

A Standardized Effect Size for Evaluating and Comparing the Strength of Phylogenetic Signal

Dean C. Adams^{a,2}, Erica K. Baken^{a,b}, and Michael L. Collyer^b

^aDepartment of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames, Iowa, 50010. USA.; ^bDepartment of Science, Chatham University, Pittsburgh, Pennsylvania, 15232. USA.

This manuscript was compiled on August 20, 2020

1 Macroevolutionary studies frequently characterize the phylogenetic
2 signal in phenotypes, however, analytical tools for comparing the
3 strength of that signal across traits remain largely underdeveloped.
4 Here we evaluate the efficacy of Pagel's λ to correctly estimate the
5 strength of phylogenetic signal in phenotypic traits across a range
6 of input values. We find that λ behaves as a Bernoulli random variable,
7 where estimates are increasingly skewed at larger and smaller
8 input levels of phylogenetic signal. Further, the precision of λ varies
9 with input signal. Another measure, Blomberg's κ , is more consistent
10 across a range of tree sizes, and exhibits a positive relationship with input levels of phylogenetic signal. However, that relationship
11 is decidedly nonlinear. Thus, neither λ nor κ are suitable as effect sizes for measuring the strength of phylogenetic signal, and
12 comparing that signal across datasets. As an alternative, we propose a standardized effect size based on κ , (Z_κ), which measures
13 the strength of phylogenetic signal more reliably than does λ , and places that signal on a common scale for statistical comparison. We
14 develop tests based on Z_κ to provide a mechanism for formally comparing the strength of phylogenetic signal across datasets, in much
15 the same manner as effect sizes may be used to summarize patterns in quantitative meta-analysis. Our approach extends the phylogenetic
16 comparative toolkit to address hypotheses that compare the strength of phylogenetic signal between various phenotypic traits,
17 even when those traits are found in different evolutionary lineages or have different units or scales.

phylogenetic signal | macroevolution | lambda | kappa

1 Investigating macroevolutionary patterns of trait variation
2 requires a phylogenetic perspective, because the shared ancestry among species violates the assumption of independence
3 among trait values that is common for statistical tests (1, 2). Accounting for this evolutionary non-independence is the
4 purview of *phylogenetic comparative methods* (PCMs): a suite of analytical tools that condition trends in the data on the
5 phylogenetic relatedness of observations (3–10). These methods are predicated on the notion that phylogenetic signal –
6 the tendency for closely related species to display similar trait values – is present in cross-species datasets (1, 11, 12). Indeed,
7 under numerous evolutionary models, phylogenetic signal is to be expected, as stochastic character change along the hierarchical
8 structure of the tree of life generates trait covariation among related taxa (1, 12, 13).

9 Several analytical tools have been developed to quantify
10 phylogenetic signal in phenotypic datasets (11, 12, 14–17),
11 and their statistical properties – namely type I error rates and
12 statistical power – have been investigated to determine under
13 what conditions phylogenetic signal can be detected (13, 16,
14 18–23). One of the most widely used methods for characterizing
15 phylogenetic signal is Pagel's λ (11), which transforms the lengths of the internal branches of the phylogeny to im-

24 prove the fit of data to the phylogeny via maximum likelihood
25 (11, 24). When incorporated in PGLS, λ serves as a tuning parameter which is optimized via log-likelihood profiling while
26 evaluating the covariation between the dependent and independent variables, given the phylogeny (11, 24). To infer whether
27 phylogenetic signal differs from no signal or a Brownian motion
28 model of evolutionary divergence, the observed model fit using
29 $\hat{\lambda}$ may be statistically compared to that using $\lambda = 0$ or $\lambda = 1$
30 via likelihood ratio tests (24–26) or confidence limits (27).

31 Another widely used measure of phylogenetic signal is
32 Blomberg's κ (12), which characterizes phylogenetic signal as the ratio of observed trait variation to the amount of variation
33 expected under Brownian motion. Blomberg's κ can be treated as a test statistic by employing a permutation test to generate its sampling distribution (12, 16) for determining
34 whether significant phylogenetic signal is present in data. Both
35 λ and κ seem intuitive to interpret, as a value of 0 for both corresponds to no phylogenetic signal, while a value of 1 corresponds to the amount of phylogenetic signal expected under
36 Brownian motion. Thus, it is tempting to regard both λ and κ as descriptive statistics that measure the relative strength
37 of phylogenetic signal, providing an estimate of its magnitude
38 for comparison.

