

<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*  
<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 various phenotypic traits, even when those traits are found  
<sup>18</sup> in different evolutionary lineages or of in different units or scale.

<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 generates statistical non-independence (Felsenstein 1985; Harvey and  
<sup>22</sup> Pagel 1991). Accounting for this evolutionary non-independence is the purview of *phylogenetic comparative*  
<sup>23</sup> *methods* (PCMs); a suite of analytical tools that condition the data on the phylogeny through the course  
<sup>24</sup> of statistical evaluations of phenotypic trends (e.g., Grafen 1989; Garland and Ives 2000; Rohlf 2001;  
<sup>25</sup> Butler and King 2004). The past several decades have witnessed a rapid expansion in the development  
<sup>26</sup> of PCMs to address an ever-growing set of macroevolutionary hypotheses (Martins and Hansen 1997;  
<sup>27</sup> O'Meara et al. 2006; Revell and Harmon 2008; Beaulieu et al. 2012; Adams 2014b,a; Adams and  
<sup>28</sup> Collyer 2018). These methods are predicated on the notion that phylogenetic signal – the tendency for  
<sup>29</sup> 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.  
<sup>37</sup> 2003; Adams 2014a), and scaling parameters used in maximum likelihood fitting of the data to the phylogeny  
<sup>38</sup> ( $\lambda$ : 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., Munkemuller et al. 2012; Pavoine and Ricotta 2012;  
<sup>41</sup> Diniz-Filho et al. 2012; Adams 2014a; Molina-Vegas and Rodriguez 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). Here, maximum likelihood is used to fit the data to the phylogeny under  
<sup>44</sup> a Brownian motion model of evolution. A parameter ( $\lambda$ ) is included, which transforms the lengths of the  
<sup>45</sup> internal branches of the phylogeny to improve the fit (Pagel 1999; Freckleton et al. 2002). Pagel's  $\lambda$  ranges  
<sup>46</sup> from 0 → 1, with larger values signifying a greater dependence of observed trait variation on the phylogeny.  
<sup>47</sup> Pagel's  $\lambda$  also has the appeal that it may be included in phylogenetic regression (PGLS) to account for the  
<sup>48</sup> degree of phylogenetic signal in comparative analyses (see Freckleton et al. 2002).

<sup>49</sup>

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

65

66

67 [insert Figure 1 here]

68

69 It seems intuitive to interpret the strength of phylogenetic signal based on the value of  $\lambda$ , as  $\lambda$  is a parameter  
70 on a bounded scale ( $0 \rightarrow 1$ ) for which interpretation of its extremal points are understood. Specifically,  
71  $\lambda = 0$  represents no phylogenetic signal, while  $\lambda = 1$  is phylogenetic signal as expected under Brownian  
72 motion. However, equating values of  $\lambda$  directly to the strength of phylogenetic signal presumes two important  
73 statistical properties that have not been fully explored. First, it presumes that values of  $\lambda$  can be precisely  
74 estimated, as biological inferences regarding the strength of phylogenetic signal depend on high accuracy in  
75 its estimation. Therefore, understanding the precision in estimating  $\lambda$  is paramount. One study (Boettiger et  
76 al. 2012) found that estimates of Pagel’s  $\lambda$  displayed less variation (i.e., greater precision) when data were  
77 simulated on a large phylogeny ( $N = 281$ ) as compared to a small one ( $N = 13$ ). From this observation it  
78 was concluded that insufficient data (i.e., the number of species) was the underlying cause of the increased  
79 variation across parameter estimates (Boettiger et al. 2012). Indeed, such a pattern is common with  
80 statistical estimators, as summary statistics and parameters are often more precise at greater sample sizes  
81 (Cohen 1988). However, this conclusion also implies that the precision of  $\lambda$  remains constant across its range

82 ( $\lambda = 0 \rightarrow 1$ ); an assumption that to date, has not been verified. Thus, despite widespread use of Pagel's  
83 (1999)  $\lambda$  in macroevolutionary studies, at present, we lack a general understanding of the precision with  
84 which  $\lambda$  can estimate levels of phylogenetic signal in phenotypic datasets.

85

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

97

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

114 hold promise for characterizing phylogenetic signal, and for comparing the strength of phylogenetic signal  
115 across datasets.

## 116 Methods and Results

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

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

129

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

141

142 All analyses were performed in R v3.6.0 (R Core Team 2019) using the packages **geiger** (Harmon et al.

<sup>143</sup> 2008), **caper** (Orme et al. 2013), **phytools** (Revell 2012), and **geomorph** 3.2.1 (Adams and Otárola-Castillo  
<sup>144</sup> 2013; Adams et al. 2020). R-scripts are found in the Supporting Information.

