

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 κ (Z_κ),
¹³ which measures the strength of phylogenetic signal, and places it on a common scale for statistical comparison.
¹⁴ Tests based on Z_κ provide a mechanism for formally comparing the strength of phylogenetic signal across
¹⁵ datasets, in much the same manner as effect sizes may be used to summarize patterns in quantitative meta-
¹⁶ analysis. Our approach extends the phylogenetic comparative toolkit to address hypotheses that compare
¹⁷ the strength of phylogenetic signal between various phenotypic traits, even when those traits are found in
¹⁸ different evolutionary lineages or have different units or [sealesscales](#).

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 (κ : Blomberg et al. 2003;
37 Adams 2014a), and scaling parameters used in maximum likelihood fitting of the data to the phylogeny (λ :
38 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., Münkemüller et al. 2012; Pavoine and Ricotta 2012;
41 Diniz-Filho et al. 2012; Adams 2014a; Molina-Vegas and Rodríguez 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 (λ) parameter transforms the lengths of the internal
44 branches of the phylogeny to improve a-a-THE fit of data to the phylogeny via maximum likelihood (Pagel
45 1999; Freckleton et al. 2002). Pagel's λ ranges from 0 → 1, with larger values signifying a greater dependence
46 of observed trait variation on the phylogeny. Pagel's λ also has the appeal that it may be included in
47 phylogenetic generalized least-squares regression (PGLS) to account for the degree of phylogenetic signal in
48 comparative analyses (see 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~~CHARACTERIZING the relative strength of phylogenetic signal
52 in phenotypic traits ~~, to determine the extent to which shared evolutionary history has influenced trait~~
53 covariation among taxa. ~~The appeal~~THIS DESCRIPTIVE UTILITY of λ ~~as a descriptive~~
54 ~~statistic for evolutionary biologists is a~~as a descriptive statistic for evolutionary biologists is a ~~THE~~ basis
55 for interpreting “weak” versus “strong” phylogenetic signal~~, i.e., small versus large values of λ , respectively,~~
56 ~~in~~in a comparative sense (e.g., De Meester et al. 2019; Pintanel et al. 2019; Su et al. 2019). Indeed,
57 statements regarding the strength of phylogenetic signal based on λ are rather common in the evolutionary
58 literature. For instance, of the 204 papers published in 2019 that estimated and reported Pagel’s λ (found
59 from a literature survey we conducted in Google.scholar), 40% EXPLICITLY interpreted the strength of
60 phylogenetic signal for at least one phenotypic trait. Further, because nearly half of the 1,572 λ values
61 reported were near 0 or 1 (Figure 1) where the biological interpretation of λ is known, this percentage is even
62 higher.

63

64 [insert Figure 1 here]

65

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

73

74 It seems intuitive to interpret the strength of phylogenetic signal based on the value of λ , as λ is a parameter
75 on a bounded scale ($0 \rightarrow 1$) for which interpretation of its extremal points are understood. Specifically,
76 $\lambda = 0$ represents no phylogenetic signal, while $\lambda = 1$ is phylogenetic signal as expected under Brownian
77 motion. However, equating values of λ directly to the strength of phylogenetic signal presumes two important
78 statistical properties that have not been fully explored. First, it presumes that values of λ can be precisely
79 estimated, as biological inferences regarding the strength of phylogenetic signal depend on high accuracy in
80 its estimation. Therefore, understanding the precision in estimating λ is paramount. One study (Boettiger et
81 al. 2012) found that estimates of Pagel’s λ displayed less variation (i.e., greater precision) when data were

82 simulated on a large phylogeny ($N = 281$) as compared to a small one ($N = 13$). From this observation it
83 was concluded that insufficient data (i.e., ~~the~~the ~~LOW~~ number of species) was the underlying cause of the
84 increased variation across parameter estimates (Boettiger et al. 2012). Indeed, such a pattern is common
85 with statistical estimators, as summary statistics and parameters are often more precise at greater sample
86 sizes (Cohen 1988). However, this conclusion also implies that the precision of λ remains constant across
87 its range ($\lambda = 0 \rightarrow 1$); an assumption that to date, has not been verified. Thus, despite widespread use of
88 Pagel's (1999) λ in macroevolutionary studies, at present, we lack a general understanding of the precision
89 with which λ can estimate levels of phylogenetic signal in phenotypic datasets.