39 The appeal of Pagel's λ and Blomberg's κ as descriptive
40 statistics is that they are based on well-defined statistical
41 procedures that are amenable to formal hypothesis testing.
42

Significance Statement

43 Evolutionary biologists wish to quantify and compare the
44 strength of phylogenetic signal across datasets, but analytical
45 tools for these comparisons are generally lacking. Here we
46 develop a standardized effect size, Z_κ , which measures
47 the strength of phylogenetic signal on a common statistical
48 scale. We also provide a test statistic, \hat{Z}_{12} , for comparing the
49 strength of phylogenetic signal across datasets. We find that
50 two commonly used parameters (Pagel's λ and Blomberg's κ),
51 not converted to effect sizes, are unsuitable for this purpose.
52 Our effect-size procedure enables biologists to quantitatively
53 address hypotheses that compare the strength of phylogenetic
54 signal between various phenotypic traits, even when those traits
55 are found in different evolutionary lineages or have different
56 units or scales.

57 D.C.A. designed the research; D.C.A., E.K.B., and M.L.C. performed the research and wrote the paper.

58 The authors declare no conflict of interest.

59 Data deposition: Data for the empirical example may be found on DRYAD: doi:10.5061/dryad.b554m44 and doi:10.5061/dryad.59zw3r23m. R-scripts for simulation tests are found on Github: XXX. Computer code for implementing the two-sample comparison of effect sizes is found in geomorph: <https://cran.r-project.org/web/packages/geomorph/index.html>

60 ²To whom correspondence should be addressed. E-mail: dcadams@iastate.edu

statistics is that they provide a basis for interpreting “weak” versus “strong” phylogenetic signal; i.e., small versus large values of $\hat{\lambda}$ or κ , respectively, in a comparative sense (28–30). Nonetheless, an important question that has yet to be considered is whether such comparisons are analytically appropriate, and whether these statistics are, or can be, converted to effect sizes for comparative analyses across datasets. To be statistics representing phylogenetic signal, they should have reliable distributional properties, which could be revealed with simulation experiments. For instance, as a proportional random variable bounded by 0 and 1, we might expect that $\hat{\lambda}$ follows a Bernoulli distribution (**add ref**); i.e., branch lengths in a tree are scaled proportionally to the probability that data arise from a BM process. Given a known λ value used to generate random data on a tree, we would also expect that the mean of an empirical sampling distribution of $\hat{\lambda}$ would approximately equal λ ; the dispersion of $\hat{\lambda}$ would be largest at intermediate values of λ , $\hat{\lambda}$ would be predictable over the range of λ with respect to tree size; the distribution of $\hat{\lambda}$ would be symmetric at intermediate values of λ and more skewed toward values of 0 or 1; and that the distribution of $\hat{\lambda}$ will be more platykurtic at intermediate values of λ , becoming more leptokurtic toward 0 and 1 (**add same ref**). Prior work (18) seems to support some of these conjectures, based superficially on statistical moments for a given tree size (mean, variance, skewness, and kurtosis; see Fig. 2 of ref. (18)). However, because the “strength of Brownian motion” was simulated as a varied weighted-average of data simulated on trees with $\lambda = 0$ and $\lambda = 1$ and not as prescribed values of λ (18), interpretation of these patterns is challenging.

By contrast, for Blomberg’s κ , which is positively unbounded, we might expect that for any λ used to generate data, estimates of κ might follow a normal distribution, with values distributed symmetrically about the input value. This attribute seemed less reasonable based on the simulations performed by Münkemüller et al. (18), which suggested that distributions were positively skewed and that Blomberg’s κ might not behave as a statistic that follows a normal distribution. However, because their simulations used a weighted combination of simulated phylogenetic signal strengths, strong inferences are not possible (and distributional attributes were not the intended result of their simulations). Thus, for both Pagel’s λ or Blomberg’s κ , evaluation of statistical moments across a range of λ used to generate data would be valuable for adjudicating the reliability of these statistics as effect sizes. Furthermore, the expected values of these statistics appear to vary with tree size (18), making comparisons across studies challenging. Therefore, transformation of these statistics into Z-scores would allow evaluation of the efficacy of each statistic to yield effect sizes that could be used for comparisons of the strength of phylogenetic signal across traits and lineages.