<sup>145</sup>

<sup>146</sup> *Results.* We found that the precision of  $\lambda_{est}$  varied widely across simulation conditions. Predictably, precision  
<sup>147</sup> improved as the number of species increased (Figure 2). This confirmed earlier findings of Boettiger et  
<sup>148</sup> al. (2012), and adhered to parametric statistical theory. However, in many cases the set of  $\lambda_{est}$  spanned  
<sup>149</sup> nearly the entire range of possible values (e.g.,  $n = 32$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.0 \rightarrow 0.985$ ), revealing that  
<sup>150</sup> estimates of  $\lambda$  were not a reliable indicator of input phylogenetic signal. Importantly, the precision of  $\lambda_{est}$   
<sup>151</sup> was not uniform across all levels of phylogenetic signal. The worst precision was observed at intermediate  
<sup>152</sup> levels of phylogenetic signal ( $\lambda_{in} \approx 0.5$ ), while precision improved as input levels approached the extremes of  
<sup>153</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>154</sup> intermediate values. Additionally, even at large levels of species richness, we found that the range of  $\lambda_{est}$  still  
<sup>155</sup> encompassed a substantial portion of possible values (e.g.,  $n = 512$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.32 \rightarrow 0.68$ ). Likewise,  
<sup>156</sup> the same  $\lambda_{est}$  could be obtained from datasets containing vastly different input levels of phylogenetic  
<sup>157</sup> signal (e.g.,  $n = 512$ ;  $\lambda_{est} = 0.5$ ;  $\lambda_{in} = 0.25 \rightarrow 0.65$ ). These findings were particularly unsettling when  
<sup>158</sup> considered in light of our literature survey. Over one quarter of the  $\lambda$  estimates obtained in empirical  
<sup>159</sup> studies (421 of 1,572) were between  $\lambda = 0.25$  and  $\lambda = 0.75$  (Figure 1). This range reflected the region  
<sup>160</sup> that our simulations identified as being the least reliable in terms of accurately characterizing levels of  
<sup>161</sup> phylogenetic signal, yet 30% of these mid-range empirical estimates were explicitly interpreted in terms of  
<sup>162</sup> the strength of phylogenetic signal that they represented (i.e., weak, intermediate, strong phylogenetic signal).

<sup>163</sup>

<sup>164</sup> Finally, when  $\lambda$  was co-estimated with regression parameters in PGLS regression and ANOVA, the results of  
<sup>165</sup> our simulations were quite similar. Here, regression parameters ( $\beta$ ) were accurately estimated, confirming  
<sup>166</sup> earlier findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal ( $\lambda$ )  
<sup>167</sup> were less precise (Figure 3; see also Supporting Information), and the spread of  $\lambda_{est}$  was similar to that  
<sup>168</sup> observed when  $\lambda$  was estimated for only the dependent variable, as in Figure 2. Taken together, these findings  
<sup>169</sup> reveal that  $\lambda_{est}$  does not precisely characterize observed levels of phylogenetic signal in phenotypic datasets,  
<sup>170</sup> and that biological interpretations of the strength of phylogenetic signal based on  $\lambda$  may be highly inaccurate.

<sup>171</sup>

<sup>172</sup> [insert Figure 2 here]

<sup>173</sup>

<sup>174</sup> [insert Figure 3 here]

176 **A Standardized Effect Size for Phylogenetic Signal**

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

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

183 where  $\theta_{obs}$  is the observed test statistic,  $E(\theta)$  is its expected value under the null hypothesis, and  $\sigma_\theta$  is its  
 184 standard error (Glass 1976; Cohen 1988; Rosenthal 1994).  $Z_\theta$  expresses the magnitude of the effect in  $\theta_{obs}$   
 185 by transforming the original test statistic to its standard normal deviate (Glass 1976; Kelley and Preacher  
 186 2012). Typically,  $\theta_{obs}$  and  $\sigma_\theta$  are estimated from the data, while  $E(\theta)$  is obtained from the distribution of  $\theta$   
 187 derived from parametric theory. However, recent advances in resampling theory (Collyer et al. 2015; Adams  
 188 and Collyer 2016, 2019a) have shown that  $E(\theta)$  and  $\sigma_\theta$  may also be obtained from an empirical sampling  
 189 distribution of  $\theta$  obtained from permutation procedures.

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

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

194 where  $K_{obs}$  is the observed phylogenetic signal, and  $\hat{\mu}_K$  and  $\hat{\sigma}_K$  are the mean and standard deviation of  
 195 the empirical sampling distribution of  $\log(Kappa)$  obtained via permutation. Here the logarithm was used  
 196 because *Kappa* takes only positive values ( $0 \rightarrow \infty$ ) and its sampling distribution is log-normally distributed

<sub>197</sub> (for a similar transformation when calculating multivariate effect sizes see: Appendix 1 of Collyer and Adams  
<sub>198</sub> 2018). Similarly, an effect size based on  $\lambda$  could be envisioned as:

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

<sub>199</sub> In this case,  $\lambda_{obs}$  and  $\hat{\sigma}_\lambda$  are empirically derived using maximum likelihood, as permutation approaches have  
<sub>200</sub> not been developed for evaluating  $\lambda$ . Note also that under the null hypothesis, no phylogenetic signal is  
<sub>201</sub> expected (Freckleton et al. 2002), and thus  $E(\lambda) = 0$  under this condition.