90

91 Second, while estimates of λ are within a bounded scale ($0 \rightarrow 1$), this does not *de-facto* imply that the
92 estimated values of this parameter correspond to the actual strength of the underlying input signal in the
93 data. For this to be the case, λ 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). Unfortunatley, not all model
98 parameters and descriptive statistics are effect sizes, 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 λ is a statistical effect size can comparisons of
101 estimates across datasets be interpretable. ~~For the ease of λ , this~~For the ease of λ , this ~~HOWEVER, THE~~
102 CALCULATION AND STATISTICAL BEHAVIOR OF λ AS AN EFFECT SIZE has not yet been explored.

103

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

114 signal: λ and κ . Through simulations we find that the precision of effect sizes based on λ (Z_λ) are less reliable
115 than that those based on κ (Z_κ), implying that Z_κ is a more robust effect size measure. We also propose a
116 two-sample test statistic that may be used to compare the strength of phylogenetic signal among datasets,
117 and provide an empirical example to demonstrate its use. We conclude that estimates of phylogenetic signal
118 using Pagel's λ are often inaccurate, and thus interpreting strength of phylogenetic signal in phenotypic
119 datasets based on this measure is compromised. By contrast, effect sizes obtained from κ hold promise for
120 characterizing phylogenetic signal, and for comparing the strength of phylogenetic signal across datasets.

121 Methods and Results

122 *The Precision of λ is Variable*

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

135

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

under a Brownian motion model of evolution on each rescaled phylogeny: $\epsilon = \mathcal{N}(\mu = 0, \sigma = \sigma^2 \mathbf{C})$: (see Revell 2010 for a similar simulation design). The fit of the phylogenetic regression was estimated using maximum likelihood, and parameter estimates (β_{est} and λ_{est}) were obtained. We then calculated precision estimates (σ_λ^2) at each input level of phylogenetic signal and level of species richness. We verified that the amount of residual variance simulated had no effect on σ_λ^2 but did influence the precision of coefficients estimated from the linear model (precision increased with smaller ϵ , as expected).

149

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

153

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

171

Finally, when λ was co-estimated with regression parameters in PGLS regression and ANOVA, the results of our simulations were quite similar. Regression parameters (β) were accurately estimated, confirming earlier findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal (λ) were less

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

179

180 [insert Figure 2 here]

181

182 [insert Figure 3 here]

183

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

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

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

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

198

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

201 that suggestion, resulting in an effect size of:

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

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

206

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

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

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

211

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

222

223 [insert Figure 4 here]

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

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

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

229 where κ_1 , κ_2 , $\hat{\mu}_{\kappa_1}$, $\hat{\mu}_{\kappa_2}$, $\hat{\sigma}_{\kappa_1}$, and $\hat{\sigma}_{\kappa_2}$ are as defined above for equation 2. The right side of the equation
230 illustrates that if Z_κ has already been calculated for two sampling distributions as in equation 2, the sampling
231 distributions have unit variance for each of the Z_κ statistics. Estimates of significance of \hat{Z}_{12} may be obtained
232 from a standard normal distribution. Typically, \hat{Z}_{12} is considered a two-tailed test, however directional
233 (one-tailed) tests may be specified should the empirical situation require it (see Adams and Collyer 2016,
234 2019a).

235

236 ***Empirical Example***

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

250 signal in each trait was then characterized using κ , which was converted to its effect size (Z_κ) using **geomorph**
251 3.2.1 (Adams and Otárola-Castillo 2013; Adams et al. 2020). Finally, the strength of phylogenetic signal was
252 compared across traits using \hat{Z}_{12} as described above (to be incorporated in **geomorph** upon manuscript
253 acceptance).

