

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 various phenotypic traits, even when those traits are found
¹⁸ in different evolutionary lineages or of in different units or scale.

¹⁹ **Introduction**

²⁰ Investigating macroevolutionary patterns of trait variation requires a phylogenetic perspective, because
²¹ the shared ancestry among species generates statistical non-independence (Felsenstein 1985; Harvey and
²² Pagel 1991). Accounting for this evolutionary non-independence is the purview of *phylogenetic comparative*
²³ *methods* (PCMs); a suite of analytical tools that condition the data on the phylogeny through the course
²⁴ of statistical evaluations of phenotypic trends (e.g., Grafen 1989; Garland and Ives 2000; Rohlf 2001;
²⁵ Butler and King 2004). The past several decades have witnessed a rapid expansion in the development
²⁶ of PCMs to address an ever-growing set of macroevolutionary hypotheses (Martins and Hansen 1997;
²⁷ O'Meara et al. 2006; Revell and Harmon 2008; Beaulieu et al. 2012; Adams 2014b,a; Adams and
²⁸ Collyer 2018). These methods are predicated on the notion that phylogenetic signal – the tendency for
²⁹ 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). Here, maximum likelihood is used to fit the data to the phylogeny under
⁴⁴ a Brownian motion model of evolution. A parameter (λ) is included, which transforms the lengths of the
⁴⁵ internal branches of the phylogeny to improve the fit (Pagel 1999; Freckleton et al. 2002). Pagel's λ ranges
⁴⁶ from 0 → 1, with larger values signifying a greater dependence of observed trait variation on the phylogeny.
⁴⁷ Pagel's λ also has the appeal that it may be included in phylogenetic regression (PGLS) to account for the
⁴⁸ degree of phylogenetic signal in comparative analyses (see Freckleton et al. 2002).

⁴⁹

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

66

67

68 [insert Figure 1 here]

69

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

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

86

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

98

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

114 in phenotypic datasets based on this measure is compromised. By contrast, effect sizes obtained from *Kappa*
115 hold promise for characterizing phylogenetic signal, and for comparing the strength of phylogenetic signal
116 across datasets.

117 Methods and Results

118 *The Precision of λ is Variable*

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

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

¹⁴³ All analyses were performed in R v3.6.0 (R Core Team 2019) using the packages **geiger** (Harmon et al.
¹⁴⁴ 2008), **caper** (Orme et al. 2013), **phytools** (Revell 2012), and **geomorph** 3.2.1 (Adams and Otárola-Castillo
¹⁴⁵ 2013; Adams et al. 2020). R-scripts are found in the Supporting Information.

¹⁴⁶

¹⁴⁷ *Results.* We found that the precision of λ_{est} varied widely across simulation conditions. Predictably, precision
¹⁴⁸ improved as the number of species increased (Figure 2). This confirmed earlier findings of Boettiger et al.
¹⁴⁹ (2012), and adhered to parametric statistical theory. However, in many cases the set of λ_{est} spanned nearly
¹⁵⁰ the entire range of possible values (e.g., $n = 32$; $\lambda_{in} = 0.5$: $\lambda_{est} = 0.0 \rightarrow 0.985$), revealing that estimates
¹⁵¹ of λ were not a reliable indicator of input phylogenetic signal. Importantly, the precision of λ_{est} was not
¹⁵² uniform across all levels of phylogenetic signal. The worst precision was observed at intermediate levels
¹⁵³ of phylogenetic signal ($\lambda_{in} \approx 0.5$), while precision improved as input levels approached the extremes of
¹⁵⁴ 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
¹⁵⁵ intermediate values. Additionally, even at large levels of species richness, we found that the range of λ_{est} still
¹⁵⁶ encompassed a substantial portion of possible values (e.g., $n = 512$; $\lambda_{in} = 0.5$: $\lambda_{est} = 0.32 \rightarrow 0.68$). Likewise,
¹⁵⁷ the same λ_{est} could be obtained from datasets containing vastly different input levels of phylogenetic
¹⁵⁸ signal (e.g., $n = 512$; $\lambda_{est} = 0.5$; $\lambda_{in} = 0.25 \rightarrow 0.65$). These findings were particularly unsettling when
¹⁵⁹ considered in light of our literature survey. Over one quarter of the λ estimates published in empirical
¹⁶⁰ studies (421 of 1,572) were between $\lambda = 0.25$ and $\lambda = 0.75$ (Figure 1). This range reflected the region
¹⁶¹ that our simulations identified as being the least reliable in terms of accurately characterizing levels of
¹⁶² phylogenetic signal, yet 30% of these mid-range empirical estimates were explicitly interpreted in terms of
¹⁶³ the strength of phylogenetic signal that they represented (i.e., weak, intermediate, strong phylogenetic signal).