<sub>202</sub>

<sub>203</sub> To evaluate the utility of  $Z_K$  and  $Z_\lambda$  we calculated both effect sizes for the simulated datasets generated  
<sub>204</sub> above, and summarized the precision of each using its variance ( $\sigma_{Z_K}^2$  and  $\sigma_{Z_\lambda}^2$ ). Results are found in Figure 4  
<sub>205</sub> (additional results are found in the Supporting Information). Here two things are evident. First, estimates of  
<sub>206</sub>  $Z_K$  track the input phylogenetic signal in a more linear fashion than do estimates of  $Z_\lambda$  (Figure 4A,B). Thus,  
<sub>207</sub> actual changes in the strength of phylogenetic signal are reflected more evenly in the corresponding values of  
<sub>208</sub> the effect size  $Z_K$ . Second, the precision of  $Z_K$  is considerably more stable as compared with  $Z_\lambda$ . This may  
<sub>209</sub> be seen by calculating the coefficients of variation for the set of precision estimates (i.e.,  $\sigma_{Z_K}^2$  and  $\sigma_{Z_\lambda}^2$ ) across  
<sub>210</sub> input levels of phylogenetic signal. Here coefficients of variation in the precision of  $Z_K$  were an order of  
<sub>211</sub> magnitude smaller for than for  $Z_\lambda$  (Figure 4C). This implied that estimates of the strength of phylogenetic  
<sub>212</sub> signal were more reliable and robust when using  $Z_K$  as compared with  $Z_\lambda$ .

<sub>213</sub>

<sub>214</sub> [insert Figure 4 here]

### <sub>215</sub> *Statistical Comparisons of Phylogenetic Signal*

<sub>216</sub> Once the magnitude of phylogenetic signal is characterized using  $Z_K$ , one may wish to compare such measures  
<sub>217</sub> across datasets, to determine whether the strength of phylogenetic signal is significantly greater in one  
<sub>218</sub> phenotypic trait as compared to another. As with other effect sizes derived from permutation distributions  
<sub>219</sub> (e.g., Adams and Collyer 2016, 2019a), a two-sample test statistic may be calculated as:

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

220 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. Estimates of significance of  
 221  $\hat{Z}_{12}$  may be obtained from a standard normal distribution. Typically,  $\hat{Z}_{12}$  is considered a two-tailed test,  
 222 however directional (one-tailed) tests may be specified should the empirical situation require it (see Adams  
 223 and Collyer 2016, 2019a).

224

225 ***Empirical Example***

226 To demonstrate the utility of  $\hat{Z}_{12}$  we performed an analysis of the strength of phylogenetic signal in two  
 227 phenotypic datasets from species of plethodontid salamander. The data were part of a series of studies  
 228 examining macroevolutionary trends in phenotypic diversification in this group (Baken and Adams 2019;  
 229 Baken et al. 2020). Our dataset contained surface area to volume ratios (SA:V) and relative body width  
 230 ( $\frac{BW}{SVL}$ ) for 305 species (Figure 5A). For SA:V, 11 linear body measurements from 2,781 individuals were  
 231 taken, from which estimates of the surface area and volume of the head, body, and tail were calculated and  
 232 subsequently combined to arrive at the SA:V for each individual (for mathematical details see Baken et al.  
 233 2020). Species means were then obtained. Likewise, body size (SVL) and body width (BW) measurements  
 234 were taken from 3,371 individuals, and species means of relative body width ( $\frac{BW}{SVL}$ ) were calculated (data  
 235 from Baken and Adams 2019). A time-dated molecular phylogeny for the group (Bonett and Blair 2017) was  
 236 then pruned to match the species in the dataset, resulting in a phylogeny and corresponding phenotypic  
 237 dataset containing 305 species. The phylogenetic signal in each trait was then characterized using  $Kappa$ ,  
 238 which was converted to its effect size ( $Z_K$ ) using geomorph 3.2.1 (Adams and Otárola-Castillo 2013; Adams  
 239 et al. 2020). Finally, the strength of phylogenetic signal was compared across traits using  $\hat{Z}_{12}$  as described  
 240 above (to be incorporated in geomorph upon manuscript acceptance).

241

242 *Results.* Both SA:V and relative body width displayed significant phylogenetic signal ( $Kappa_{SA:V} = 0.7608$ ;  
 243  $P = 0.001$ ;  $Kappa_{BW/SVL} = 0.2515$ ;  $P = 0.001$ ). For both phenotypic traits,  $K_{obs}$  differed markedly from  
 244 their corresponding permutation distributions, which were found to overlap almost perfectly (Figure 5B).  
 245 However, while both traits displayed significant phylogenetic signal, there was nearly a four-fold difference

246 in the magnitude of their effect sizes, with SA:V displaying the greater phylogenetic signal (Figure 5C).  
247 Using the two-sample test statistic above, this difference was found to be highly significant ( $\hat{Z}_{12} = 4.13$ ;  
248  $P = 0.000036$ ). Thus it may be concluded that SA:V displays significantly stronger phylogenetic signal  
249 than does relative body width, and that shared evolutionary history has strongly influenced trait covariation  
250 among taxa for SA:V. Biologically, this observation corresponds with the fact that tropical species – which  
251 form a monophyletic group within plethodontids – display greater variation in SA:V which covaries with  
252 disparity in their climatic niches (Baken et al. 2020). Thus, because of this macroevolutionary association,  
253 strong phylogenetic signal in SA:V is to be expected.

