

¹ A Standardized Effect Size for Evaluating the Strength of Phylo-
² genetic Signal, and Why Lambda is not Appropriate

³

⁴

⁵ **Abstract**

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

¹⁹ **Introduction**

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

³³

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

⁴⁹

50 In addition to functioning as a parameter that is tuned for appropriate analysis, λ can function as a
51 descriptive statistic to describe the relative strength of phylogenetic signal in phenotypic traits, to determine
52 the extent to which shared evolutionary history has influenced trait covariation among taxa. The appeal
53 of λ as a descriptive statistic for evolutionary biologists is a basis for interpreting “weak” versus “strong”
54 phylogenetic signal; i.e., small versus large values of λ , respectively, in a comparative sense (e.g., De
55 Meester et al. 2019; Pintanel et al. 2019; Su et al. 2019). Indeed, statements regarding the strength of
56 phylogenetic signal based on λ are rather common in the evolutionary literature. For instance, of the 204
57 papers published in 2019 that estimated and reported Pagel’s λ (found from a literature survey we conducted
58 in Google.scholar), 40% interpreted the strength of phylogenetic signal for at least one phenotypic trait.
59 Further, because nearly half of the 1,572 λ values reported were near 0 or 1 (Figure 1) where the biological
60 interpretation of λ is known, this percentage is even higher.

61

62 [insert Figure 1 here]

63

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

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

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

87

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

99

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

114 phylogenetic signal using Pagel's λ are often inaccurate, and thus interpreting strength of phylogenetic signal
115 in phenotypic datasets based on this measure is compromised. By contrast, effect sizes obtained from *Kappa*
116 hold promise for characterizing phylogenetic signal, and for comparing the strength of phylogenetic signal
117 across datasets.

118 **note from Mike** Throughout the ms thus far *Kappa* is used. Why not κ , the actual greek symbol for
119 *Kappa*?

120 Methods and Results

121 *The Precision of λ is Variable*

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

134

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

143 Revell 2010 for a similar simulation design). The fit of the phylogenetic regression was estimated using
144 maximum likelihood, and parameter estimates (β_{est} and λ_{est}) were obtained. We then calculated precision
145 estimates (σ_λ^2) at each input level of phylogenetic signal and level of species richness. We verified that the
146 amount of residual variance simulated had no effect on σ_λ^2 but did influence the precision of coefficients
147 estimated from the linear model (precision increased with smaller ϵ , as expected).

148

149 All analyses were performed in R v3.6.0 (R Core Team 2019) using the packages **geiger** (Harmon et al.
150 2008), **caper** (Orme et al. 2013), **phytools** (Revell 2012), and **geomorph** 3.2.1 (Adams and Otárola-Castillo
151 2013; Adams et al. 2020). R-scripts are found in the Supporting Information.

152

153 *Results.* We found that the precision of λ_{est} varied widely across simulation conditions. Predictably, precision
154 improved as the number of species increased (Figure 2). This confirmed earlier findings of Boettiger et al.
155 (2012), and adhered to parametric statistical theory. However, in many cases the set of λ_{est} spanned nearly
156 the entire range of possible values (e.g., $n = 32$; $\lambda_{in} = 0.5$: $\lambda_{est} = 0.0 \rightarrow 0.985$), revealing that estimates
157 of λ were not a reliable indicator of input phylogenetic signal. Importantly, the precision of λ_{est} was not
158 uniform across all levels of phylogenetic signal. The worst precision was observed at intermediate levels
159 of phylogenetic signal ($\lambda_{in} \approx 0.5$), while precision improved as input levels approached the extremes of
160 λ 's range (i.e., $\lambda_{in} \rightarrow 0$ & $\lambda_{in} \rightarrow 1$). Thus, estimates of λ were least reflective of the true input signal at
161 intermediate values. Additionally, even at large levels of species richness, we found that the range of λ_{est} still
162 encompassed a substantial portion of possible values (e.g., $n = 512$; $\lambda_{in} = 0.5$: $\lambda_{est} = 0.32 \rightarrow 0.68$). Likewise,
163 the same λ_{est} could be obtained from datasets containing vastly different input levels of phylogenetic
164 signal (e.g., $n = 512$; $\lambda_{est} = 0.5$; $\lambda_{in} = 0.25 \rightarrow 0.65$). These findings were particularly unsettling when
165 considered in light of our literature survey. Over one quarter of the λ estimates published in empirical
166 studies (421 of 1,572) were between $\lambda = 0.25$ and $\lambda = 0.75$ (Figure 1). This range reflected the region
167 that our simulations identified as being the least reliable in terms of accurately characterizing levels of
168 phylogenetic signal, yet 30% of these mid-range empirical estimates were explicitly interpreted in terms of
169 the strength of phylogenetic signal that they represented (i.e., weak, intermediate, strong phylogenetic signal).

