

<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 PGLS analysis,  $\lambda$  is a parameter that is tuned via log-likelihood profiling for its estimation. After obtaining  
51 an optimized  $\lambda$  value,  $\hat{\lambda}$ , likelihood ratio tests can be performed to compare observed model fits using  $\hat{\lambda}$  to  
52 those obtained when  $\lambda = 0$  or  $\lambda = 1$  to infer whether phylogenetic signal differs from no signal or a Brownian  
53 motion (BM) model of evolutionary divergence, respectively (Freckleton et al. 2002; Cooper et al. 2010; Bose  
54 et al. 2019). Confidence limits on  $\hat{\lambda}$  can be also estimated based on percentiles of a  $\chi^2$  distribution and  
55 where these values are mapped onto a log-likelihood profile (*sensu* Freckleton et al. 2002). With respect  
56 to confidence limits, one can ascertain whether  $\lambda = 0$  or  $\lambda = 1$  is contained within an interval of  $\hat{\lambda}$ , as an  
57 analog to likelihood ratio tests (see, e.g., Vandeloek et al. 2019). It is, therefore, tempting to regard  $\hat{\lambda}$  as a  
58 descriptive statistic that measures the relative strength of phylogenetic signal, for inferring the extent to  
59 which shared evolutionary history has influenced trait covariation among taxa, is common. The appeal of  $\hat{\lambda}$   
60 as a descriptive statistic is that it serves as the basis for interpreting “weak” versus “strong” phylogenetic  
61 signal; i.e., small versus large values of  $\hat{\lambda}$ , respectively, in a comparative sense (e.g., De Meester et al. 2019;  
62 Pintanel et al. 2019; Su et al. 2019). For example, qualitative comparisons of  $\hat{\lambda}$  for multiple traits on the  
63 same phylogenetic tree provide a basis for inferring whether the strength of phylogenetic signal is greater in  
64 one trait as compared to another (e.g., Liu et al. 2019; Bai et al. 2019). Indeed, statements regarding the  
65 strength of phylogenetic signal based on  $\hat{\lambda}$  are rather common in the evolutionary literature. For instance, of  
66 the 204 papers published in 2019 that estimated and reported Pagel’s  $\lambda$  (found from a literature survey we  
67 conducted in Google.scholar), 40% explicitly interpreted the strength of phylogenetic signal for at least one  
68 phenotypic trait. Further, because nearly half of the 1,572  $\hat{\lambda}$  values reported were near 0 or 1 (Figure 1)  
69 where the biological interpretation of  $\lambda$  is known ( $\lambda = 0$  represents no phylogenetic signal, while  $\lambda = 1$  is  
70 phylogenetic signal as expected under BM), this percentage is even higher.

71

72 [insert Figure 1 here]

73

74 It seems intuitive to interpret the strength of phylogenetic signal based on the value of  $\hat{\lambda}$ , as  $\lambda$  is a parameter  
75 on a bounded scale ( $0 \rightarrow 1$ ) for which interpretation of its extremal points are understood. However, equating  
76 values of  $\hat{\lambda}$  directly to the strength of phylogenetic signal presumes two important statistical properties that  
77 have not been fully explored. First, it presumes that values of  $\hat{\lambda}$  can be precisely estimated, as biological  
78 inferences regarding the strength of phylogenetic signal depend on high accuracy in its estimation. (In order  
79 to be precise, the log-likelihood profile would have to be peaked rather than flat.) Therefore, understanding  
80 the precision in estimating  $\hat{\lambda}$  is paramount. One study (Boettiger et al. 2012) found that estimates of Pagel’s  
81  $\lambda$  displayed less variation (i.e., greater precision) when data were simulated on a large phylogeny ( $N = 281$ )

82 as compared to a small one ( $N = 13$ ). From this observation it was concluded that insufficient data (i.e.,  
83 low number of species) was the underlying cause of the increased variation across parameter estimates  
84 (Boettiger et al. 2012). Indeed, such a pattern is common with statistical estimators, as summary statistics  
85 and parameters are often more precise at greater sample sizes (Cohen 1988). However, this conclusion  
86 could imply that the shape of log-likelihood profiles are assumed consistent with species number. Such an  
87 assumption and the importance of tree shape have to date not been verified. Thus, despite widespread use of  
88 Pagel's (1999)  $\lambda$  in macroevolutionary studies, at present, we lack a general understanding of the precision  
89 with which  $\hat{\lambda}$  can estimate levels of phylogenetic signal in phenotypic datasets.

90

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

103

104 As a proportional random variable bounded by 0 and 1, we might expect that  $\hat{\lambda}$  follows A Bernoulli distribution;  
105 i.e., branch lengths in a tree are scaled proportionally to the probability that data arise from a BM process.  
106 As such, we would also expect in simulation experiment that given a known  $\lambda$  value used to generate random  
107 data on a tree: the mean of an empirical sampling distribution of  $\hat{\lambda}$  would approximately equal  $\lambda$ ; the standard  
108 error of  $\hat{\lambda}$  would be largest at intermediate values of  $\lambda$ , and predictable over the range of  $\lambda$  with respect  
109 to treesize; the distribution of  $\hat{\lambda}$  would be symmetric at intermediate values of  $\lambda$  and more skewed toward  
110 values of 0 or 1; and that the distribution of  $\hat{\lambda}$  will be more platykurtic at intermediate values of  $\lambda$ , becoming  
111 more leptokurtic toward 0 and 1. These properties seemed to be superficially reasonable from the simulations  
112 performed by Münkemüller et al.(2012, see Fig. 2). Their study simulated strength of phylogenetic signal  
113 as a varied weighted-average of data simulated on trees with  $\lambda = 0$  and  $\lambda = 1$ , rather than from prescribed

114 values of  $\lambda$ . Nonetheless, their results revealed that as a random variable,  $\hat{\lambda}$  resembled a Bernoulli random  
115 variable based superficially on statistical moments (mean, variance, skewness, and kurtosis; see their Fig.  
116 2) with respect to strength of phylogenetic signal, for a given tree size. However, for comparisons among  
117 tree sizes for a given strength of phylogenetic signal, median values of  $\hat{\lambda}$  were slightly, positively associated  
118 with the number of species in a tree, especially at weaker strength of signals, suggesting  $\hat{\lambda}$  might not be an  
119 accurate statistic across a range of tree sizes. In order for  $\hat{\lambda}$  to be considered a reliable statistic as a measure  
120 of phylogenetic signal, as a minimum, it should be accurate and precise, properties that could be measured  
121 from simulation experiments for known values of  $\lambda$ .