## 254 Discussion

255 It is common in comparative evolutionary studies to characterize the phylogenetic signal in phenotypic traits  
256 to determine the extent to which shared evolutionary history has generated trait covariation among taxa.  
257 However, while numerous analytical approaches may be used to quantify phylogenetic signal (e.g., Abouheif  
258 1999; Gittleman and Kot 1990; Pagel 1999; Blomberg et al. 2003; Adams 2014a), methods that explicitly  
259 measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained  
260 underdeveloped. In this study, we evaluated the precision of one common measure, Pagel's  $\lambda$ , and explored its  
261 efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations,  
262 we found that the precision of  $\lambda$  increased with increasing sample sizes; a pattern noted previously (Boettiger  
263 et al. 2012), and one that conformed with parametric statistical theory (Cohen 1988). However, we also found  
264 that vastly different  $\lambda$  estimates could be obtained from data containing the same level of phylogenetic signal,  
265 and that similar  $\lambda$  estimates may be obtained from data containing differing levels of phylogenetic signal.  
266 Further, the precision of  $\lambda$  varied with the strength of phylogenetic signal, where lower precision was observed  
267 when in data whose phylogenetic signal was of intermediate strength. From these findings we conclude that  $\lambda$  is  
268 not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and that biologi-  
269 cal interpretations of the strength of signal based on this parameter may inaccurately characterize such effects.

270

271 As an alternative, we described a standardized effect size ( $Z$ ) for assessing the strength of phylogenetic signal.  
272  $Z$  expresses the magnitude of phylogenetic signal as a standard normal deviate, which is easily interpretable  
273 as the strength of phylogenetic signal relative to the mean. We applied this concept to both  $\lambda$  and  $Kappa$ ,  
274 and found that  $Z_K$  was a better estimate of the strength of phylogenetic signal in phenotypic data. First,  $Z_K$   
275 was more precise than  $Z_\lambda$ , and variation in this precision was more consistent across the range of input levels

276 of phylogenetic signal. Additionally, values of  $Z_K$  more accurately tracked known levels of phylogenetic signal,  
277 with changes in the actual strength of phylogenetic signal reflected in a more linear fashion by concomite  
278 changes in the values of  $Z_K$ . Thus,  $Z_K$  holds promise as a measure of the relative strength of phylogenetic  
279 signal that reflects the magnitude of this effect in phenotypic data. We therefore recommend that future  
280 studies interested in the strength of phylogenetic signal incorporate  $Z_K$  as a statistical measure of this effect.

281

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

298

299 One important advantage of the approach advocated here is that the resulting effect sizes ( $Z_K$ ) are  
300 dimensionless, as the units of measurement cancel out during the calculation of  $Z$  (Sokal and Rohlf 2012).  
301 Thus,  $Z_K$  represents the strength of phylogenetic signal on a common and comparable scale – measured  
302 in standard deviation units – regardless of the initial units and original scale of the phenotypic variables  
303 under investigation. This means that the strength of phylogenetic signal may be compared across datasets  
304 for continuous phenotypic traits measured in different units and scale, because those units have been  
305 standardized through their conversion to  $Z_K$ . For example, our approach could be utilized to determine  
306 whether the strength of phylogenetic signal (say, in response to ecological differentiation) is stronger in  
307 morphological traits (linear traits:  $mm$ ), physiological traits (metabolic rate:  $\frac{O^2}{min}$ ), or behavioral traits

308 (aggression:  $\frac{\#displays}{second}$ ). Additionally, our method is capable of comparing the strength of phylogenetic signal  
309 in traits of different dimensionality, as estimates of phylogenetic signal using *Kappa* have been generalized for  
310 multivariate data ( $K_{mult}$ : see Adams 2014a). Likewise, as  $Z_K$  is a standard deviation unit (which accounts  
311 for variation, sample size, and the phylogeny), tests based on  $\hat{Z}_{12}$  may be utilized for comparing the strength  
312 of phylogenetic signal across datasets containing a different number of species, and even for phenotypes  
313 obtained from species in different lineages, because their phylogenetic non-independence, and observed  
314 variation, are taken into account in the generation of the empirical sampling distribution via permutation.

315

316 \*\* Finally... Need Closing paragraph. \*\*

317 **References**

- 318 Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.  
319 Evolutionary Ecology Research 1:895–909.
- 320 Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other  
321 high-dimensional multivariate data. Systematic Biology 63:685–697.
- 322 Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other  
323 high-dimensional multivariate data. Evolution 68:2675–2688.
- 324 Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative  
325 modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- 326 Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration  
327 across morphometric datasets. Evolution 70:2623–2631.
- 328 Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,  
329 and a refined permutation procedure. Evolution 72:1204–1215.
- 330 Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate  
331 phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- 332 Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric  
333 analyses. R package version 3.2.1.
- 334 Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of  
335 geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- 336 Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.  
337 Trends in Ecology and Evolution 10:236–240.
- 338 Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with  
339 phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil  
340 434:305–326.
- 341 Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution  
342 9:7005–7016.