170

171 Finally, when λ was co-estimated with regression parameters in PGLS regression and ANOVA, the results of
172 our simulations were quite similar. Regression parameters (β) were accurately estimated, confirming earlier
173 findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal (λ) were
174 less precise (Figure 3; see also Supporting Information), and the spread of λ_{est} was similar to that observed

175 when λ was estimated for only the dependent variable, as in Figure 2. Taken together, these findings
176 reveal that λ_{est} does not precisely characterize observed levels of phylogenetic signal in phenotypic datasets,
177 and that biological interpretations of the strength of phylogenetic signal based on λ may be highly inaccurate.

178

179 [insert Figure 2 here]

180

181 [insert Figure 3 here]

182

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

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

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

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

197

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

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

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

205

206 An effect size based on λ could be envisioned, which is found as:

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

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

210

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

221

222 [insert Figure 4 here]

223 ***Statistical Comparisons of Phylogenetic Signal***

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

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

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

234

235 ***Empirical Example***

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

249 et al. 2020). Finally, the strength of phylogenetic signal was compared across traits using \hat{Z}_{12} as described
250 above (to be incorporated in geomorph upon manuscript acceptance).

251

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

264 Discussion

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

279 cal interpretations of the strength of signal based on this parameter may innacurately characterize such effects.

280

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

291

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

307

308 One important advantage of the approach advocated here is that the resulting effect sizes (Z_K) are
309 dimensionless, as the units of measurement cancel out during the calculation of Z (Sokal and Rohlf 2012).
310 Thus, Z_K represents the strength of phylogenetic signal on a common and comparable scale – measured

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

325

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

³⁴² **References**

- ³⁴³ Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.
³⁴⁴ Evolutionary Ecology Research 1:895–909.
- ³⁴⁵ Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other
³⁴⁶ high-dimensional multivariate data. Systematic Biology 63:685–697.
- ³⁴⁷ Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other
³⁴⁸ high-dimensional multivariate data. Evolution 68:2675–2688.
- ³⁴⁹ Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative
³⁵⁰ modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- ³⁵¹ Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration
³⁵² across morphometric datasets. Evolution 70:2623–2631.
- ³⁵³ Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,
³⁵⁴ and a refined permutation procedure. Evolution 72:1204–1215.
- ³⁵⁵ Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate
³⁵⁶ phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- ³⁵⁷ Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric
³⁵⁸ analyses. R package version 3.2.1.
- ³⁵⁹ Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of
³⁶⁰ geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- ³⁶¹ Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.
³⁶² Trends in Ecology and Evolution 10:236–240.
- ³⁶³ Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with
³⁶⁴ phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil
³⁶⁵ 434:305–326.
- ³⁶⁶ Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution
³⁶⁷ 9:7005–7016.

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

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

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

- ⁴⁴⁹ Vandelook, F., S. Janssens, P. Gijbels, E. Fischer, W. Van den Ende, O. Honnay, and S. Abrahamczyk. 2019.
⁴⁵⁰ Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124:269–279.

451 **Figure Legends**

452 **Figure 1.** Frequency distribution of λ estimates published in 2019. The majority of these values were
453 close to 0 or 1, and from phylogenies with fewer than 200 taxa.