¹⁶⁴

¹⁶⁵ Finally, when λ was co-estimated with regression parameters in PGLS regression and ANOVA, the results of
¹⁶⁶ our simulations were quite similar. Here, regression parameters (β) were accurately estimated, confirming
¹⁶⁷ earlier findings of Revell (2010) (see Supporting Information). However, estimates of phylogenetic signal (λ)
¹⁶⁸ were less precise (Figure 3; see also Supporting Information), and the spread of λ_{est} was similar to that
¹⁶⁹ observed when λ was estimated for only the dependent variable, as in Figure 2. Taken together, these findings
¹⁷⁰ reveal that λ_{est} does not precisely characterize observed levels of phylogenetic signal in phenotypic datasets,
¹⁷¹ and that biological interpretations of the strength of phylogenetic signal based on λ may be highly inaccurate.

¹⁷²

¹⁷³ [insert Figure 2 here]

¹⁷⁴

175 [insert Figure 3 here]

176

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

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

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

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

191

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

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

195 where K_{obs} is the observed phylogenetic signal, and $\hat{\mu}_K$ and $\hat{\sigma}_K$ are the mean and standard deviation of the
196 empirical sampling distribution of $\log(Kappa)$ obtained via permutation. Note that the logarithm was used be-

197 cause $Kappa$ takes only positive values ($0 \rightarrow \infty$) and its sampling distribution is log-normally distributed (for a
198 similar transformation when calculating multivariate effect sizes see: Appendix 1 of Collyer and Adams 2018).

199

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

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

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

204

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

215

216 [insert Figure 4 here]

217 ***Statistical Comparisons of Phylogenetic Signal***

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

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

222 where K_1 , K_2 , $\hat{\mu}_{K_1}$, $\hat{\mu}_{K_2}$, $\hat{\sigma}_{K_1}$, and $\hat{\sigma}_{K_2}$ are as defined above for equation 2. Estimates of significance of
 223 \hat{Z}_{12} may be obtained from a standard normal distribution. Typically, \hat{Z}_{12} is considered a two-tailed test,
 224 however directional (one-tailed) tests may be specified should the empirical situation require it (see Adams
 225 and Collyer 2016, 2019a).

226

227 ***Empirical Example***

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

243

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

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

256 Discussion

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

272

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

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

283

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

299

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

₃₁₀ mm^{-1} while relative body size is a unitless ratio ($\frac{BW}{SVL}$). Additionally, our method is capable of comparing
₃₁₁ the strength of phylogenetic signal in traits of different dimensionality, as estimates of phylogenetic signal
₃₁₂ using *Kappa* have been generalized for multivariate data (K_{mult} : see Adams 2014a). Furthermore, tests
₃₁₃ based on \hat{Z}_{12} may be utilized for comparing the strength of phylogenetic signal among datasets containing a
₃₁₄ different number of species, and even for phenotypes obtained from species in different lineages, because
₃₁₅ their phylogenetic non-independence and observed variation are taken into account in the generation of the
₃₁₆ empirical sampling distribution via permutation.