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

- 370 Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* 125:1–15.
- 371 Freckleton, R. P., P. H. Harvey, and M. Pagel. 2002. Phylogenetic analysis and comparative data: A test and  
372 review of evidence. *American Naturalist* 160:712–726.
- 373 Garland, T. J., and A. R. Ives. 2000. Using the past to predict the present: Confidence intervals for regression  
374 equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
- 375 Gittleman, J. L., and M. Kot. 1990. Adaptation: Statistics and a null model for estimating phylogenetic  
376 effects. *Systematic Zoology* 39:227–241.
- 377 Glass, G. V. 1976. Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
- 378 Grafen, A. 1989. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London B,*  
379 *Biological Sciences* 326:119–157.
- 380 Harmon, L. J., J. T. Weir, C. D. Brock, R. E. Glor, and W. Challenger. 2008. GEIGER: Investigating  
381 evolutionary radiations. *Bioinformatics* 24:129–131.
- 382 Harvey, P. H., and M. D. Pagel. 1991. The comparative method in evolutionary biology. Oxford University  
383 Press, Oxford.
- 384 Hedges, L. V., and I. Olkin. 1985. Statistical methods for meta-analysis. Elsevier.
- 385 Kelley, K., and K. J. Preacher. 2012. On effect size. *Psychological Methods* 17:137–152.
- 386 Klingenberg, C. P., and N. A. Gidaszewski. 2010. Testing and quantifying phylogenetic signals and homoplasy  
387 in morphometric data. *Systematic biology* 59:245–261.
- 388 Liu, H., C. P. Osborne, D. Yin, R. P. Freckleton, G. Jiang, and M. Liu. 2019. Phylogeny and ecological  
389 processes influence grass coexistence at different spatial scales within the steppe biome. *Oecologia*  
390 191:25–38.
- 391 Martins, E. P., and T. F. Hansen. 1997. Phylogenies and the comparative method: A general approach  
392 to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist*  
393 149:646–667.
- 394 Molina-Venegas, R., and M. A. Rodriguez. 2017. Revisiting phylogenetic signal; strong or negligible impacts  
395 of polytomies and branch length information? *BMC evolutionary biology* 17:53.
- 396 Munkemuller, T., S. Lavergne, B. Bzeznik, S. Dray, T. Jombart, K. Schiffers, and W. Thuiller. 2012. How to

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

- <sup>424</sup> Vandelook, F., S. Janssens, P. Gijbels, E. Fischer, W. Van den Ende, O. Honnay, and S. Abrahamczyk. 2019.  
<sup>425</sup> Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124:269–279.