254

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

267 Discussion

268 It is common in comparative evolutionary studies to characterize the phylogenetic signal in phenotypic traits
269 to determine the extent to which shared evolutionary history has generated trait covariation among taxa.
270 However, while numerous analytical approaches may be used to quantify phylogenetic signal (e.g., Abouheif
271 1999; Gittleman and Kot 1990; Pagel 1999; Blomberg et al. 2003; Adams 2014a), methods that explicitly
272 measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained
273 underdeveloped. In this study, we evaluated the precision of one common measure, Pagel's λ , and explored its
274 efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations,
275 we found that the precision of λ increased with increasing sample sizes; a pattern noted previously (Boettiger
276 et al. 2012), and one that conformed with parametric statistical theory (Cohen 1988). However, we also found
277 that vastly different λ estimates could be obtained from data containing the same level of phylogenetic signal,
278 and that similar λ estimates may be obtained from data containing differing levels of phylogenetic signal.
279 Further, the precision of λ varied with the strength of phylogenetic signal, where lower precision was observed

when in data whose phylogenetic signal was of intermediate strength. From these findings we conclude that λ is not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and that biological interpretations of the strength of signal based on this parameter may inaccurately characterize such effects.

283

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

294

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

310

One important advantage of the approach advocated here is that the resulting effect sizes (Z_κ) are

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

328

This study is not the first to compare λ and κ for their ability as statistics to measure phylogenetic signal. Our results for λ and κ values are consistent with those found in the simulations performed by Münkemüller et al. (2012), but that study investigated type I error rates and statistical power, finding that λ performed better in both regards, irrespective of species number in trees. Although not the central focus of their study, the same tendency for variable λ and consistent κ at intermediate phylogenetic signal strengths was observed (see Fig. 2, Münkemüller et al. 2012). Recent work by Molina-Venegas and Rodríguez (2017) found that κ but not ~~lambda~~ λ tended to inflate the estimate of phylogenetic signal, leading to moderate type I and type II biases, if polytomous chronograms were used. Their work more thoroughly addressed previous observations of inflated κ for incompletely resolved phylogenetic trees (Davies et al. 2012; Münkemüller et al. 2012). An interesting question is whether an inflated κ value leads to an inflated Z_κ or does a tendency of a particular tree to inflate estimates of κ also inflate the values in random permutations of a test, in which case Z_κ is robust to polytomies? We repeated the analyses in Figure 4, adjusting trees to have 50% collapsed nodes, per the technique of Molina-Venegas and Rodríguez (2017), and found results were consistent (see Supporting Information). This confirms that any tendency of incompletely resolved trees to inflate κ as a descriptive statistic does not inflate Z_κ as an effect size, ~~and~~, and, FURTHERMORE, because comparison of effect

344 sizes in a test is a comparison of locations of observed values in their sampling distributions, ~~if the sampling~~
345 ~~distributions also shift if the sampling distributions also~~ WHICH WOULD shift EQUALLY because of this
346 tendency, the Z_{12} test statistic in equation 4 appears to be robust in spite of unresolved trees.

347

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

365 **References**

- 366 Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.
367 Evolutionary Ecology Research 1:895–909.
- 368 Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other
369 high-dimensional multivariate data. Systematic Biology 63:685–697.
- 370 Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other
371 high-dimensional multivariate data. Evolution 68:2675–2688.
- 372 Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative
373 modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- 374 Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration
375 across morphometric datasets. Evolution 70:2623–2631.
- 376 Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,
377 and a refined permutation procedure. Evolution 72:1204–1215.
- 378 Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate
379 phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- 380 Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric
381 analyses. R package version 3.2.1.
- 382 Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of
383 geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- 384 Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.
385 Trends in Ecology and Evolution 10:236–240.
- 386 Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with
387 phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil
388 434:305–326.
- 389 Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution
390 9:7005–7016.