454

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

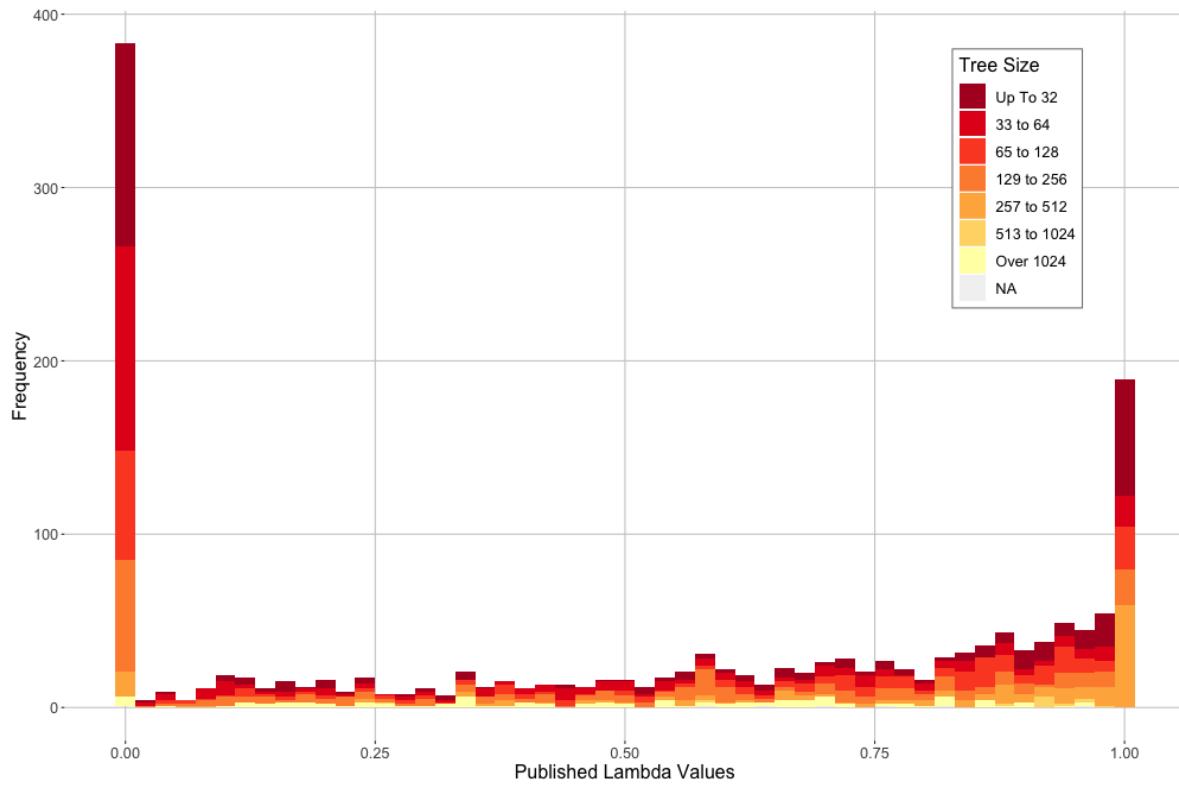
459

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

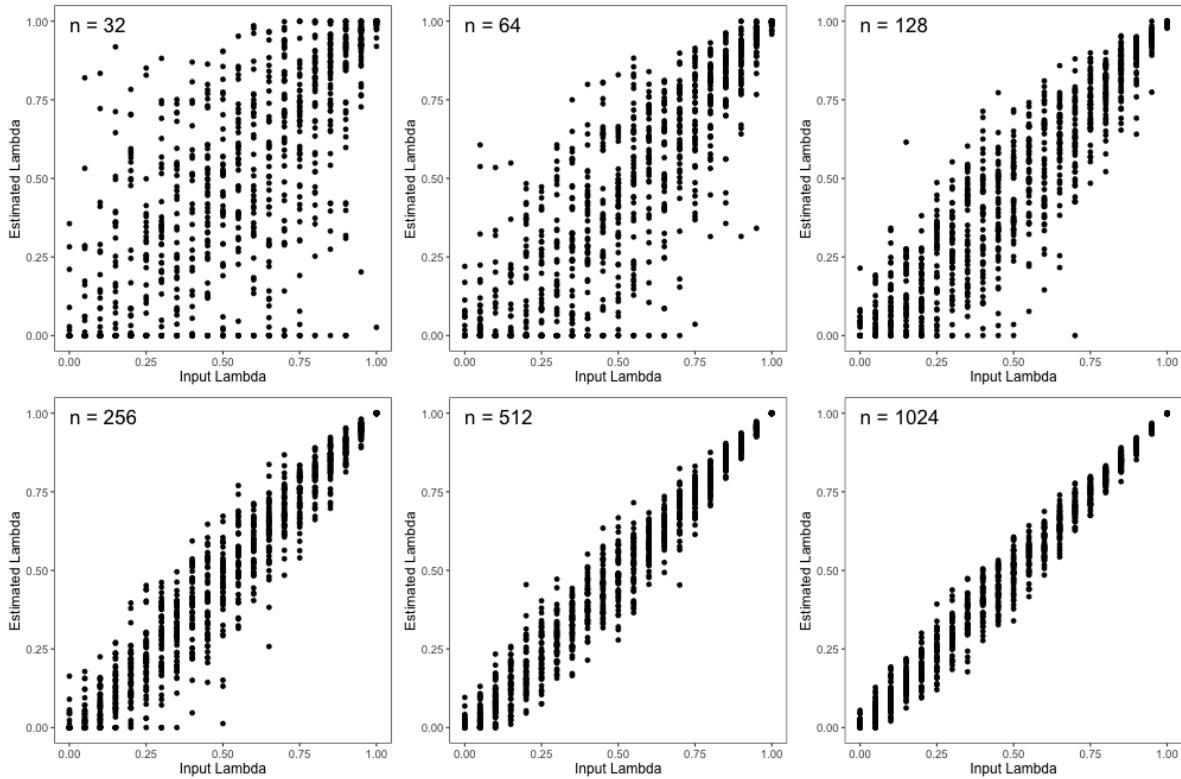
464

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