122 In this study, we evaluate the precision of  $\hat{\lambda}$  as an optimized parameter, used as a statistic, for estimating  
123 known levels of phylogenetic signal in phenotypic data (data simulated with known  $\lambda$ ). We use computer  
124 simulations with differing numbers of species, differently shaped phylogenies, and differing input levels of  $\lambda$ ,  
125 to explore the degree to which  $\hat{\lambda}$  correctly identifies known levels of phylogenetic signal,  $\lambda$ , and under which  
126 circumstances. We find that estimates of  $\hat{\lambda}$  vary widely for a given input value of  $\lambda$ , and that the precision of  
127  $\hat{\lambda}$  is not constant across its range. Rather, there is decreased precision when input levels of phylogenetic signal  
128 are of intermediate strength. Additionally, the same values of  $\hat{\lambda}$  may be obtained from datasets containing  
129 vastly different input levels of  $\lambda$ . Thus,  $\hat{\lambda}$  is not a reliable indicator of the strength of phylogenetic signal  
130 in phenotypic data. We then describe a standardized effect size for measuring the strength of phylogenetic  
131 signal in phenotypic datasets and apply the concept to two common measures of phylogenetic signal:  $\hat{\lambda}$  and  
132  $\kappa$ . Through simulations we find that the precision of effect sizes based on  $\hat{\lambda}$  ( $Z_{\lambda}$ ) are less reliable than that  
133 those based on  $\kappa$  ( $Z_{\kappa}$ ), implying that  $Z_{\kappa}$  is a more robust effect size measure. We also propose a two-sample  
134 test statistic that may be used to compare the strength of phylogenetic signal among datasets, and provide  
135 an empirical example to demonstrate its use. We conclude that estimates of phylogenetic signal using Pagel's  
136  $\lambda$  (as a descriptive statistic, not as a PGLS parameter) are often inaccurate, and thus interpreting strength  
137 of phylogenetic signal in phenotypic datasets based on this measure is compromised. By contrast, effect  
138 sizes obtained from  $\kappa$  hold promise for characterizing phylogenetic signal, and for comparing the strength of  
139 phylogenetic signal across datasets.

<sup>140</sup> **Methods and Results**

<sup>141</sup> ***The Precision of  $\hat{\lambda}$  is Variable***

<sup>142</sup> We conducted a series of computer simulations to evaluate the precision of  $\hat{\lambda}$  based on varied input of  
<sup>143</sup> Pagel's  $\lambda$  to generate data. Our primary simulations were based on pure-birth phylogenies; however, we  
<sup>144</sup> also evaluated patterns on both balanced and pectinate trees to determine whether tree shape affected our  
<sup>145</sup> findings (see Supporting Information). First we generated 50 pure-birth phylogenies at each of six different  
<sup>146</sup> tree sizes, ranging from 32 to 1024 taxa ( $n = 2^5 - 2^{10}$ ). For each simulation run, we rescaled the simulated  
<sup>147</sup> phylogenies by multiplying the internal branches by  $\lambda_{in}$ , using 21 intervals of 0.05 units across its range  
<sup>148</sup> ( $\lambda_{in} = 0.0 \rightarrow 1.0$ ), resulting in 1050 scaled phylogenies at each level of species richness ( $n$ ). Continuous traits  
<sup>149</sup> were then simulated on each phylogeny under a BM model of evolution to obtain datasets with differing  
<sup>150</sup> levels of phylogenetic signal, that ranged from no phylogenetic signal (when  $\lambda_{in} = 0$ ), to phylogenetic signal  
<sup>151</sup> reflecting a BM model of evolution (when  $\lambda_{in} = 1$ ). For each dataset we then estimated phylogenetic  
<sup>152</sup> signal ( $\hat{\lambda}$ ), and calculated its variance ( $\sigma_{\hat{\lambda}}^2$ ) at each input level of phylogenetic signal and level of species  
<sup>153</sup> richness as an estimate of precision. Random data samples were generated from a normal distribution, i.e.,  
<sup>154</sup>  $\mathcal{N}(\mu = 0, \sigma = \sigma^2 \mathbf{C})$ . We verified that the level of  $\sigma^2$  had no effect on  $\hat{\lambda}$  or  $\sigma_{\hat{\lambda}}^2$ , but the results we present used  
<sup>155</sup>  $\sigma^2 = 1$ .

<sup>156</sup>

<sup>157</sup> We also evaluated the precision of  $\hat{\lambda}$  when estimated in PGLS regression and ANOVA (i.e.,  $Y \sim X$ ). For  
<sup>158</sup> these simulations,  $Y = \beta X + \epsilon$  was used to generate data with varied phylogenetic signal and variable  
<sup>159</sup> associations (*sensu* Revell 2010). In the case of PGLS regression,  $X$  was a variable simulated under a BM  
<sup>160</sup> model of evolution, constrained by rescaling phylogeny according to  $\lambda$ , and in the case of ANOVA,  $X$  was a  
<sup>161</sup> balanced binary variable (values of 0 and 1), randomly ordered. The strength of variable associations ranged  
<sup>162</sup> as  $\beta = 0.0, 0.25, 0.5, 0.75$  and 1.0. Residuals,  $\epsilon$ , were simulated under a BM model of evolution, constrained  
<sup>163</sup> by rescaling phylogeny according to  $\lambda$ . Simulated independent variables or residuals were generated from a  
<sup>164</sup> normal distribution, i.e.,  $\mathcal{N}(\mu = 0, \sigma = \sigma^2 \mathbf{C})$ , with the phylogenetic covariance matrix,  $\mathbf{C}$ , implicitly varied  
<sup>165</sup> for each tree by the level of  $\lambda$ . The fit of the phylogenetic regression was estimated using maximum likelihood,  
<sup>166</sup> and parameter estimates ( $\hat{\lambda}$  and  $\hat{\beta}$ ) were obtained. We then calculated precision estimates ( $\sigma_{\hat{\lambda}}^2$ ) at each input  
<sup>167</sup> level of phylogenetic signal and level of species richness. We verified that the level of  $\sigma^2$  had no effect on  $\sigma_{\hat{\lambda}}^2$   
<sup>168</sup> but did influence the precision of  $\hat{\beta}$  estimated from the linear model (precision increased with smaller  $\epsilon$ , as  
<sup>169</sup> expected). The results we present used  $\sigma^2 = 1$ .