₃₁₇

₃₁₈ **Finally... Need Closing paragraph.**

319 **References**

- 320 Abouheif, E. 1999. A method for testing the assumption of phylogenetic independence in comparative data.
321 Evolutionary Ecology Research 1:895–909.
- 322 Adams, D. C. 2014a. A generalized Kappa statistic for estimating phylogenetic signal from shape and other
323 high-dimensional multivariate data. Systematic Biology 63:685–697.
- 324 Adams, D. C. 2014b. A method for assessing phylogenetic least squares models for shape and other
325 high-dimensional multivariate data. Evolution 68:2675–2688.
- 326 Adams, D. C., and M. L. Collyer. 2019a. Comparing the strength of modular signal, and evaluating alternative
327 modular hypotheses, using covariance ratio effect sizes with morphometric data. Evolution 73:2352–2367.
- 328 Adams, D. C., and M. L. Collyer. 2016. On the comparison of the strength of morphological integration
329 across morphometric datasets. Evolution 70:2623–2631.
- 330 Adams, D. C., and M. L. Collyer. 2018. Phylogenetic anova: Group-clade aggregation, biological challenges,
331 and a refined permutation procedure. Evolution 72:1204–1215.
- 332 Adams, D. C., and M. L. Collyer. 2019b. Phylogenetic comparative methods and the evolution of multivariate
333 phenotypes. Annual Review of Ecology, Evolution, and Systematics 50:405–425.
- 334 Adams, D. C., M. L. Collyer, and A. Kaliontzopoulou. 2020. Geomorph: Software for geometric morphometric
335 analyses. R package version 3.2.1.
- 336 Adams, D. C., and E. Otárola-Castillo. 2013. Geomorph: An r package for the collection and analysis of
337 geometric morphometric shape data. Methods in Ecology and Evolution 4:393–399.
- 338 Arnqvist, G., and D. Wooster. 1995. Meta-analysis: Synthesizing research findings in ecology and evolution.
339 Trends in Ecology and Evolution 10:236–240.
- 340 Bai, K., S. Lv, S. Ning, D. Zeng, Y. Guo, and B. Wang. 2019. Leaf nutrient concentrations associated with
341 phylogeny, leaf habit and soil chemistry in tropical karst seasonal rainforest tree species. Plant and Soil
342 434:305–326.
- 343 Baken, E. K., and D. C. Adams. 2019. Macroevolution of arboreality in salamanders. Ecology and Evolution
344 9:7005–7016.

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

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

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

- ⁴²⁶ Vandelook, F., S. Janssens, P. Gijbels, E. Fischer, W. Van den Ende, O. Honnay, and S. Abrahamczyk. 2019.
⁴²⁷ Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124:269–279.