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

- 418 Diniz-Filho, J. A. F., T. Santos, T. F. Rangel, and L. M. Bini. 2012. A comparison of metrics for estimating
419 phylogenetic signal under alternative evolutionary models. *Genetics and Molecular Biology* 35:673–679.
- 420 Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* 125:1–15.
- 421 Freckleton, R. P., P. H. Harvey, and M. Pagel. 2002. Phylogenetic analysis and comparative data: A test and
422 review of evidence. *American Naturalist* 160:712–726.
- 423 Garland, T. J., and A. R. Ives. 2000. Using the past to predict the present: Confidence intervals for regression
424 equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
- 425 Gittleman, J. L., and M. Kot. 1990. Adaptation: Statistics and a null model for estimating phylogenetic
426 effects. *Systematic Zoology* 39:227–241.
- 427 Glass, G. V. 1976. Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
- 428 Grafen, A. 1989. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London B,*
429 *Biological Sciences* 326:119–157.
- 430 Harmon, L. J., J. T. Weir, C. D. Brock, R. E. Glor, and W. Challenger. 2008. GEIGER: Investigating
431 evolutionary radiations. *Bioinformatics* 24:129–131.
- 432 Harvey, P. H., and M. D. Pagel. 1991. The comparative method in evolutionary biology. Oxford University
433 Press, Oxford.
- 434 Hedges, L. V., and I. Olkin. 1985. Statistical methods for meta-analysis. Elsevier.
- 435 Kelley, K., and K. J. Preacher. 2012. On effect size. *Psychological Methods* 17:137–152.
- 436 Klingenberg, C. P., and N. A. Gidaszewski. 2010. Testing and quantifying phylogenetic signals and homoplasy
437 in morphometric data. *Systematic biology* 59:245–261.
- 438 Liu, H., C. P. Osborne, D. Yin, R. P. Freckleton, G. Jiang, and M. Liu. 2019. Phylogeny and ecological
439 processes influence grass coexistence at different spatial scales within the steppe biome. *Oecologia*
440 191:25–38.
- 441 Martins, E. P., and T. F. Hansen. 1997. Phylogenies and the comparative method: A general approach
442 to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist*
443 149:646–667.
- 444 Molina-Venegas, R., and M. A. Rodríguez. 2017. Revisiting phylogenetic signal; strong or negligible impacts

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

⁴⁷² Su, G., S. Villéger, and S. Brosse. 2019. Morphological diversity of freshwater fishes differs between realms,
⁴⁷³ but morphologically extreme species are widespread. *Global ecology and biogeography* 28:211–221.

⁴⁷⁴ Vandelook, F., S. Janssens, P. Gijbels, E. Fischer, W. Van den Ende, O. Honnay, and S. Abrahameczyk. 2019.
⁴⁷⁵ Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124:269–279.