<sup>170</sup>

171 All analyses were performed in R v3.6.0 (R Core Team 2019) using the packages **geiger** (Harmon et al.  
172 2008), **caper** (Orme et al. 2013), **phytools** (Revell 2012), and **geomorph** 3.2.1 (Adams and Otárola-Castillo  
173 2013; Adams et al. 2020). R-scripts are found in the Supporting Information.

174

175 *Results.* We found that the precision of  $\lambda_{est}$  varied widely across simulation conditions. Predictably, precision  
176 improved as the number of species increased (Figure 2). This confirmed earlier findings of Boettiger et al.  
177 (2012), and adhered to parametric statistical theory. However, in many cases the set of  $\lambda_{est}$  spanned nearly  
178 the entire range of possible values (e.g.,  $n = 32$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.0 \rightarrow 0.985$ ), revealing that estimates  
179 of  $\lambda$  were not a reliable indicator of input phylogenetic signal. Importantly, the precision of  $\lambda_{est}$  was not  
180 uniform across all levels of phylogenetic signal. The worst precision was observed at intermediate levels  
181 of phylogenetic signal ( $\lambda_{in} \approx 0.5$ ), while precision improved as input levels approached the extremes of  
182  $\lambda$ '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  
183 intermediate values. Additionally, even at large levels of species richness, we found that the range of  $\lambda_{est}$  still  
184 encompassed a substantial portion of possible values (e.g.,  $n = 512$ ;  $\lambda_{in} = 0.5$ :  $\lambda_{est} = 0.32 \rightarrow 0.68$ ). Likewise,  
185 the same  $\lambda_{est}$  could be obtained from datasets containing vastly different input levels of phylogenetic  
186 signal (e.g.,  $n = 512$ ;  $\lambda_{est} = 0.5$ ;  $\lambda_{in} = 0.25 \rightarrow 0.65$ ). These findings were particularly unsettling when  
187 considered in light of our literature survey. Over one quarter of the  $\lambda$  estimates published in empirical  
188 studies (421 of 1,572) were between  $\lambda = 0.25$  and  $\lambda = 0.75$  (Figure 1). This range reflected the region  
189 that our simulations identified as being the least reliable in terms of accurately characterizing levels of  
190 phylogenetic signal, yet 30% of these mid-range empirical estimates were explicitly interpreted in terms of  
191 the strength of phylogenetic signal that they represented (i.e., weak, intermediate, strong phylogenetic signal).

192

193 When phylogenetic signal was investigated on balanced or pectinate trees, patterns in the precision of  $\lambda$  were  
194 largely the same, with decreased precision at intermediate levels of phylogenetic signal (Supporting Informa-  
195 tion). Likewise, when  $\lambda$  was co-estimated with regression parameters in PGLS regression and ANOVA, the  
196 results of our simulations were quite similar. Regression parameters ( $\beta$ ) were accurately estimated, confirming  
197 earlier findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal ( $\lambda$ )  
198 were less precise (Figure 3; see also Supporting Information), and the spread of  $\lambda_{est}$  was similar to that  
199 observed when  $\lambda$  was estimated for only the dependent variable, as in Figure 2. Taken together, these findings  
200 reveal that  $\lambda_{est}$  does not precisely characterize known levels of phylogenetic signal in phenotypic datasets,  
201 and that biological interpretations of the strength of phylogenetic signal based on  $\lambda$  may be highly inaccurate.

202

203 [insert Figure 2 here]

204

205 [insert Figure 3 here]

206

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

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

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

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

221

222 Adams and Collyer (2019b) suggested that the strength of phylogenetic signal could be represented as an effect  
223 size based on the  $\kappa$  statistic and its empirical sampling distribution from permutation. Here we formalize  
224 that suggestion, resulting in an effect size of:

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

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

229

230 An effect size based on  $\lambda$  is a little more challenging because, as a statistic, it does not have a sampling  
 231 distribution on which to measure standard error. Confidence intervals for  $\lambda$  are not generated from standard  
 232 error calculations, but rather for the values of lambda that intersect the log-likelihood profile for corresponding  
 233 percentiles of the  $\chi^2$  distribution used to compare the putative appropriate model to a null  
 234 model with  $\lambda = 0$  [citation needed, think Boettiger has a paper on this]. The simplest similar effect size is

$$|Z_\lambda| = \sqrt{\chi_{\hat{\lambda}}^2} \quad (3)$$

235 where,  $\hat{\lambda}$  is the maximized likelihood value of  $\lambda$  and  $\chi_{\hat{\lambda}}^2$  is the likelihood ratio statistic for the value.  
 236 Hypothetically, this equation would be better represented as  $Z_\lambda = d\sqrt{\chi_{\hat{\lambda}}^2}$ , where  $d$  is a binary value  $(-1, 1)$   
 237 to indicate a direction based on whether  $\hat{\lambda}$  is below or above a critical value of  $\lambda$ , for a quantile from a  $\chi^2$   
 238 distribution at a probability of 0.5. The critical value would correspond to a value of  $Z = 0$ ; i.e., the strength  
 239 of phylogenetic signal that might be observed at the 50th percentile in a distribution of random outcomes,  
 240 generated from a model where  $\lambda = 0$ . However, in some preliminary simulations, we found that mapping  $\chi_{\hat{\lambda}}^2$   
 241 values this way did not produce effect sizes that were symmetrical about  $Z = 0$ , as the mapping is not linear  
 242 and, especially for small trees, the log-likelihood profiles can be rather flat, making estimation of  $\hat{\lambda}$  rather  
 243 imprecise (as we elaborate more below.) We will henceforth use equation 3, recognizing that large values of  $Z$   
 244 would likely be positive, anyway, based on the reference model for the likelihood ratio test.