428 **Figure Legends**

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

431

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

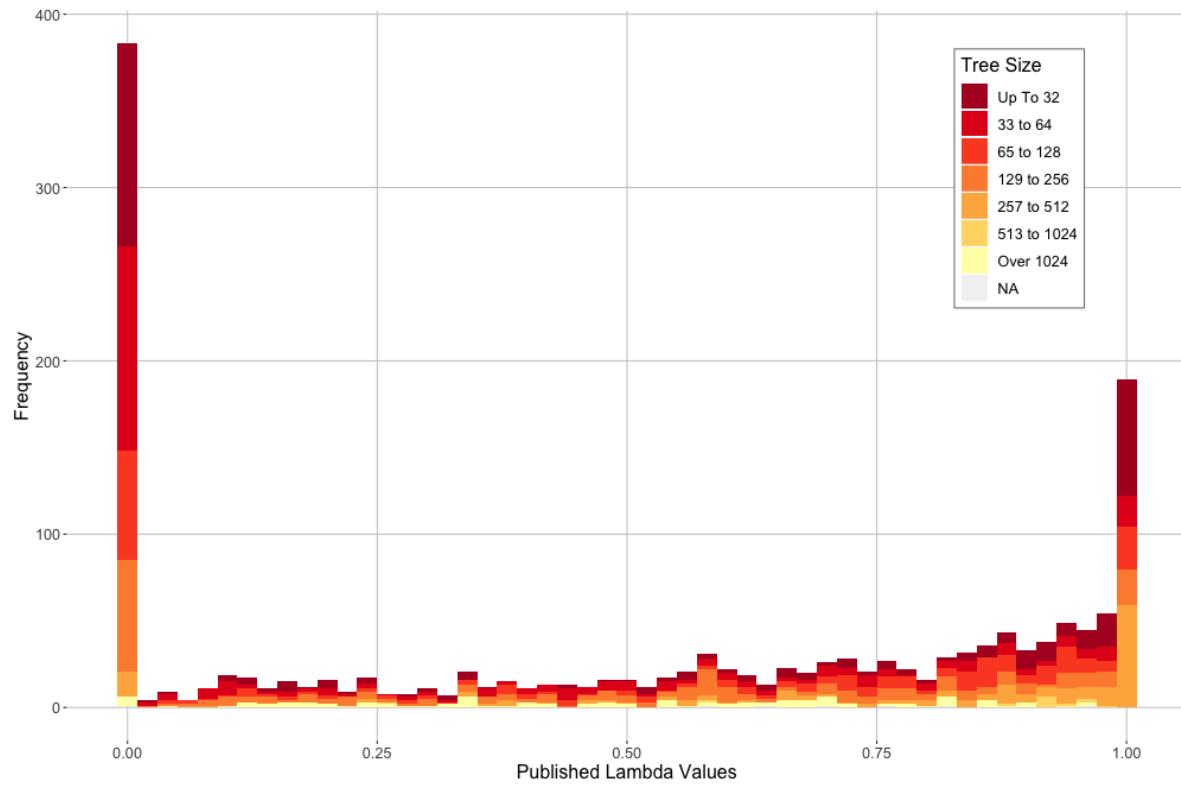
436

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

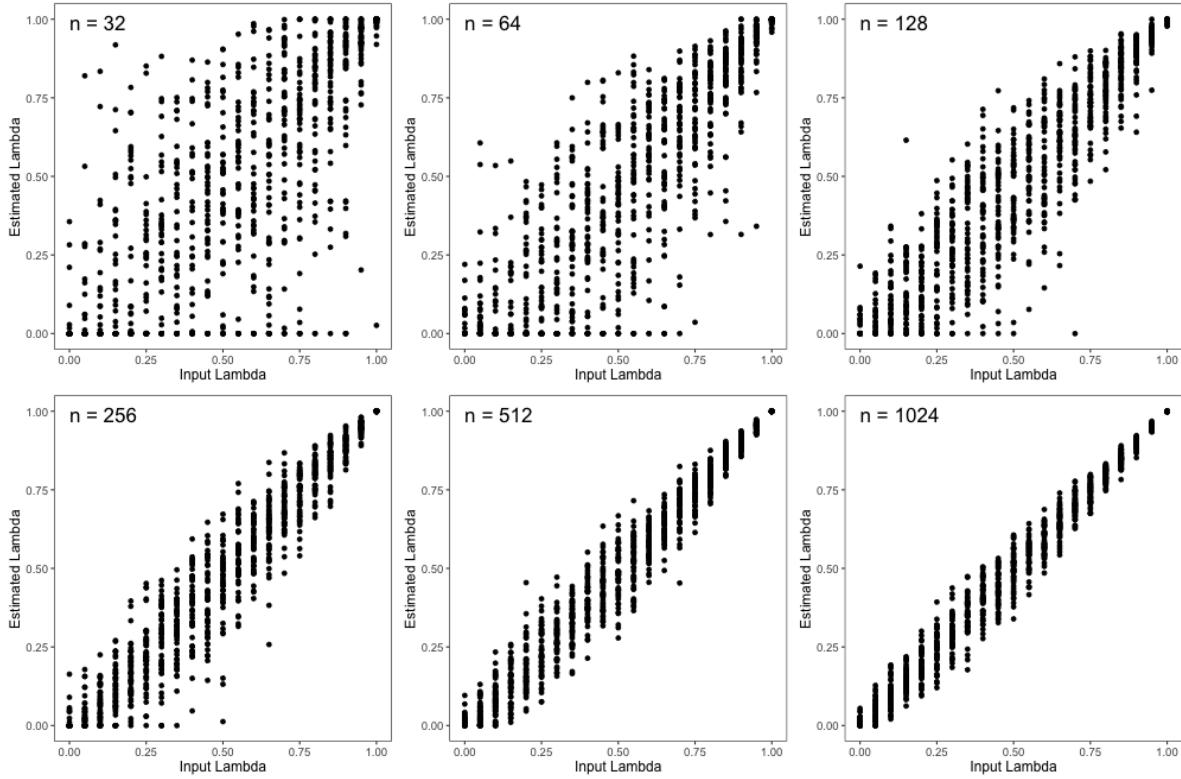
441