476 **Figure Legends**

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

479

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

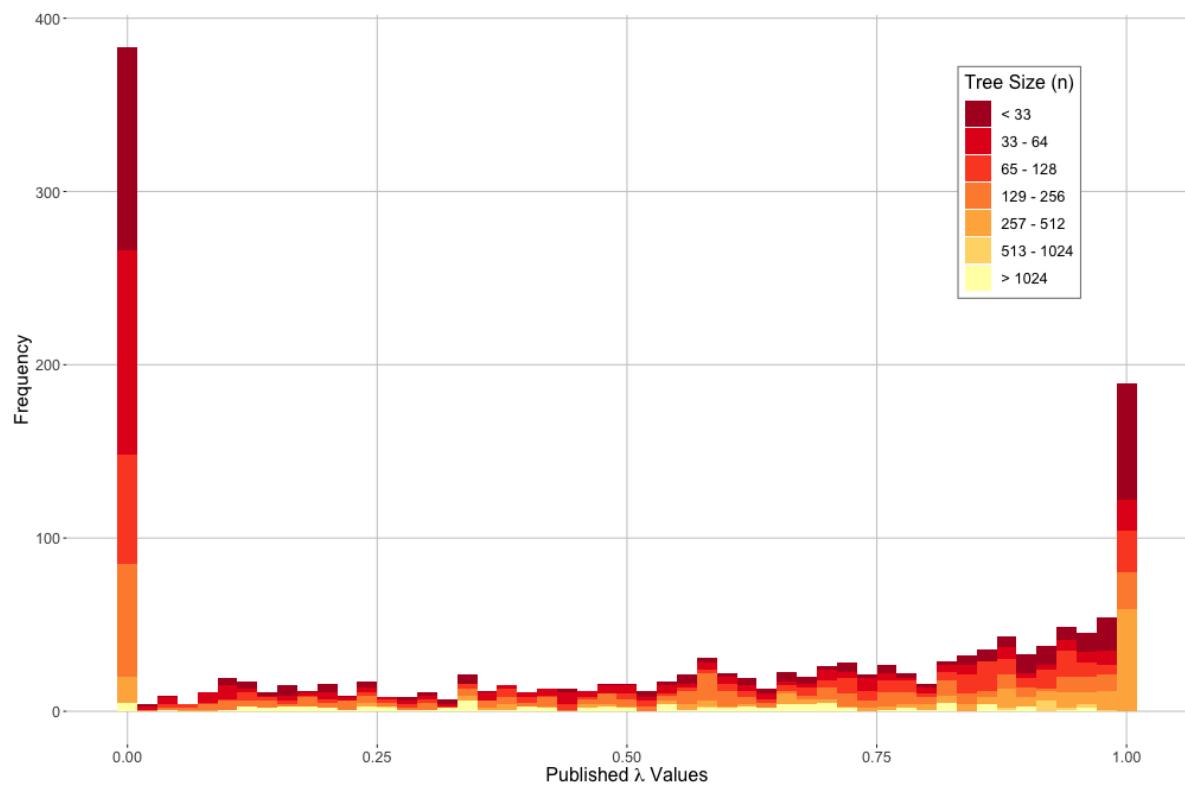
484

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

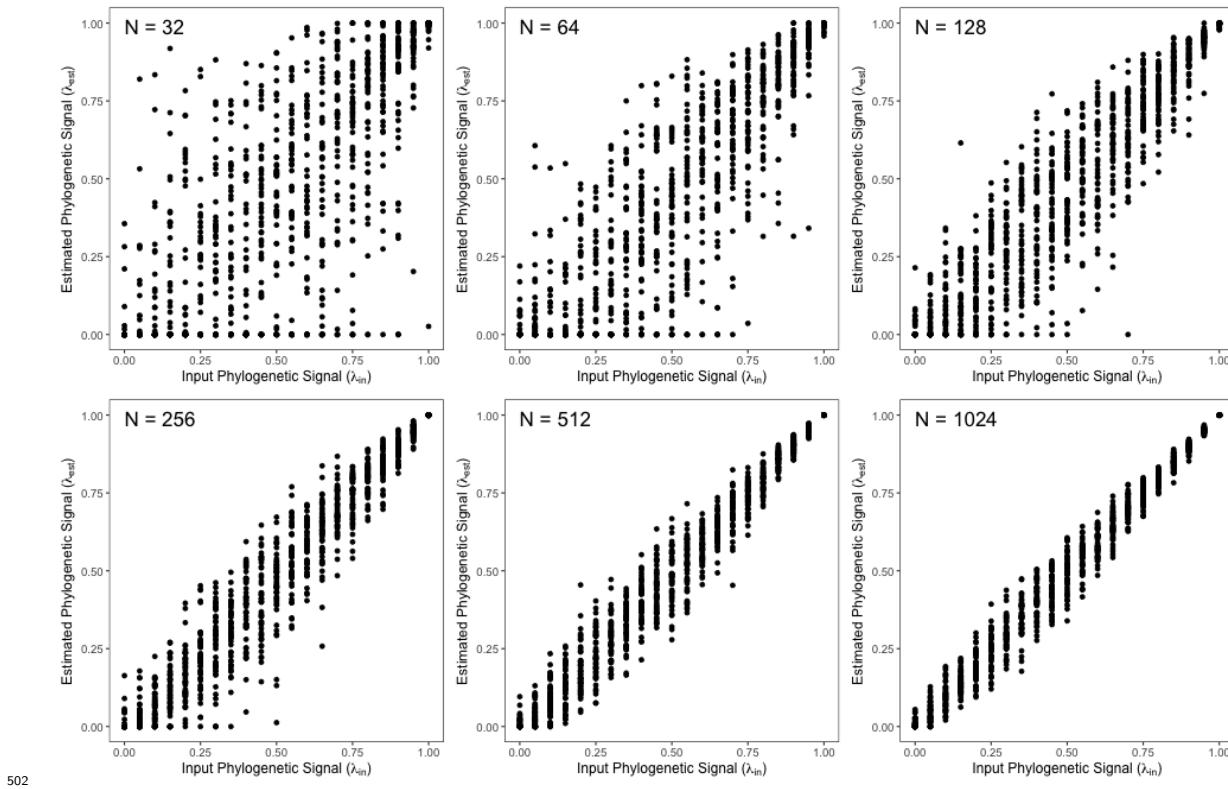
489

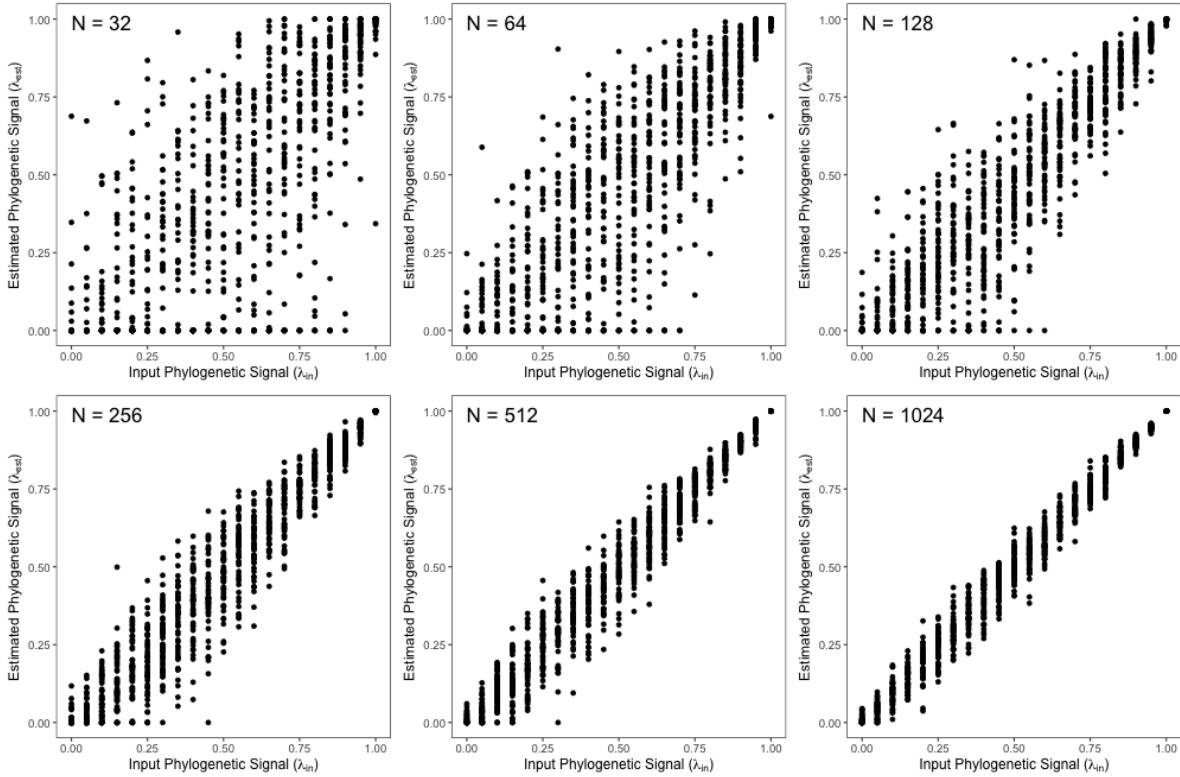
490 **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.
491 (A) Estimates Z_λ for data simulated on phylogenies with 32 taxa ($n = 32$), (B) Estimates of Z_κ for data
492 simulated on phylogenies with 32 taxa ($n = 32$), (C) Coefficients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