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

248

249 To evaluate the utility of  $Z_\kappa$  and  $Z_\lambda$  we calculated both effect sizes for the simulated datasets generated  
 250 above, and summarized the precision of each using its variance ( $\sigma_{Z_\kappa}^2$  and  $\sigma_{Z_\lambda}^2$ , Figure 4: additional results in  
 251 the Supporting Information). Here two things are evident. First, estimates of  $Z_\kappa$  linearly track the input

252 phylogenetic signal whereas estimates of  $Z_\lambda$  do not (Figure 4A, B). Thus, actual changes in the strength  
253 of phylogenetic signal are reflected more evenly in the corresponding values of the effect size  $Z_\kappa$ . Second,  
254 the precision of  $Z_\kappa$  is considerably more stable as compared with  $Z_\lambda$ . This may be seen by calculating the  
255 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  
256 signal. Coefficients of variation in the precision of  $Z_\kappa$  were up to an order of magnitude smaller than for  $Z_\lambda$   
257 (Figure 4C), implying that estimates of the strength of phylogenetic signal were more reliable and robust  
258 when using  $Z_\kappa$ .

259

260 [insert Figure 4 here]

## 261 *Statistical Comparisons of Phylogenetic Signal*

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

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

266 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  
267 illustrates that if  $Z_\kappa$  has already been calculated for two sampling distributions as in equation 2, the sampling  
268 distributions have unit variance for each of the  $Z_\kappa$  statistics. Estimates of significance of  $\hat{Z}_{12}$  may be obtained  
269 from a standard normal distribution. Typically,  $\hat{Z}_{12}$  is considered a two-tailed test, however directional  
270 (one-tailed) tests may be specified should the empirical situation require it (see Adams and Collyer 2016,  
271 2019a).

272

## 273 *Empirical Example*

274 To demonstrate the utility of  $\hat{Z}_{12}$  we quantified and compared the strength of phylogenetic signal of two  
275 phenotypic traits across species of plethodontid salamander. The data were part of a series of studies

276 examining macroevolutionary trends in phenotypic diversification in this group (Baken and Adams 2019;  
277 Baken et al. 2020). Our dataset contained surface area to volume ratios (SA:V) and relative body width  
278 ( $\frac{BW}{SVL}$ ) for 305 species (Figure 5A). For SA:V, 11 linear body measurements from 2,781 individuals were  
279 taken, from which estimates of the surface area and volume of the head, body, and tail were calculated and  
280 subsequently combined to arrive at the SA:V for each individual (for mathematical details see Baken et al.  
281 2020). Species means were then obtained. Likewise, body size (SVL) and body width (BW) measurements  
282 were taken from 3,371 individuals, and species means of relative body width ( $\frac{BW}{SVL}$ ) were calculated (data  
283 from Baken and Adams 2019). A time-dated molecular phylogeny for the group (Bonett and Blair 2017)  
284 was then pruned to match the species in the phenotypic dataset. The phylogenetic signal in each trait was  
285 then characterized using  $\kappa$ , which was converted to its effect size ( $Z_\kappa$ ) using geomorph 3.2.1 (Adams and  
286 Otárola-Castillo 2013; Adams et al. 2020). Finally, the strength of phylogenetic signal was compared across  
287 traits using  $\hat{Z}_{12}$  as described above (to be incorporated in geomorph upon manuscript acceptance).

288

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

## 301 Discussion

302 It is common in comparative evolutionary studies to characterize the phylogenetic signal in phenotypic traits  
303 to determine the extent to which shared evolutionary history has generated trait covariation among taxa.  
304 However, while numerous analytical approaches may be used to quantify phylogenetic signal (e.g., Abouheif  
305 1999; Gittleman and Kot 1990; Pagel 1999; Blomberg et al. 2003; Adams 2014a), methods that explicitly

306 measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained  
307 underdeveloped. In this study, we evaluated the precision of one common measure, Pagel's  $\lambda$ , and explored its  
308 efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations,  
309 we found that the precision of  $\lambda$  increased with increasing sample sizes; a pattern noted previously (Boettiger  
310 et al. 2012), and one that conformed with parametric statistical theory (Cohen 1988). However, we also found  
311 that vastly different  $\lambda$  estimates could be obtained from data containing the same level of phylogenetic signal,  
312 and that similar  $\lambda$  estimates may be obtained from data containing differing levels of phylogenetic signal.  
313 Further, the precision of  $\lambda$  varied with the strength of phylogenetic signal, where lower precision was observed  
314 when in data whose phylogenetic signal was of intermediate strength. From these findings we conclude that  $\lambda$  is  
315 not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and that biological  
316 interpretations of the strength of signal based on this parameter may inaccurately characterize such effects.

317

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

328

329 Based on the effect size  $Z_\kappa$ , we then proposed a two-sample test, which provides means of determining  
330 whether the strength of phylogenetic signal is greater in one phenotypic trait as compared to another, via a  
331 hypothesis test. Prior studies have summarized patterns of variation in phylogenetic signal across datasets  
332 using summary test values, such as  $\kappa$  (e.g., Blomberg et al. 2003). However,  $\kappa$  does not scale linearly with  
333 input levels of phylogenetic signal, and its variance increases (i.e., precision decreases) with increasing  
334 strength of phylogenetic signal (Münkemüller et al. 2012; Diniz-Filho et al. 2012; see also Supporting  
335 Information). Thus,  $\kappa$  should not be considered an effect size that measures the strength of phylogenetic  
336 signal on a common scale. By contrast, standardizing  $\kappa$  ( $Z_\kappa$ , via equation 2) alleviates these concerns, and  
337 facilitates formal statistical comparisons of the strength of signal across datasets. Thus when viewed from this

338 perspective, the approach developed here aligns well with other statistical approaches such as meta-analysis  
339 (sensu Hedges and Olkin 1985; Glass 1976; Arnqvist and Wooster 1995), where summary statistics across  
340 datasets are converted to standardized effect sizes for subsequent “higher order” statistical summaries or  
341 comparisons. As such, our approach enables evolutionary biologists to quantitatively examine the relative  
342 strength of phylogenetic signal across a wide range of phenotypic traits, and thus opens the door for future dis-  
343 coveries that inform on how phenotypic diversity accumulates in macroevolutionary time across the tree of life.