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

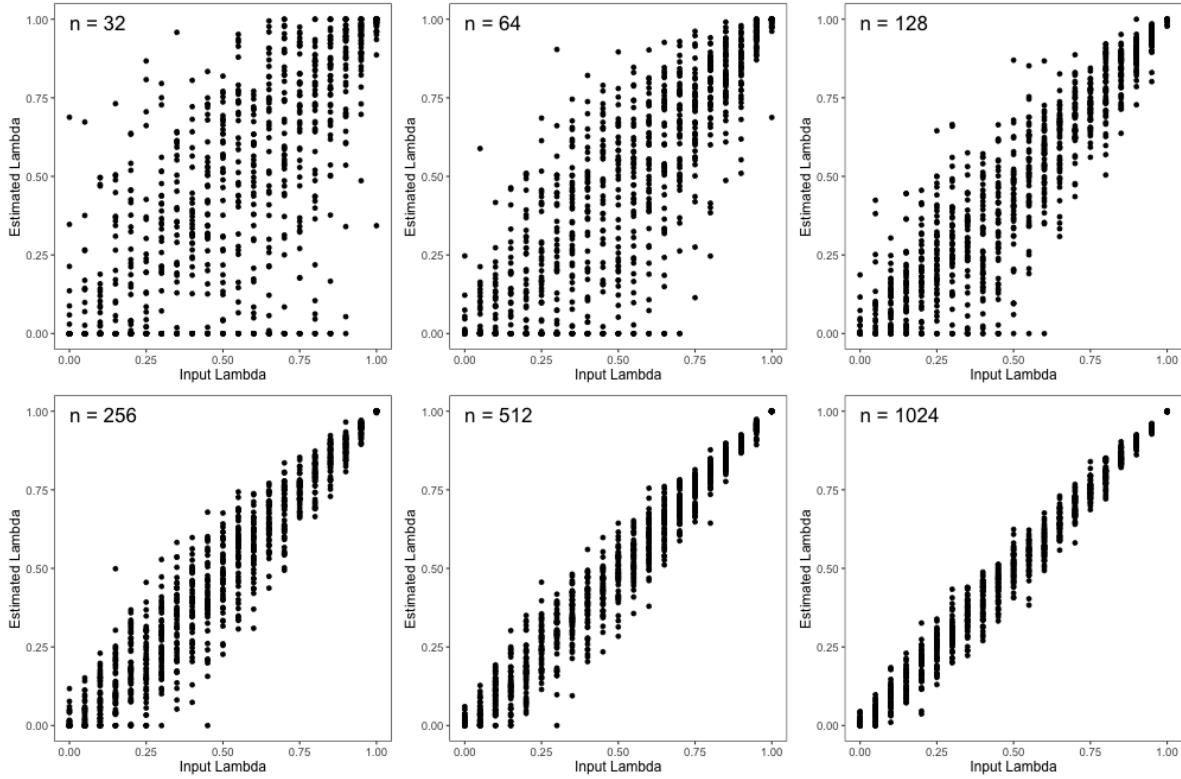


475 **Figure 1.** Frequency distribution of λ estimates published in 2019. The majority of these values were close
476 to 0 or 1, and from phylogenies with fewer than 200 taxa.