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

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

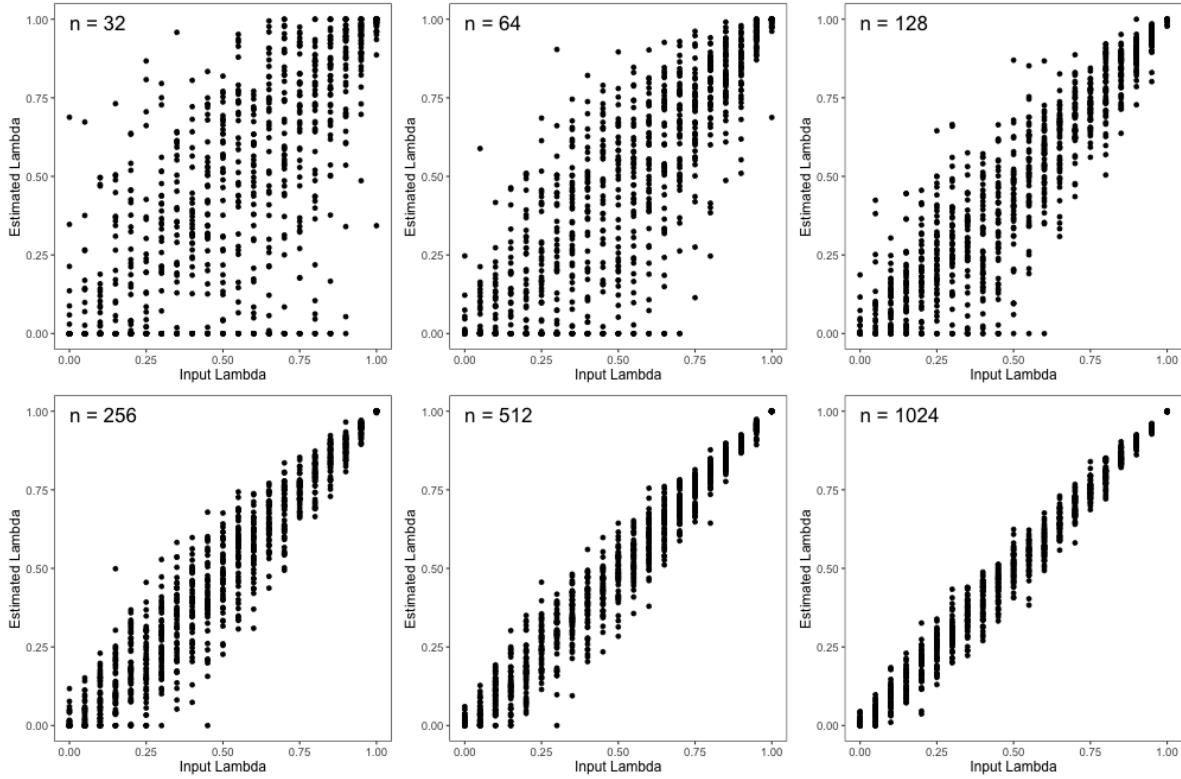


452 **Figure 1.** Frequency distribution of λ estimates published in 2019. The majority of these values were close
453 to 0 or 1, and from phylogenies with fewer than 200 taxa.