344

345 One important advantage of the approach advocated here is that the resulting effect sizes ( $Z_\kappa$ ) are  
346 dimensionless, as the units of measurement cancel out during the calculation of  $Z$  (Sokal and Rohlf 2012).  
347 Thus,  $Z_\kappa$  represents the strength of phylogenetic signal on a common and comparable scale – measured  
348 in standard deviations – regardless of the initial units and original scale of the phenotypic variables under  
349 investigation. This means that the strength of phylogenetic signal may be compared across datasets for  
350 continuous phenotypic traits measured in different units and scale, because those units have been standardized  
351 through their conversion to  $Z_\kappa$ . For example, our approach could be utilized to determine whether the  
352 strength of phylogenetic signal (say, in response to ecological differentiation) is stronger in morphological  
353 traits (linear traits:  $mm$ ), physiological traits (metabolic rate:  $\frac{O^2}{min}$ ), or behavioral traits (aggression:  
354  $\frac{\#displays}{second}$ ). In fact, our empirical example provided such a comparison, as SA:V is represented in  $mm^{-1}$   
355 while relative body size is a unitless ratio ( $\frac{BW}{SVL}$ ). Additionally, our method is capable of comparing the  
356 strength of phylogenetic signal in traits of different dimensionality, as estimates of phylogenetic signal using  $\kappa$   
357 have been generalized for multivariate data ( $\kappa_{mult}$ : see Adams 2014a). Furthermore, tests based on  $\hat{Z}_{12}$  may  
358 be utilized for comparing the strength of phylogenetic signal among datasets containing a different number  
359 of species, and even for phenotypes obtained from species in different lineages, because their phylogenetic  
360 non-independence and observed variation are taken into account in the generation of the empirical sampling  
361 distribution via permutation.

362

363 This study is not the first to compare  $\lambda$  and  $\kappa$  for their ability as statistics to measure phylogenetic signal.  
364 Our results for  $\lambda$  and  $\kappa$  values are consistent with those found in the simulations performed by Münkemüller  
365 et al. (2012), but that study investigated type I error rates and statistical power, finding that  $\lambda$  performed  
366 better in both regards, irrespective of species number in trees. Although not the central focus of their study,  
367 the same tendency for variable  $\lambda$  and consistent  $\kappa$  at intermediate phylogenetic signal strengths was observed  
368 (see Fig. 2, Münkemüller et al. 2012). Recent work by Molina-Venegas and Rodríguez (2017) found that  
369  $\kappa$  but not  $\lambda$  tended to inflate the estimate of phylogenetic signal, leading to moderate type I and type II

370 biases, if polytomic chronograms were used. Their work more thoroughly addressed previous observations of  
371 inflated  $\kappa$  for incompletely resolved phylogenetic trees (Davies et al. 2012; Münkemüller et al. 2012). An  
372 interesting question is whether an inflated  $\kappa$  value leads to an inflated  $Z_\kappa$  or does a tendency of a particular  
373 tree to inflate estimates of  $\kappa$  also inflate the values in random permutations of a test, in which case  $Z_\kappa$  is  
374 robust to polytomies? We repeated the analyses in Figure 4, adjusting trees to have 50% collapsed nodes, per  
375 the technique of Molina-Venegas and Rodríguez (2017), and found results were consistent (see Supporting  
376 Information). This confirms that any tendency of incompletely resolved trees to inflate  $\kappa$  as a descriptive  
377 statistic does not inflate  $Z_\kappa$  as an effect size. Furthermore, because comparison of effect sizes in a test is a  
378 comparison of locations of observed values in their sampling distributions, which would shift concomitantly  
379 because of this tendency, the  $Z_{12}$  test statistic in equation 4 appears to be robust in spite of unresolved trees.

380

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

398    **References**

- 399    Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.  
400         Evolutionary Ecology Research 1:895–909.
- 401    Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other  
402         high-dimensional multivariate data. Systematic Biology 63:685–697.
- 403    Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other  
404         high-dimensional multivariate data. Evolution 68:2675–2688.
- 405    Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative  
406         modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- 407    Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration  
408         across morphometric datasets. Evolution 70:2623–2631.
- 409    Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,  
410         and a refined permutation procedure. Evolution 72:1204–1215.
- 411    Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate  
412         phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- 413    Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric  
414         analyses. R package version 3.2.1.
- 415    Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of  
416         geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- 417    Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.  
418         Trends in Ecology and Evolution 10:236–240.
- 419    Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with  
420         phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil  
421         434:305–326.
- 422    Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution  
423         9:7005–7016.

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

- 451 Diniz-Filho, J. A. F., T. Santos, T. F. Rangel, and L. M. Bini. 2012. A comparison of metrics for estimating  
452 phylogenetic signal under alternative evolutionary models. *Genetics and Molecular Biology* 35:673–679.
- 453 Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* 125:1–15.
- 454 Freckleton, R. P., P. H. Harvey, and M. Pagel. 2002. Phylogenetic analysis and comparative data: A test and  
455 review of evidence. *American Naturalist* 160:712–726.
- 456 Garland, T. J., and A. R. Ives. 2000. Using the past to predict the present: Confidence intervals for regression  
457 equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
- 458 Gittleman, J. L., and M. Kot. 1990. Adaptation: Statistics and a null model for estimating phylogenetic  
459 effects. *Systematic Zoology* 39:227–241.
- 460 Glass, G. V. 1976. Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
- 461 Grafen, A. 1989. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London B,*  
462 *Biological Sciences* 326:119–157.
- 463 Harmon, L. J., J. T. Weir, C. D. Brock, R. E. Glor, and W. Challenger. 2008. GEIGER: Investigating  
464 evolutionary radiations. *Bioinformatics* 24:129–131.
- 465 Harvey, P. H., and M. D. Pagel. 1991. The comparative method in evolutionary biology. Oxford University  
466 Press, Oxford.
- 467 Hedges, L. V., and I. Olkin. 1985. Statistical methods for meta-analysis. Elsevier.
- 468 Kelley, K., and K. J. Preacher. 2012. On effect size. *Psychological Methods* 17:137–152.
- 469 Klingenberg, C. P., and N. A. Gidaszewski. 2010. Testing and quantifying phylogenetic signals and homoplasy  
470 in morphometric data. *Systematic biology* 59:245–261.
- 471 Liu, H., C. P. Osborne, D. Yin, R. P. Freckleton, G. Jiang, and M. Liu. 2019. Phylogeny and ecological  
472 processes influence grass coexistence at different spatial scales within the steppe biome. *Oecologia*  
473 191:25–38.
- 474 Martins, E. P., and T. F. Hansen. 1997. Phylogenies and the comparative method: A general approach  
475 to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist*  
476 149:646–667.
- 477 Molina-Venegas, R., and M. A. Rodríguez. 2017. Revisiting phylogenetic signal; strong or negligible impacts