Here we use simulation experiments to compare the distributional attributes of $\hat{\lambda}$ and κ , plus their effect sizes (Z-scores), across a range of tree size and phylogenetic signal strength. We find that estimates of $\hat{\lambda}$ are increasingly skewed at larger and smaller input levels of phylogenetic signal and at smaller tree sizes, vary widely for a given input value of λ , and that the precision of $\hat{\lambda}$ is not constant across its range. By contrast, estimates of κ are more consistent across tree sizes, and are normally distributed across the range of input levels of λ , making κ a more reliable statistic. We then propose an effect size based on κ , (Z_κ), which provides consistent estimates

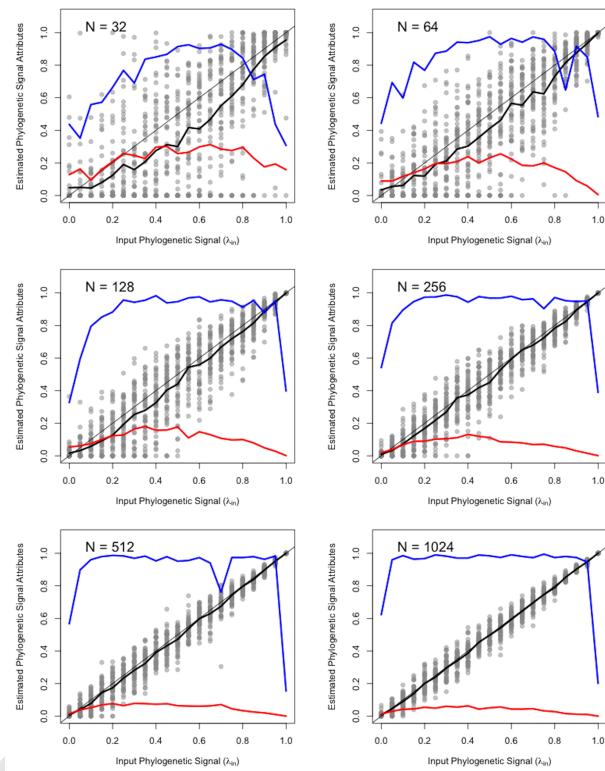


Fig. 1. Response of Pagel’s λ to increasing strength of Brownian motion. Gray line signifies the 1:1 line where the input value matches the estimate $\hat{\lambda}$. At each input level, the dark black line represents the empirically derived expected value (mean) of $\hat{\lambda}$, the red line is the standard deviation of $\hat{\lambda}$, and the blue line is Shapiro Wilks statistic of $\hat{\lambda}$ ($W = 1.0$ signifies normality, $W < 1.0$ represent skewed distributions).

of the strength of phylogenetic signal across tree sizes and signal strength, and facilitates quantitative comparisons of the relative strength of phylogenetic signal across datasets.

1. Results

Lambda (λ) estimates of phylogenetic signal are inaccurate. Computer simulations reveal that for $\hat{\lambda}$, the distributional expectations of a Bernoulli variable were mostly upheld. First, the mean value of $\hat{\lambda}$ increases as λ increases, but it is negatively-biased (particularly for small tree sizes), and is consistently less than the input λ value across most of its range (Fig. 1 black line). Second, the standard deviation of $\hat{\lambda}$ is largest at intermediate values of λ and smallest at extreme values (Fig. 1 red line), implying that the precision in estimating λ varies across the range of input values. Additionally, standard deviations of $\hat{\lambda}$ are negatively associated with tree size, and for trees of 128 species or less, $\hat{\lambda}$ are quite variable, except for cases when λ is near or equal to 1. Third, the distributions of $\hat{\lambda}$ are not normal across its range, but become increasingly skewed at more extreme values of λ (Fig. 1 blue line). For small tree sizes, it is also clear that distributions are more platykurtic at intermediate values of $\hat{\lambda}$. Taken together these results reveal that $\hat{\lambda}$ inconsistently estimates phylogenetic signal, both across tree sizes and across the range of input values. Additional simulations (Supplemental Information) reveal that incorporating $\hat{\lambda}$ in PGLS anova and regression does not adversely affect the statistical properties of PGLS parameter estimation or model evaluation (type I error, power,