426      **Figure Legends**

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

429

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

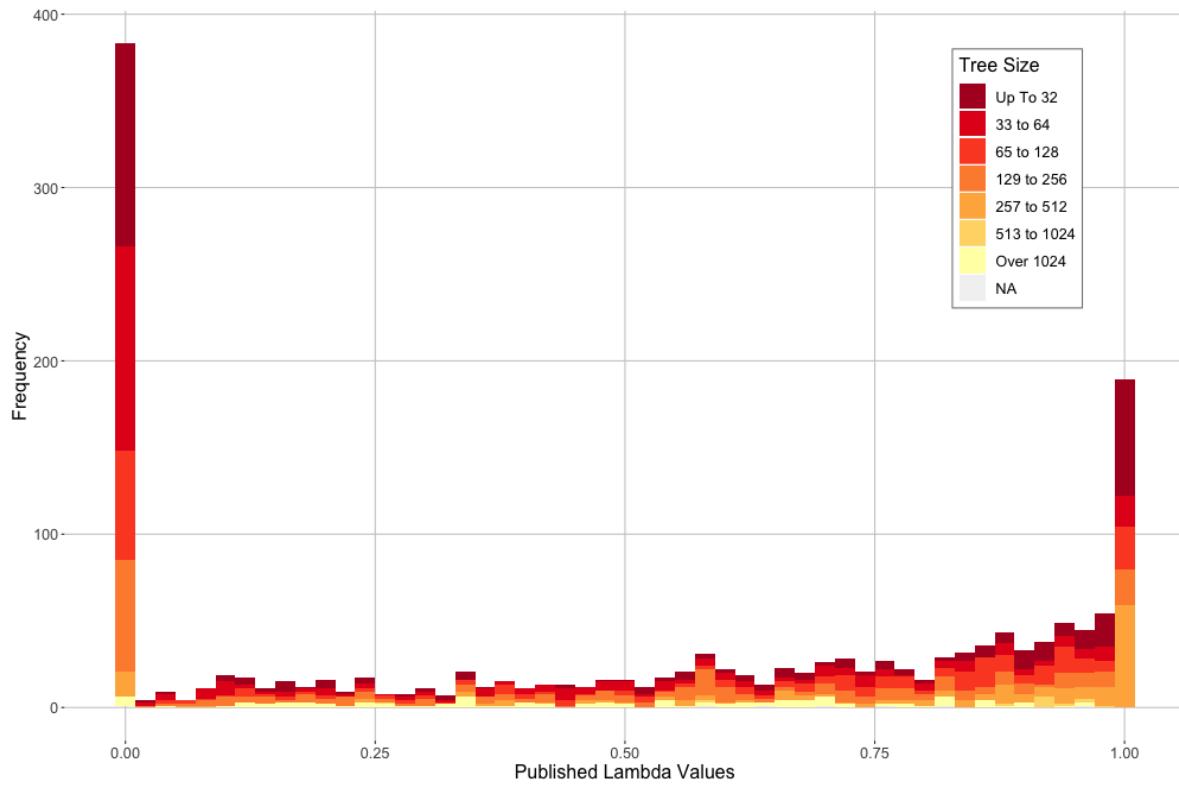
434

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

439

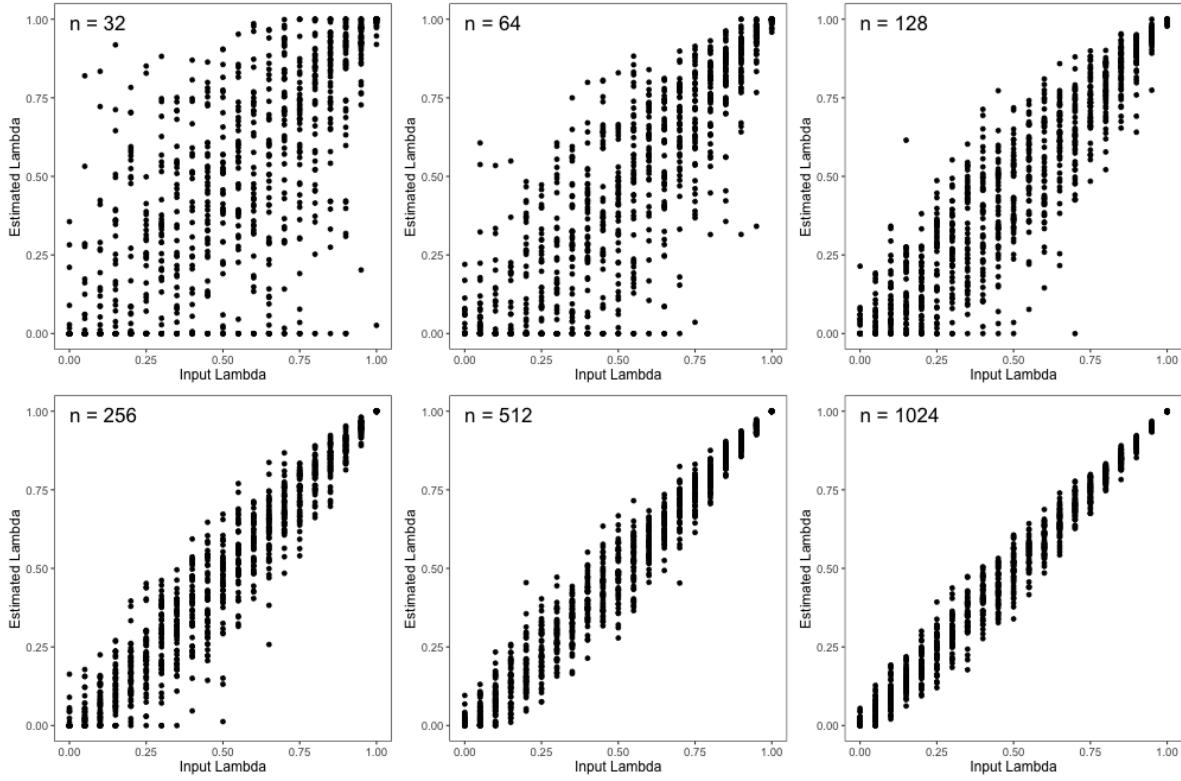
440      **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.  
441      (A) Estimates  $Z_\lambda$  for data simulated on phylogenies with 32 taxa ( $n = 32$ ), (B) Estimates of  $Z_K$  for data  
442      simulated on phylogenies with 32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