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

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

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

509      **Figure Legends**

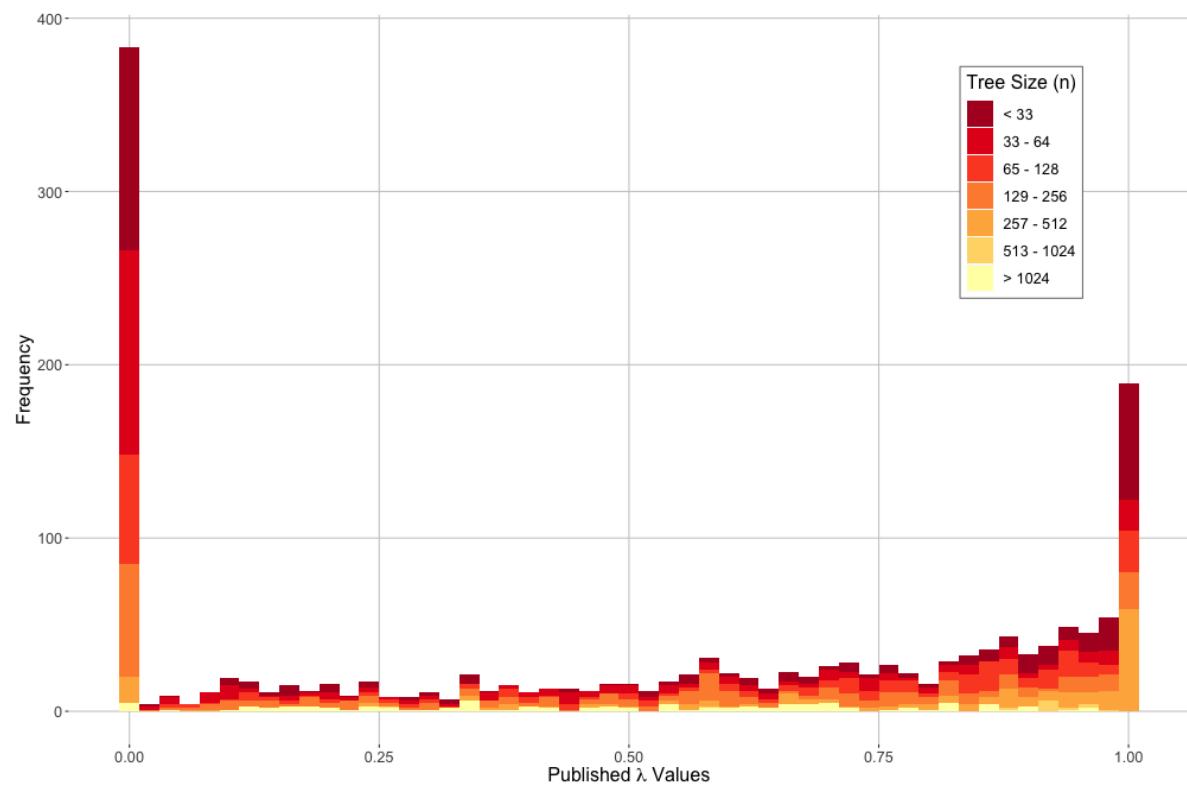
510      **Figure 1.** Frequency distribution of  $\lambda$  estimates published in 2019. The majority of these values were  
511      close to 0 or 1, and from phylogenies with fewer than 200 taxa.

512  
513      **Figure 2.** Precision of Pagel's  $\lambda$  across known levels of input phylogenetic signal ( $\lambda_{in}$ ) on phylogenies  
514      of various sizes. As phylogenies increase in size, variation in  $\lambda_{in}$  decreases; however the precision is  
515      not constant across the range of input levels ( $\lambda_{in} : 0 \rightarrow 1$ ), and is highest at intermediate levels of  
516      phylogenetic signal.

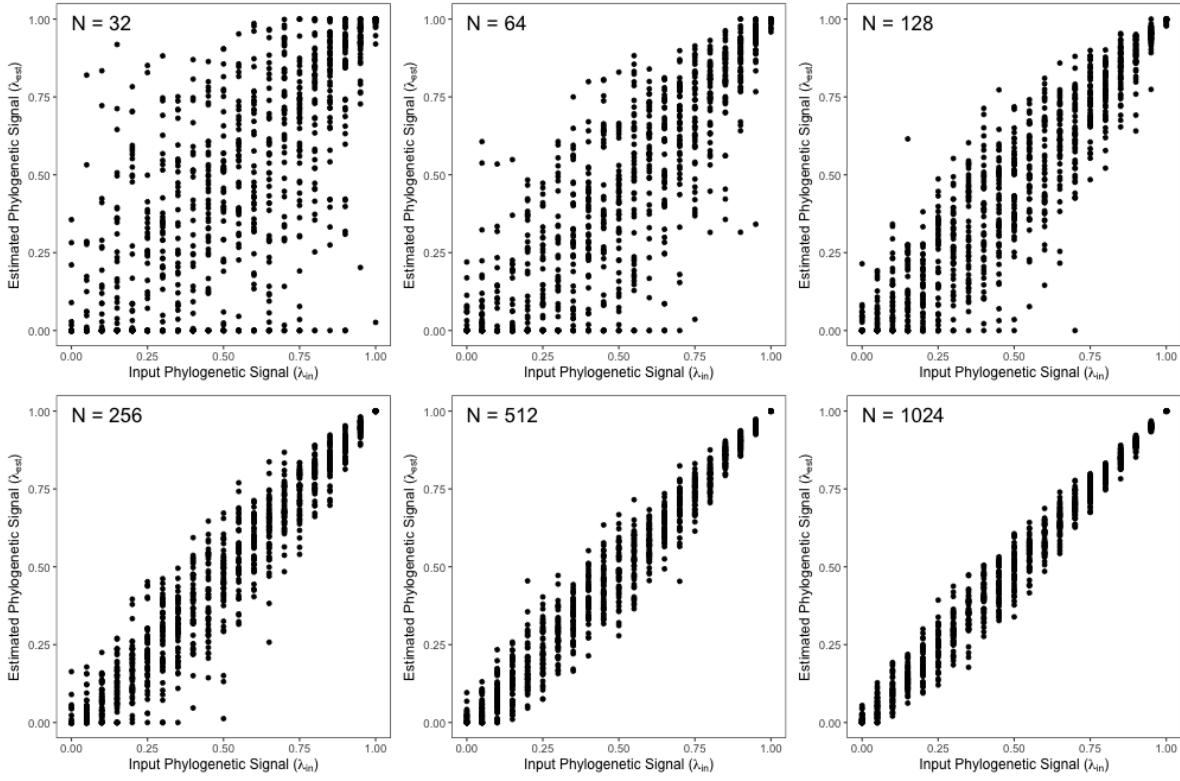
517  
518      **Figure 3.** Precision of Pagel's  $\lambda$  when incorporated in phylogenetic regression ( $Y \sim X$ ), across known  
519      levels of input phylogenetic signal ( $\lambda_{in}$ ) on phylogenies of various sizes. As phylogenies increase in  
520      size, variation in  $\lambda_{in}$  decreases; however the precision is not constant across the range of input levels  
521      ( $\lambda_{in} : 0 \rightarrow 1$ ), and is highest at intermediate levels of phylogenetic signal.