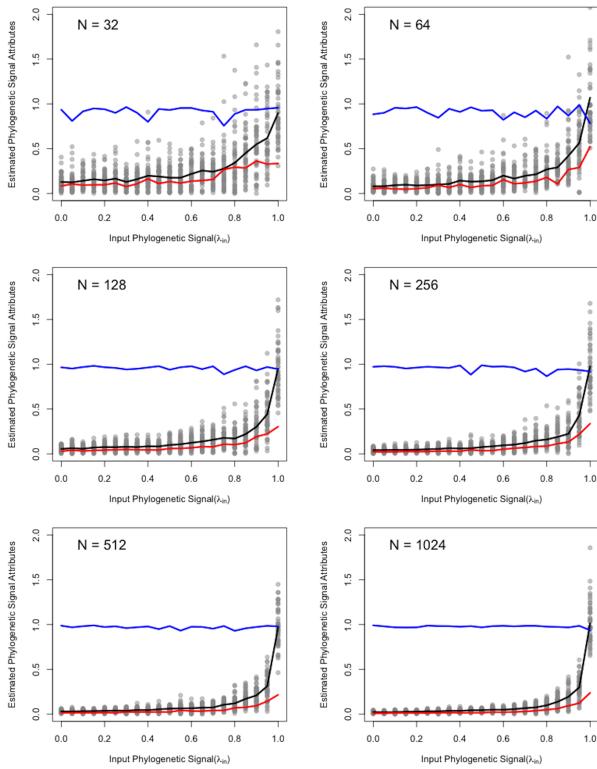


Fig. 2. Response of Blomberg's $\hat{\kappa}$ to increasing strength of Brownian motion. Gray line signifies the 1:1 line where the input value matches the estimate $\hat{\kappa}$. At each input level, the dark black line represents the empirically derived expected value (mean) of $\hat{\kappa}$, the red line is the standard deviation of $\hat{\kappa}$, and the blue line is Shapiro Wilks statistic of $\hat{\kappa}$ ($W = 1.0$ signifies normality, $W < 1.0$ represent skewed distributions).

Effect sizes from κ (Z_κ) better characterize phylogenetic signal. To measure the strength of phylogenetic signal on a common scale, we propose effect sizes (Z-scores) for both λ and κ . Statistically, a standardized effect size may be found as:

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

where θ_{obs} is the observed test statistic, $E(\theta)$ is its expected value under the null hypothesis, and σ_θ is its standard error (31–33). Typically, θ_{obs} and σ_θ are estimated from the data, while $E(\theta)$ is obtained from the distribution of θ derived from parametric theory. However, recent advances in resampling theory (34–37) have shown that $E(\theta)$ and σ_θ may also be obtained from an empirical sampling distribution of θ obtained from permutation procedures.

Formalizing the suggestion of Adams and Collyer (38), an effect size for κ may be found as:

$$Z_\kappa = \frac{\kappa_{obs} - \hat{\mu}_\kappa}{\hat{\sigma}_\kappa}, \quad [2]$$

where κ_{obs} is the observed phylogenetic signal, and $\hat{\mu}_\kappa$ and $\hat{\sigma}_\kappa$ are the mean and standard deviation of the empirical sampling distribution of κ obtained via permutation. The empirical sampling distribution of κ is first transformed via Box-Cox to better adhere to the assumption of normality.

For λ , deriving an effect size is more challenging, as λ does not have a sampling distribution from which the standard error and confidence intervals may be obtained, and estimates from the Hessian matrix from PGLS are unreliable (23). Confidence intervals are therefore generated for the values of λ that intersect the log-likelihood profile for corresponding percentiles of the χ^2 distribution used to compare the putative model to a null model with $\lambda = 0$ [add ref: MLC thinks Boettiger paper?]. Thus, an effect size for λ may be found as:

$$|Z_\lambda| = \sqrt{\chi^2_\lambda} \quad [3]$$

where, $\hat{\lambda}$ is the maximized likelihood value of λ and χ^2_λ is the likelihood ratio statistic for the value.

Here we evaluate the ability of Z_λ and Z_κ to characterize known levels of phylogenetic signal. Both Z_λ and Z_κ are associated with input phylogenetic signal (λ), indicating that both statistics capture the observed signal (Fig. 3). However, effect sizes from $\hat{\lambda}$ made little sense, as they are more strongly associated with tree size than they are with the actual phylogenetic signal in the data (Fig. 3). By contrast, Z_κ is much more consistent across tree sizes, and increases more linearly with increasing levels of phylogenetic signal. Additionally, Z_κ exhibits a much stronger association with phylogenetic signal strength as compared to tree size (Fig. 3), and its standard deviation across input signal is more consistent. This implies that similar levels of precision are found with Z_κ across the range of input values. Thus between the two statistics, Z_κ is a more reliable measure of the strength of phylogenetic signal, and may be used to compare levels of phylogenetic signal across datasets.