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



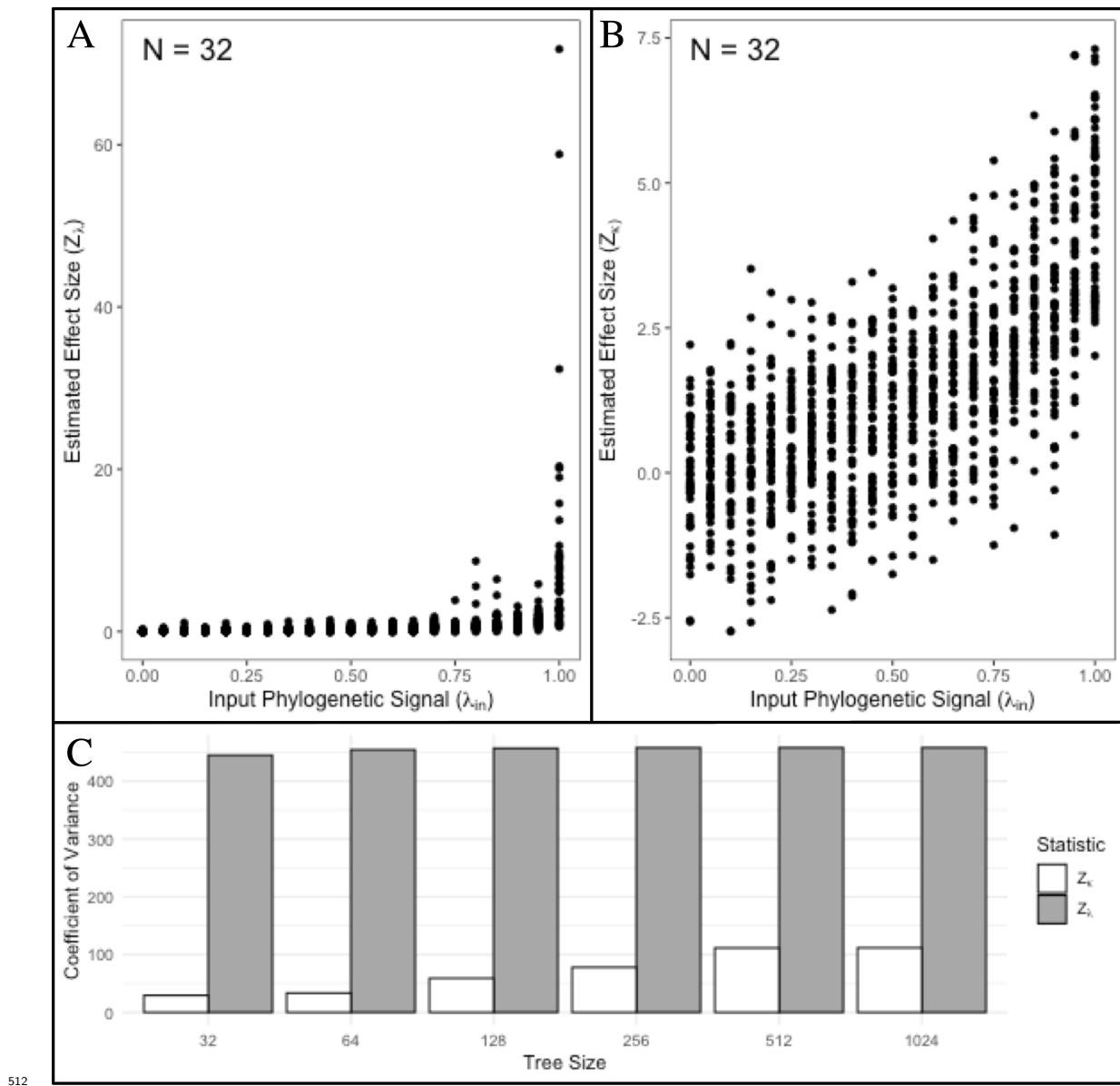
500 **Figure 1.** Frequency distribution of λ estimates published in 2019. The majority of these values were close
501 to 0 or 1, and from phylogenies with fewer than 200 taxa.