522  
523      **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.  
524      (A) Estimates  $Z_\lambda$  for data simulated on phylogenies with 32 taxa ( $n = 32$ ), (B) Estimates of  $Z_\kappa$  for data  
525      simulated on phylogenies with 32 taxa ( $n = 32$ ), (C) Coefficients of variation of precision estimates of  $Z_\lambda$   
526      and  $Z_\kappa$  across input levels of phylogenetic signal, estimated on phylogenies containing differing numbers  
527      of species.

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

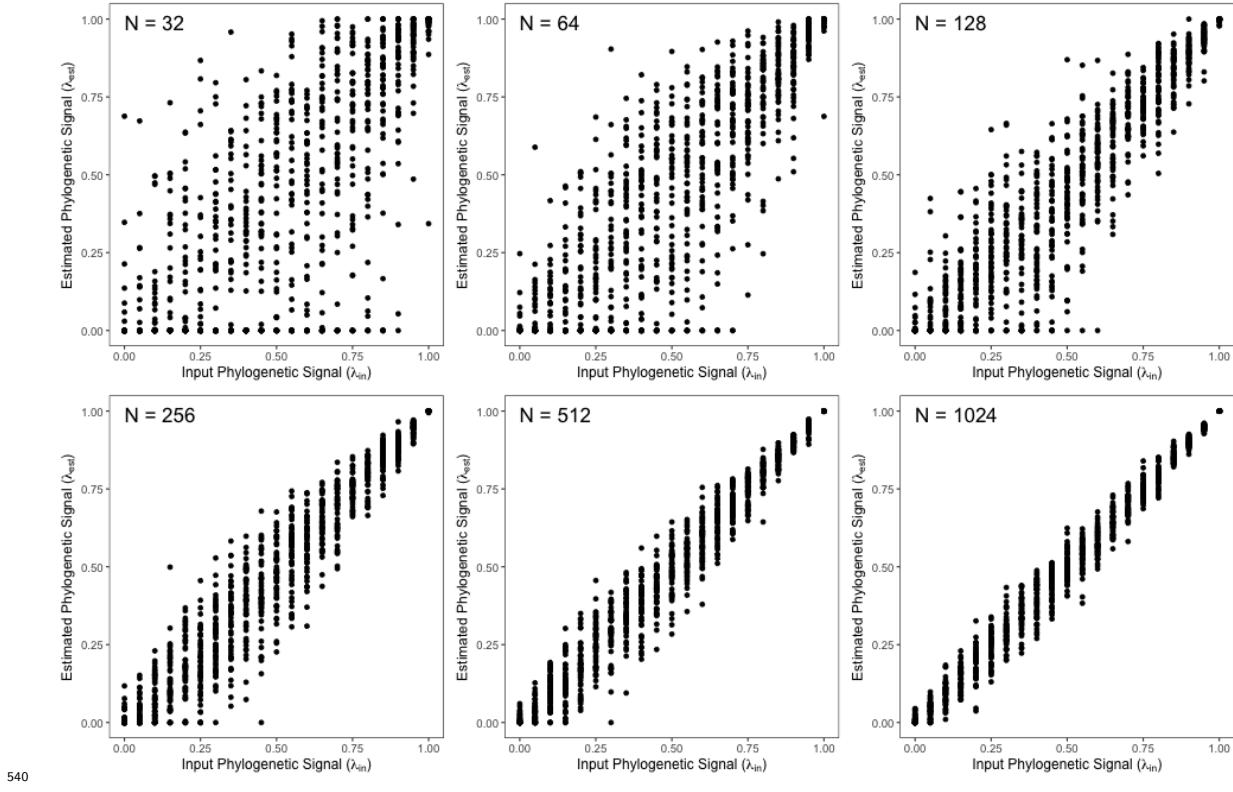


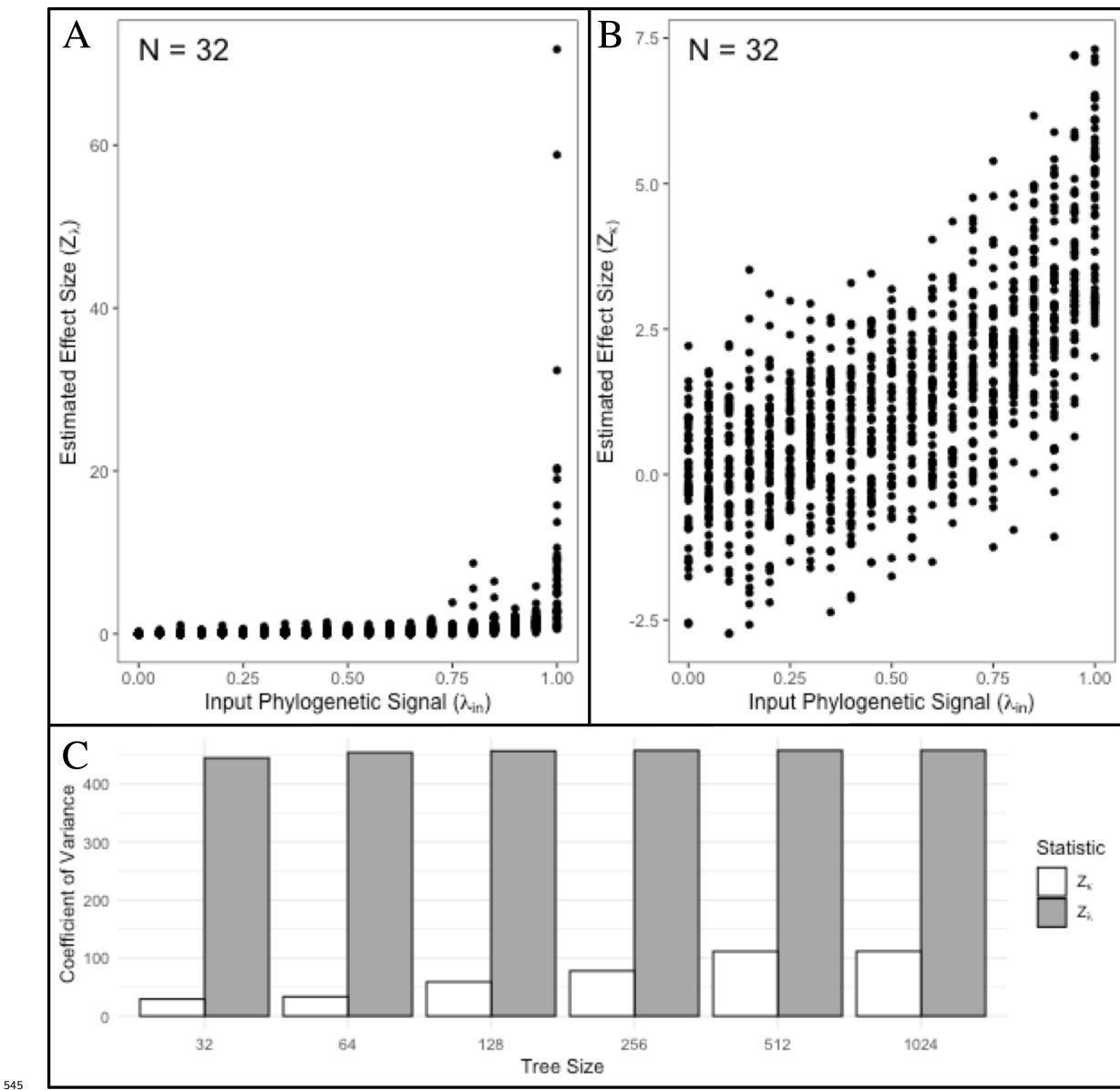
533 **Figure 1.** Frequency distribution of  $\lambda$  estimates published in 2019. The majority of these values were close  
 534 to 0 or 1, and from phylogenies with fewer than 200 taxa.



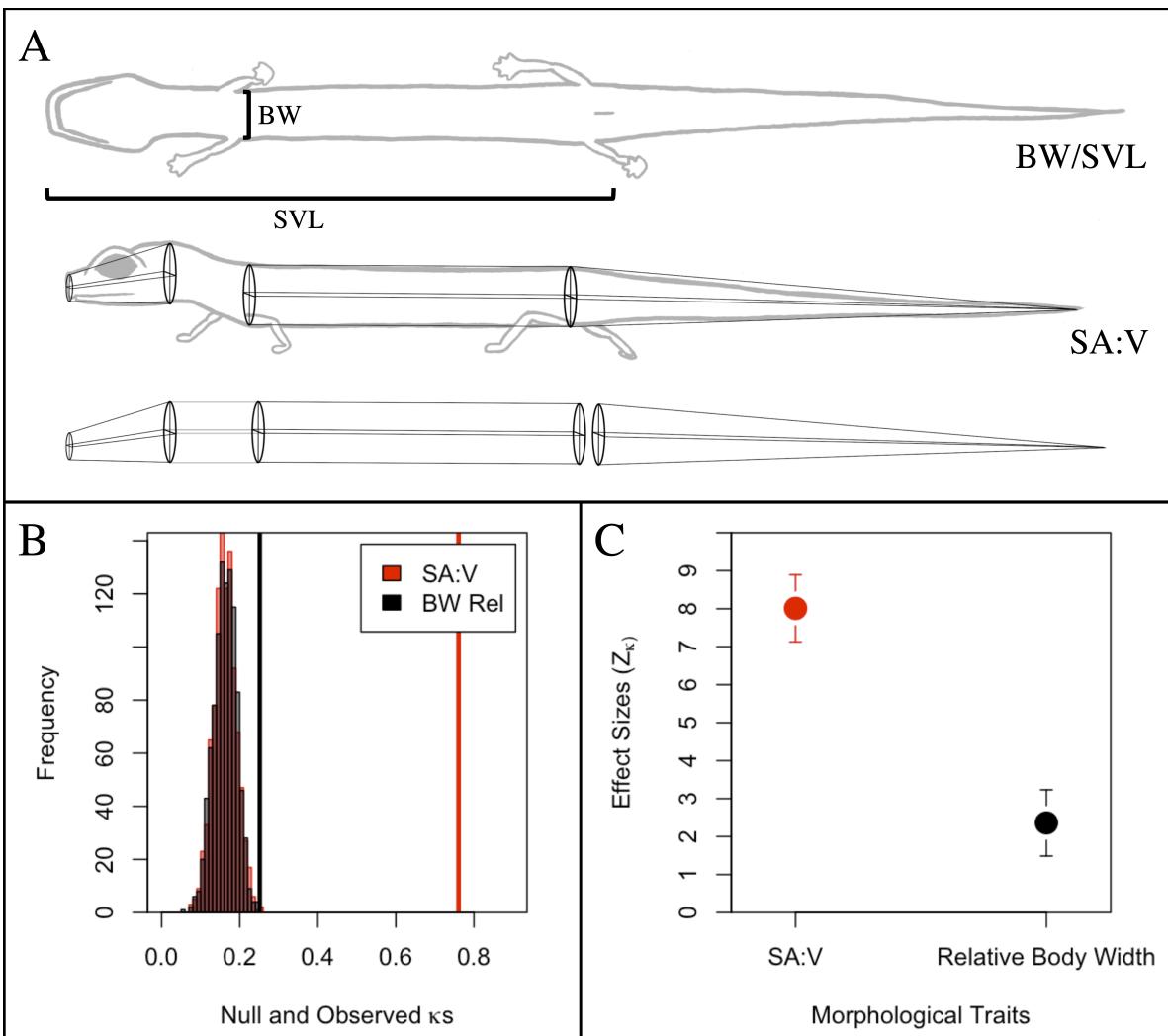
535

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





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



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