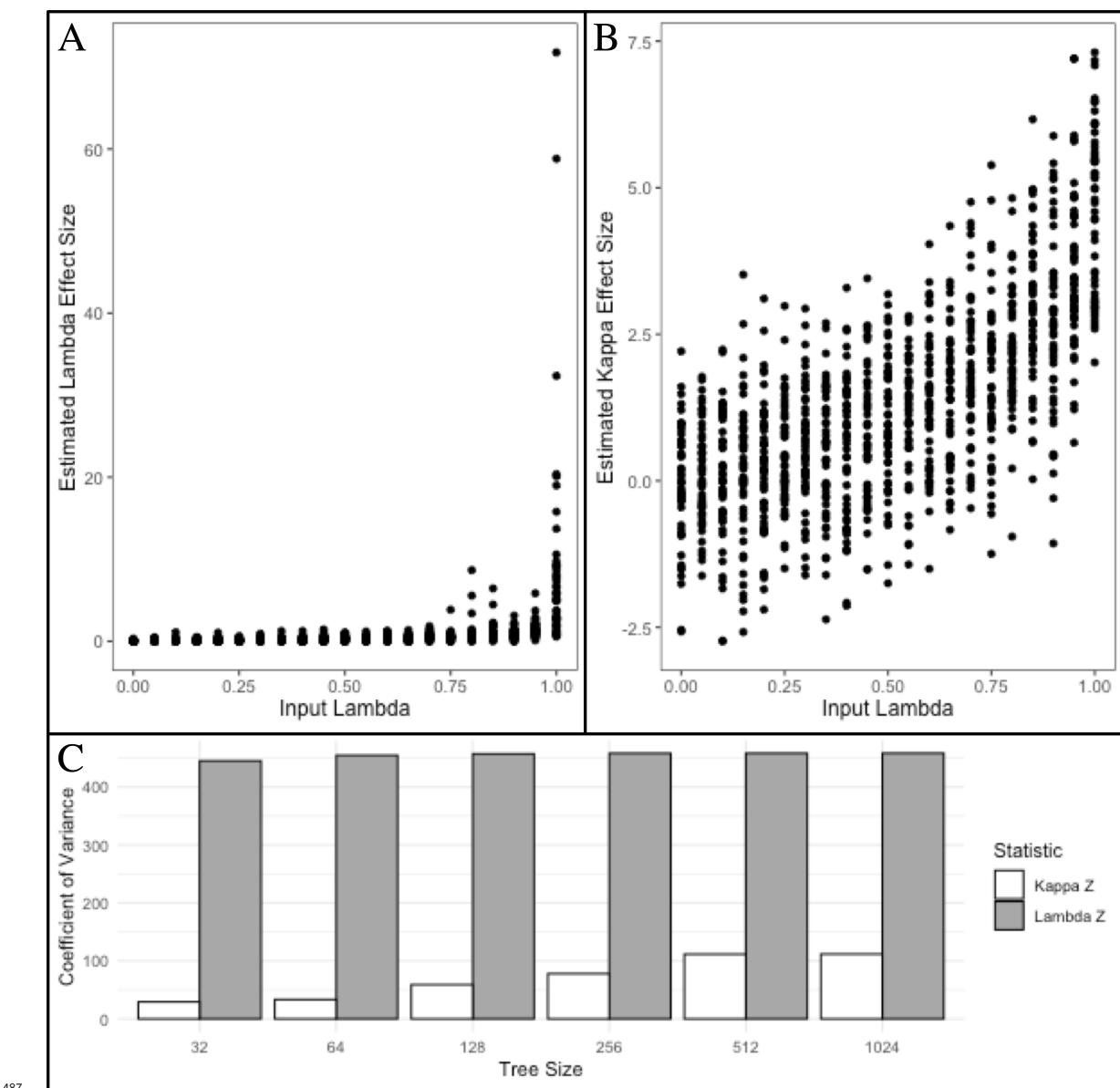
477

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

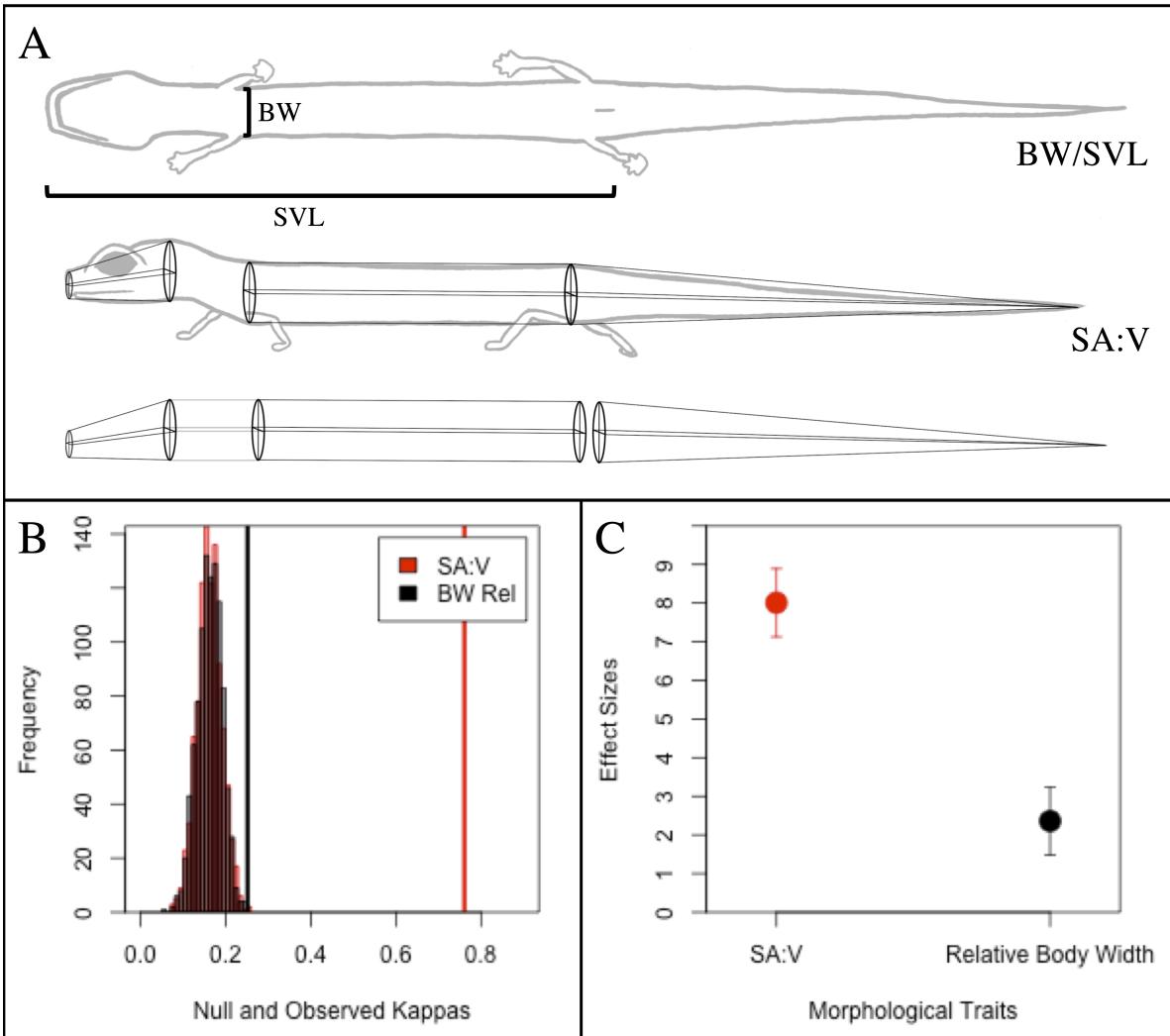


482

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



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



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