A test statistic (\hat{Z}_{12}) allows meaningful comparisons across datasets. To statistically compare the strength

bias in coefficients). Thus, it is reasonable to incorporate $\hat{\lambda}$ in PGLS as a parameter for tuning the degree of phylogenetic signal in the dependent variables during the analysis. However, the statistical properties shown in Fig. 1 demonstrate that λ is unsuitable as an effect size for measuring the strength of phylogenetic signal in data, and thus λ should not be used for comparing phylogenetic signal across datasets.

Kappa (κ) estimates of phylogenetic signal are more reliable. Simulation results for $\hat{\kappa}$ demonstrate that $\hat{\kappa}$ displays better statistical properties. First, as expected, mean values of $\hat{\kappa}$ increase with increasing signal (λ) irrespective of tree size, though the increase does not scale linearly with input levels of phylogenetic signal (Fig. 2 black line). Additionally, the standard deviation of $\hat{\kappa}$ is consistent across tree sizes (Fig. 2 red line), and while it increases with λ , it is always less than the corresponding mean. This finding is perhaps unsurprising, as $\hat{\kappa}$ is lower-bounded by 0, and is never large for small values of λ . Importantly, $\hat{\kappa}$ is normally distributed across the range of input λ , and remains consistent in this pattern regardless of tree size (Fig. 2 blue line). This result differs from those of (18), where the skewing appears to be the result of combining random values generated independently, rather than being a property of κ itself. Overall, these findings reveal that while κ is more reliable as an estimate of phylogenetic signal, the non-linear scaling with input signal implies that it should not be considered an effect size that measures the strength of phylogenetic signal on a common scale for comparison across datasets.

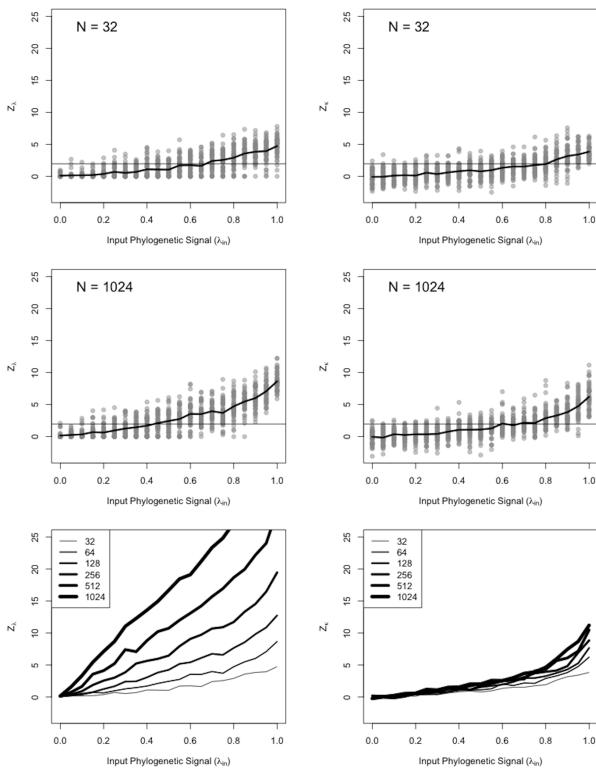


Fig. 3. Response of effect sizes Z_λ and Z_κ to increasing strength of Brownian motion.

of phylogenetic signal across datasets we propose a two-sample test statistic (\hat{Z}_{12}). Based on statistical theory, a two-sample test statistic may be calculated as:

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

where $\kappa_1, \kappa_2, \hat{\mu}_{\kappa_1}, \hat{\mu}_{\kappa_2}, \hat{\sigma}_{\kappa_1}$, and $\hat{\sigma}_{\kappa_2}$ are as defined above for equation 2. The right side of the equation illustrates that if Z_κ has already been calculated for two sampling distributions as in equation 2, the sampling distributions have unit variance for each of the Z_κ statistics. Estimates of significance of \hat{Z}_{12} may be obtained from a standard normal distribution. Typically, \hat{Z}_{12} is considered a two-tailed test, however directional (one-tailed) tests may be specified should the empirical situation require it (35, 37).