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



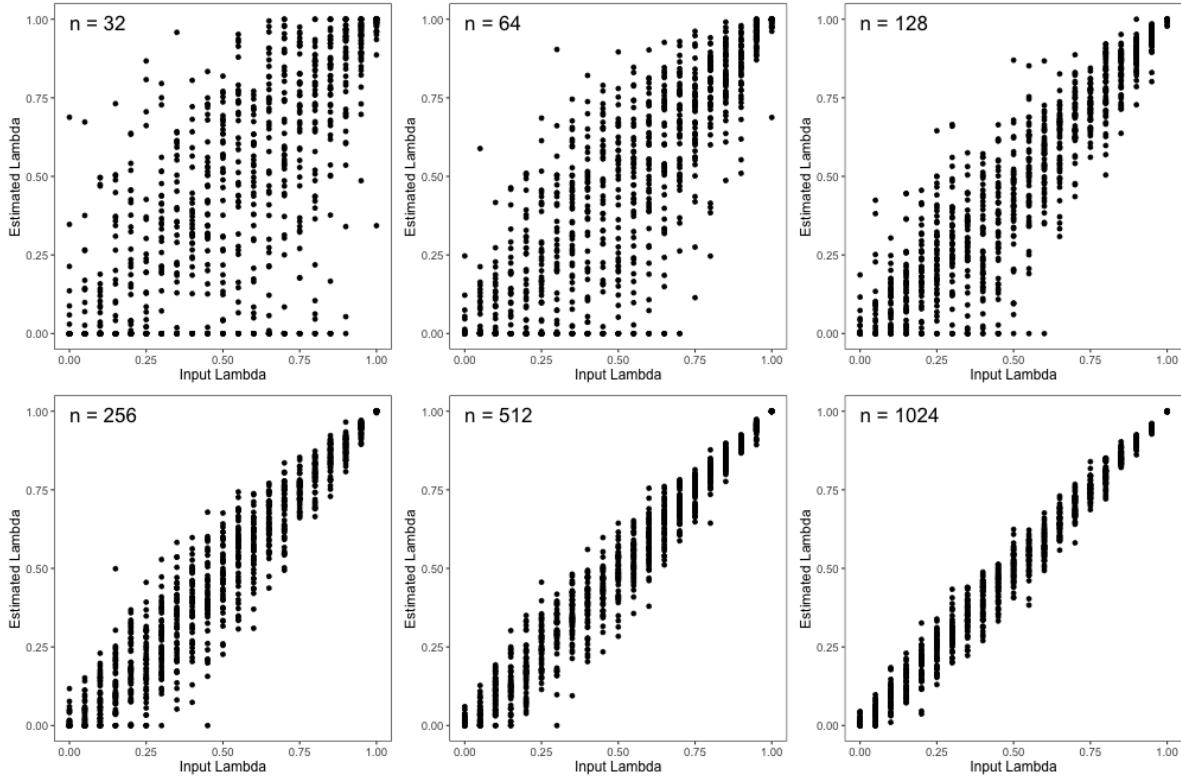
449

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



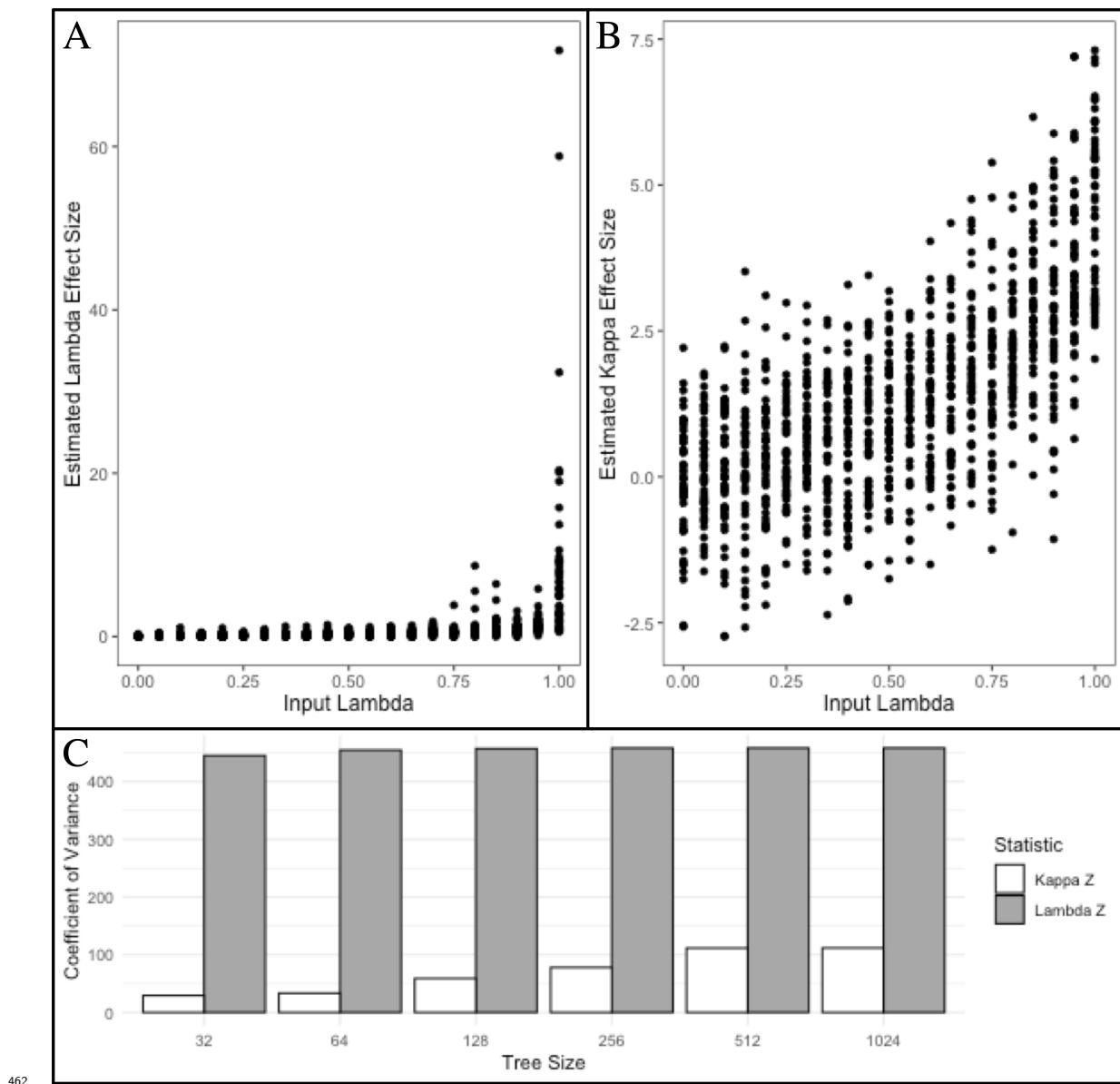
452

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