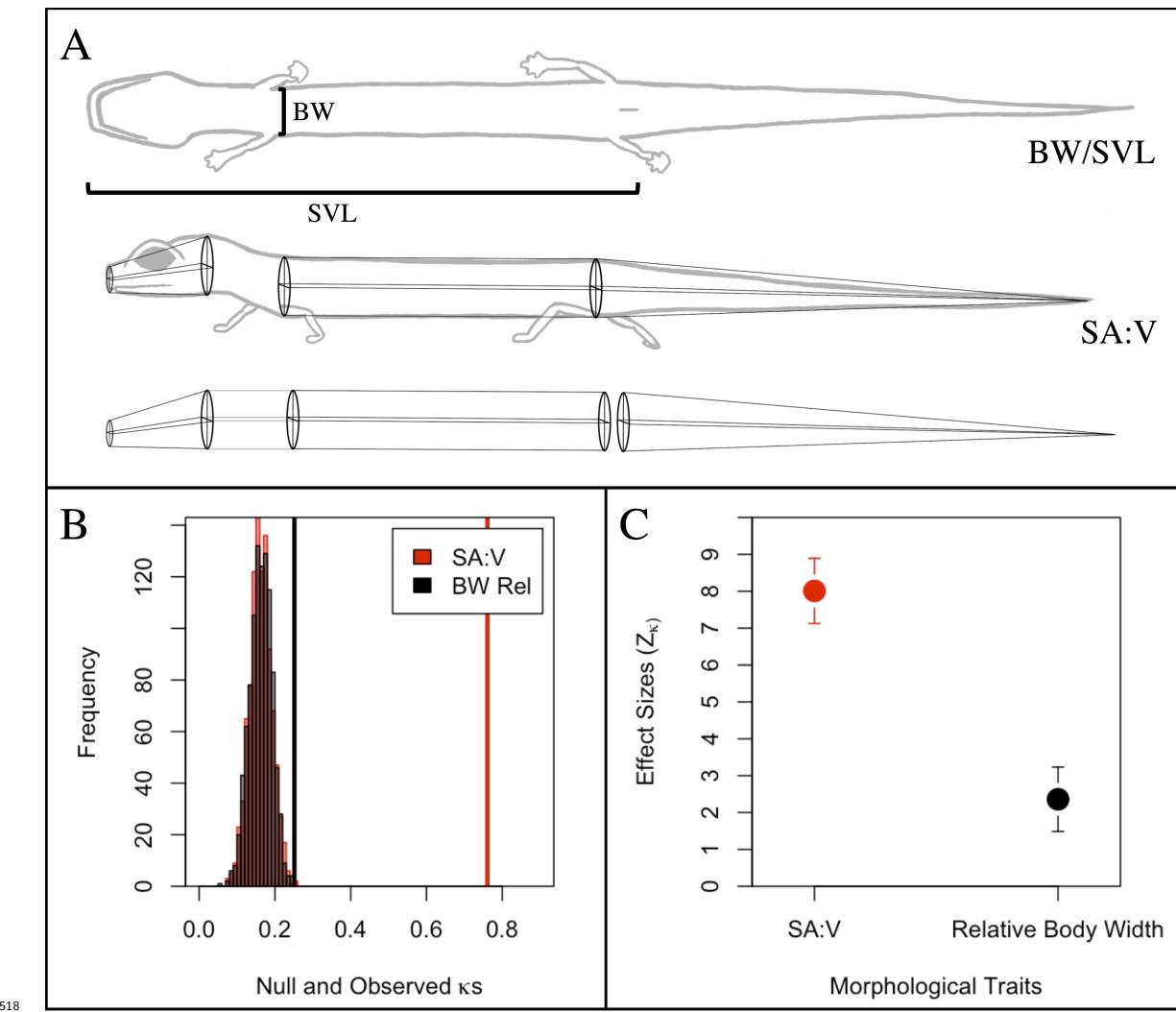


507

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



513 **Figure 4.** Variation in effect size estimates of phylogenetic signal across input levels of phylogenetic signal.
 514 (A) Estimates Z_λ for data simulated on phylogenies with 32 taxa ($n = 32$), (B) Estimates of Z_κ for data
 515 simulated on phylogenies with 32 taxa ($n = 32$), (C) Coefficients of variation of precision estimates of Z_λ
 516 and Z_κ across input levels of phylogenetic signal, estimated on phylogenies containing differing numbers
 517 of species.



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