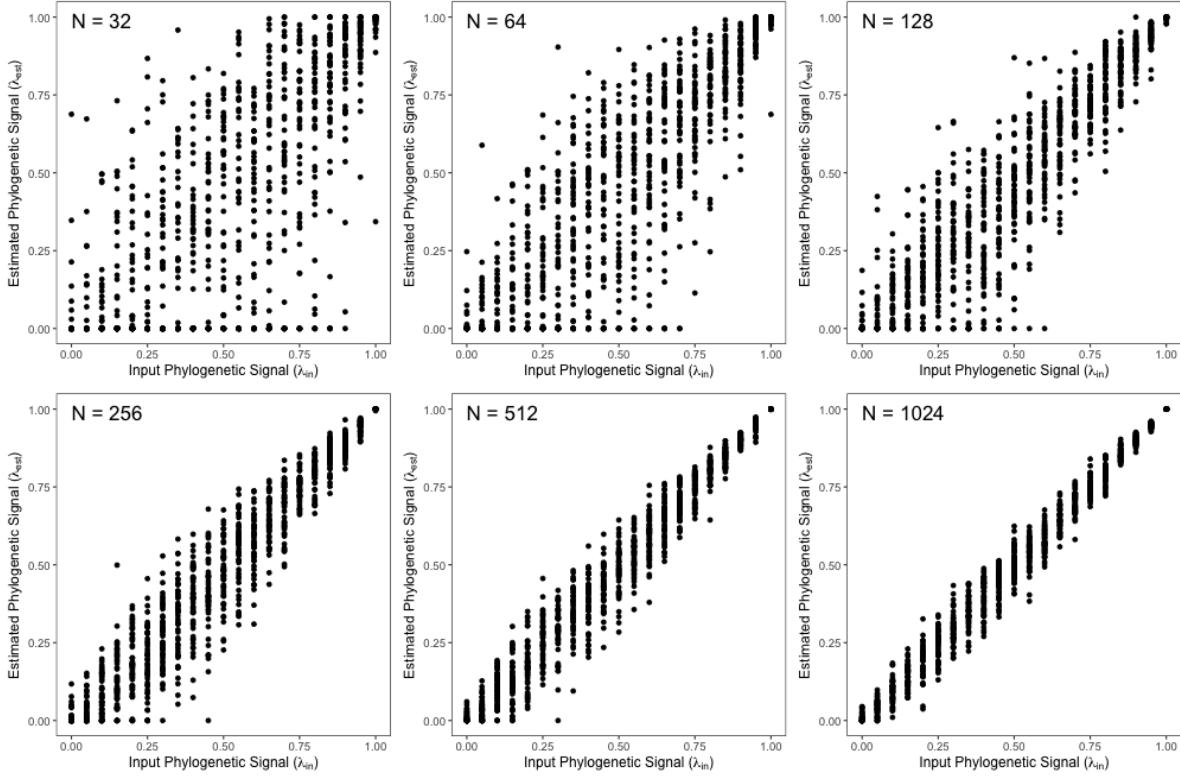
495

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



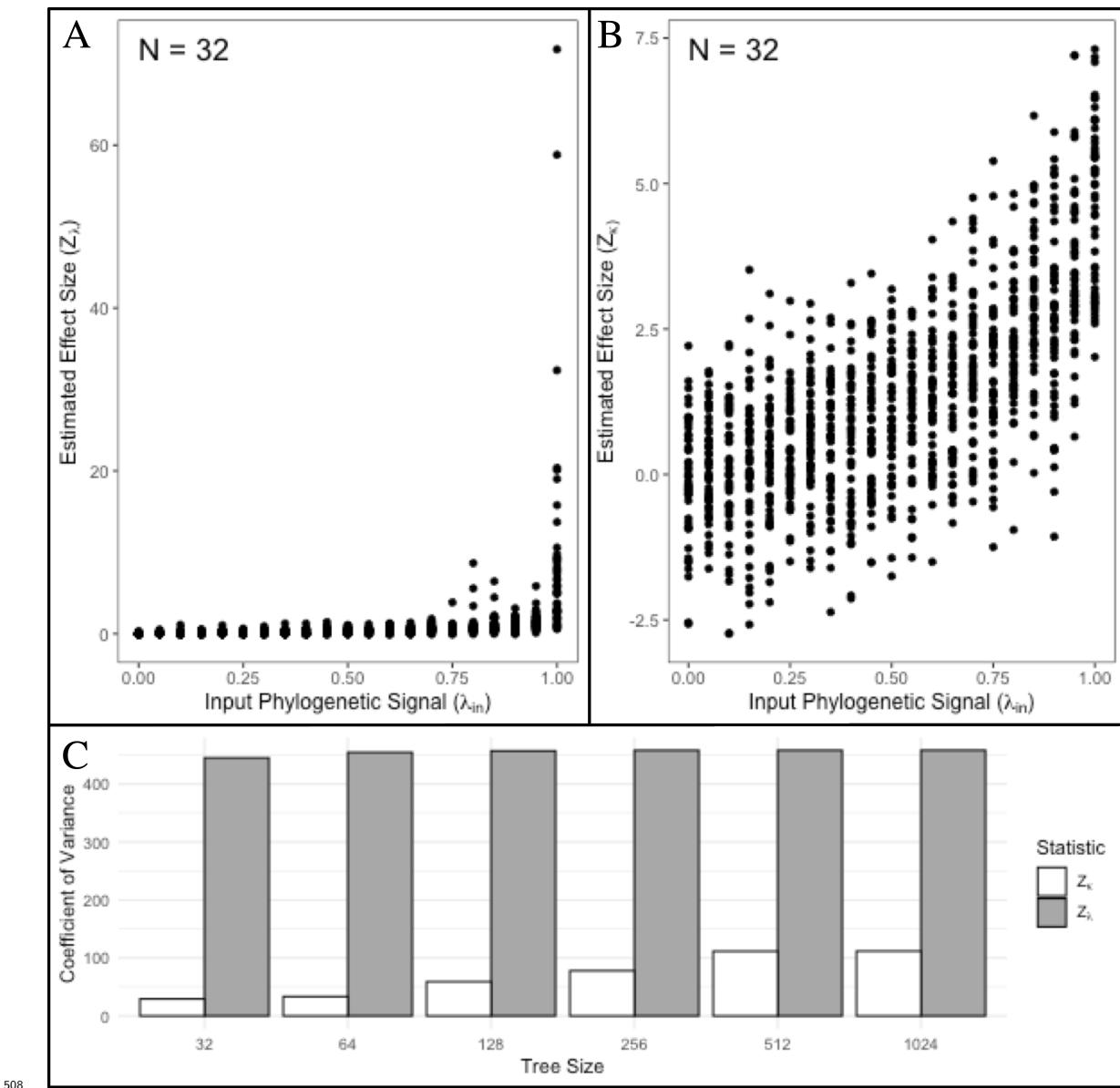
498

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

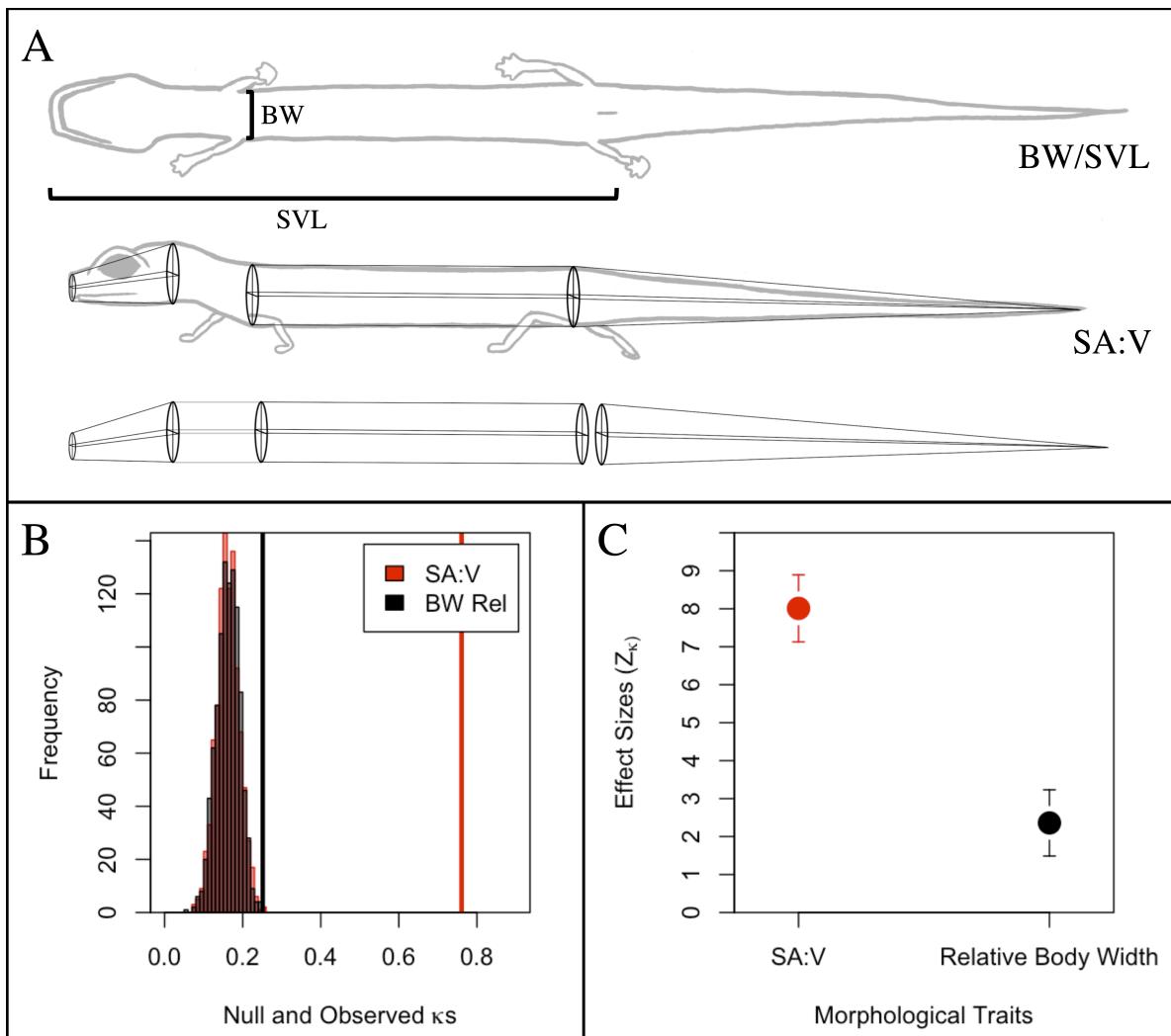


503

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



509 **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.  
 510 (A) Estimates  $Z_\lambda$  for data simulated on phylogenies with 32 taxa ( $n = 32$ ), (B) Estimates of  $Z_\kappa$  for data  
 511 simulated on phylogenies with 32 taxa ( $n = 32$ ), (C) Coefficients of variation of precision estimates of  $Z_\lambda$   
 512 and  $Z_\kappa$  across input levels of phylogenetic signal, estimated on phylogenies containing differing numbers  
 513 of species.



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