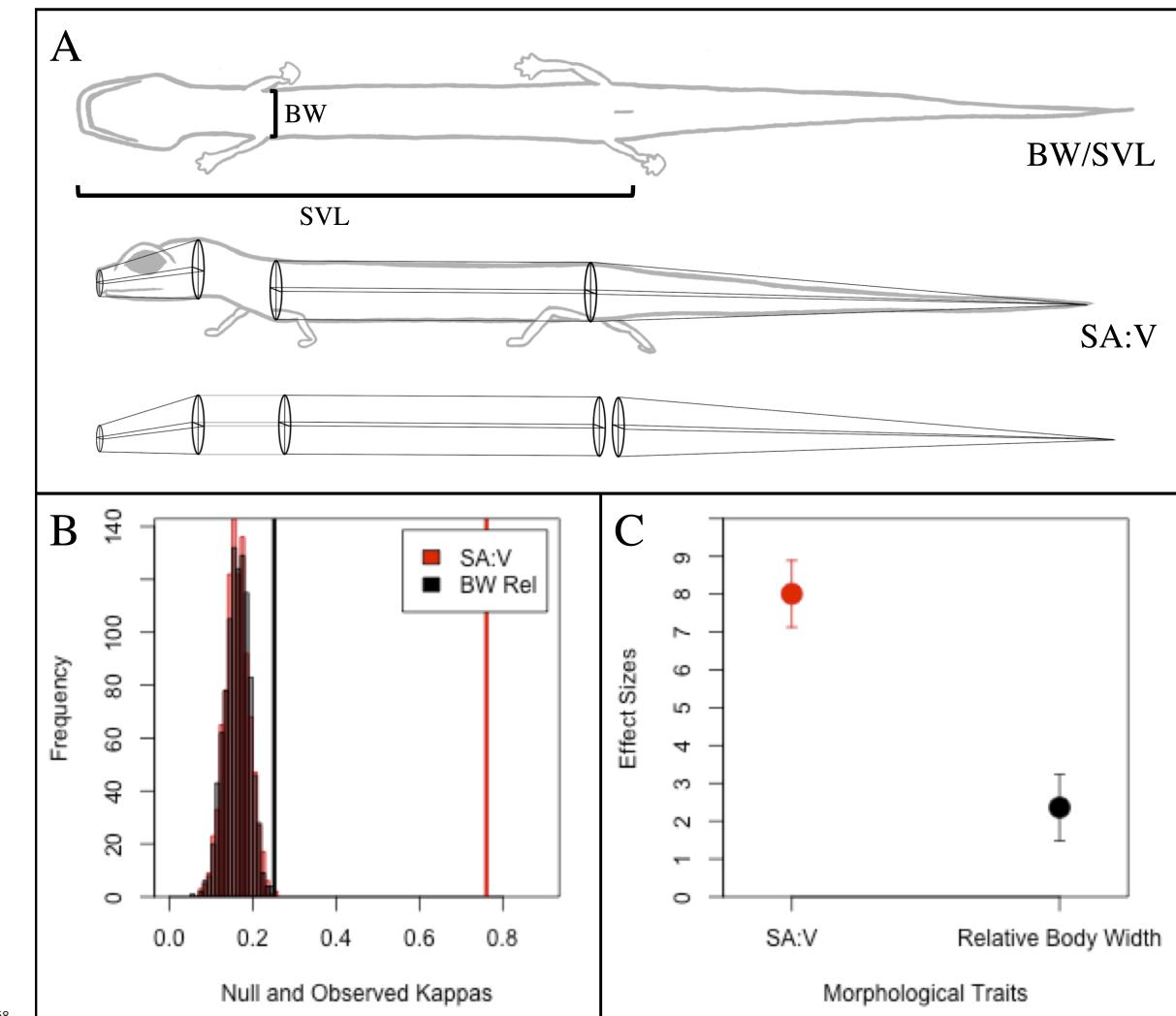


457

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



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



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