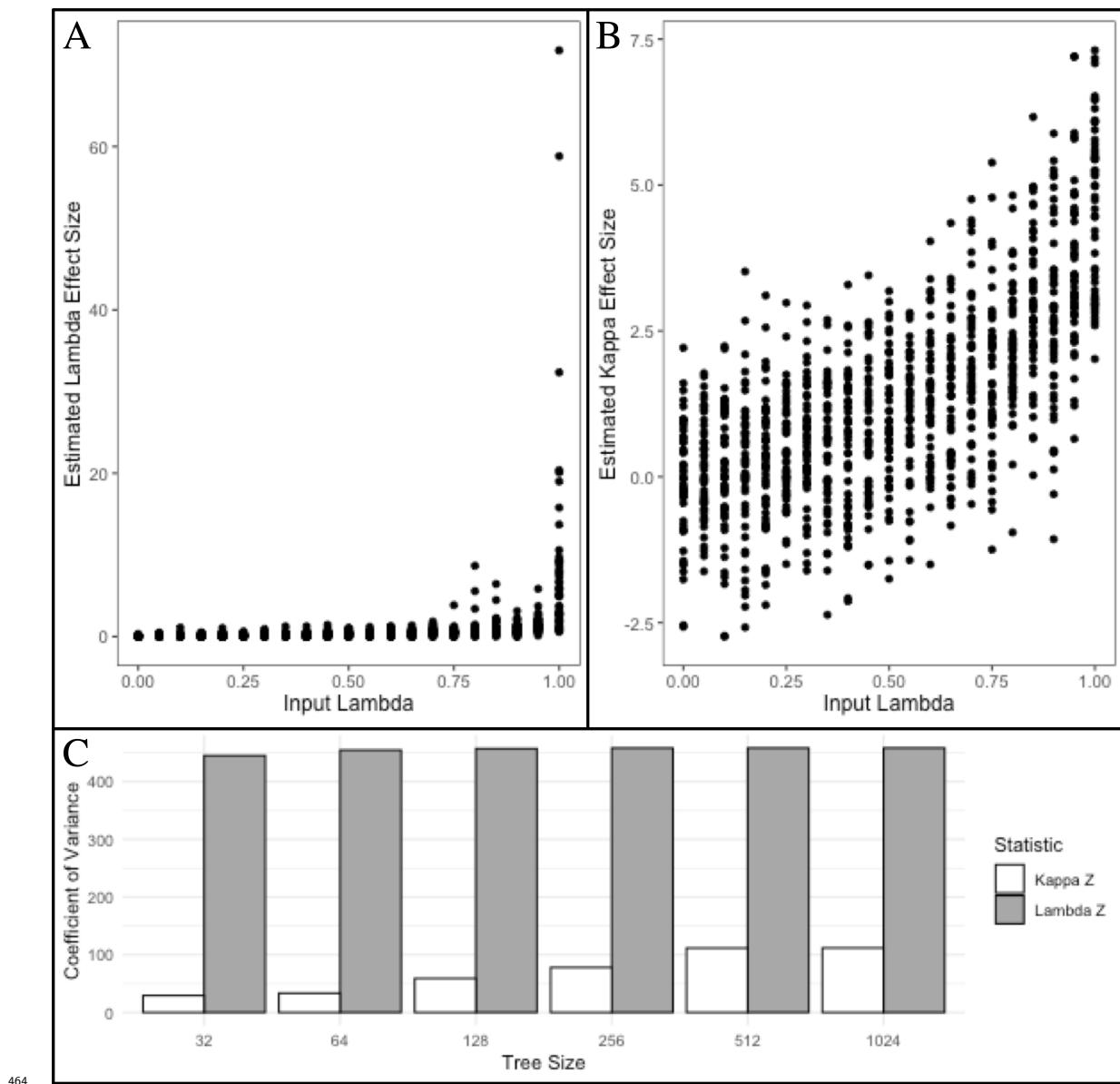
454

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

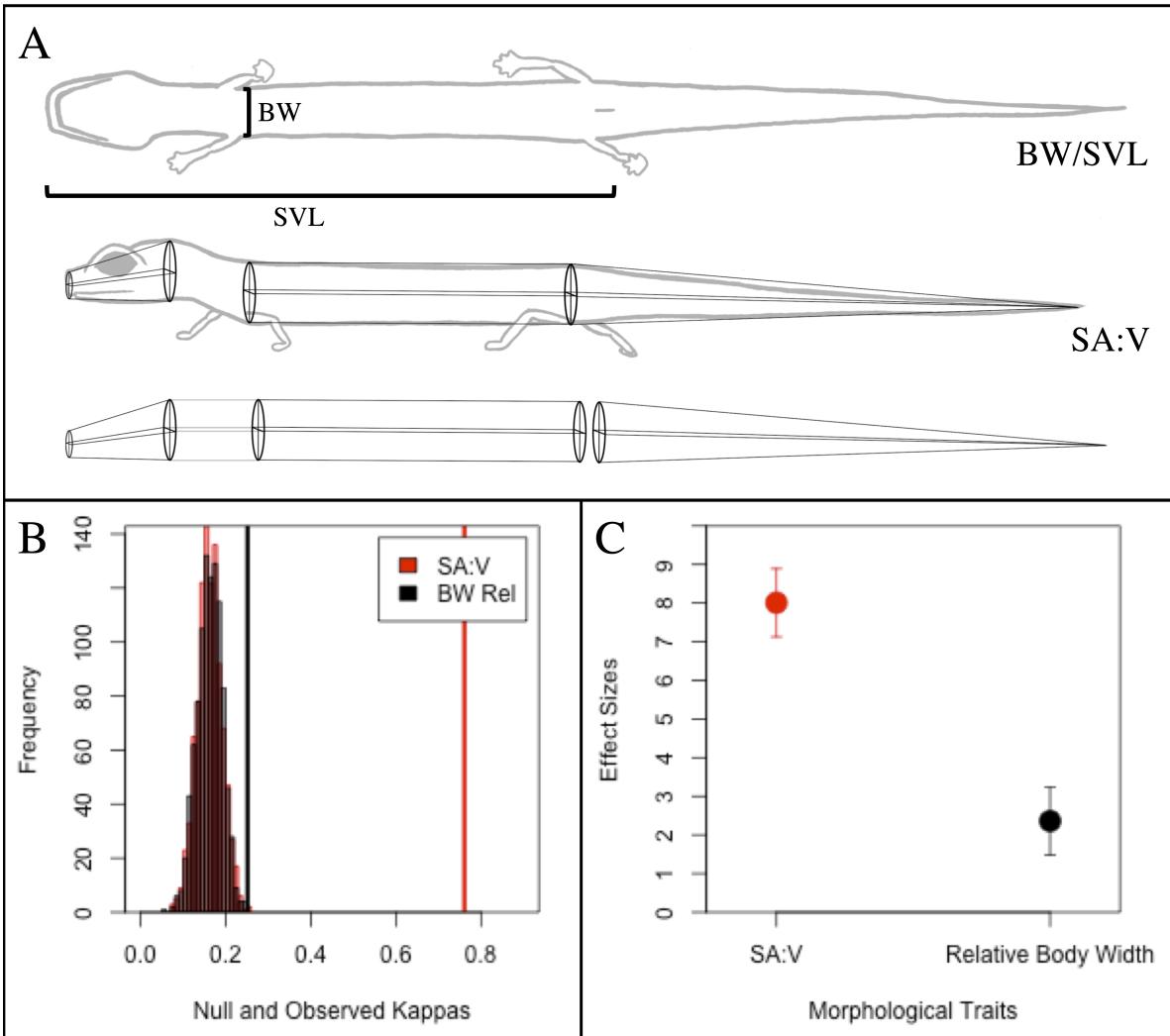


459

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



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.



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