To demonstrate the utility of \hat{Z}_{12} , we compared Z_κ for two ecologically-relevant traits in plethodontid salamander (Fig. 4): surface area to volume ratios (SA:V) and relative body width ($\frac{BW}{SVL}$) (39, 40). While both traits contained significant phylogenetic signal, tests based on \hat{Z}_{12} revealed that the degree of phylogenetic signal was significantly stronger in SA:V ($\hat{Z}_{12} = 4.13; P = 0.000036$; Fig. 4). Biologically, this observation may be interpreted by the fact that tropical species – which form a monophyletic group within plethodontids – display greater variation in SA:V, which covaries with disparity in their climatic niches (40). Thus, greater phylogenetic signal in SA:V is to be expected.

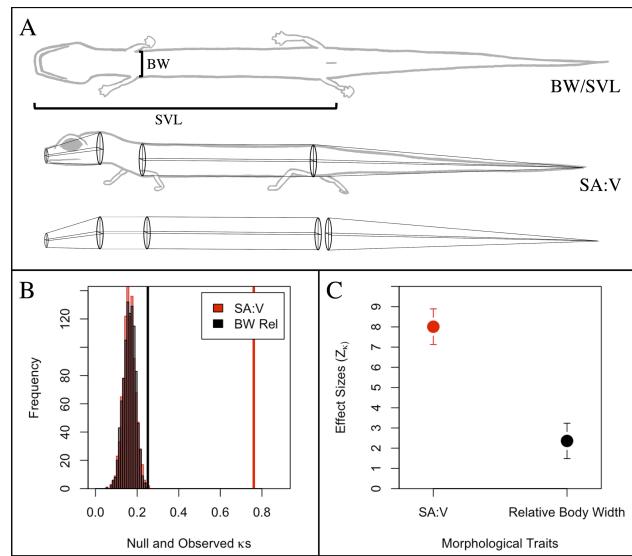


Fig. 4. (A) Linear measures for relative body size, and regions of the body used to estimate surface area to volume (SA:V) ratios. (B) Permutation distributions of phylogenetic signal for SA:V and $\frac{BW}{SVL}$, with observed values shown as vertical bars. (C) Effect sizes (Z_κ) for SA:V and $\frac{BW}{SVL}$, with their 95% confidence intervals (CI not standardized by $\sqrt{(n)}$).

2. Discussion

It is common in comparative evolutionary studies to characterize the phylogenetic signal in phenotypic traits to determine the extent to which shared evolutionary history has generated trait covariation among taxa. However, while numerous analytical approaches may be used to quantify phylogenetic signal (11, 12, 14–16), methods that explicitly measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained underdeveloped. We evaluated the precision of one common measure, Pagel's λ , and explored its efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations, we found that λ behaves as a Bernoulli random variable, with estimates that are increasingly skewed at larger and smaller input levels of phylogenetic signal. Further, the precision of λ in estimating actual levels of phylogenetic signal varies with both tree size (see also ref. (23)) and input levels of phylogenetic signal. From these findings we conclude that λ is not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and should not be used as an effect size for comparing the degree of phylogenetic signal between datasets.

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, values of Z_κ more accurately tracked known changes in the magnitude of phylogenetic signal, as demonstrated by the linear relationship between Z_κ and input signal. Additionally, the precision of Z_κ was more consistent across the range of input levels of phylogenetic signal. Thus, Z_κ is a more reliable measure of the relative strength of phylogenetic signal, and places that effect on a common and comparable scale. We

273 therefore recommend that future studies interested in evaluating
274 the strength of phylogenetic signal incorporate Z_κ as a
275 statistical measure of this effect.

276 Next we proposed a two-sample test (\hat{Z}_{12}), which provides
277 a formal statistical procedure for determining whether the
278 strength of phylogenetic signal is greater in one phenotypic
279 trait as compared to another. Prior studies have summarized
280 patterns of variation in phylogenetic signal across datasets
281 using summary test values, such as κ (12). However, because κ
282 does not scale linearly with input levels of phylogenetic signal
283 (Fig. 2), and its variance increases with increasing strength of
284 phylogenetic signal (18, 20), it should not be considered an
285 effect size that measures the strength of phylogenetic signal
286 on a common scale. By contrast, standardizing κ to Z_κ via
287 equation 2 alleviates these concerns, and facilitates formal sta-
288 tistical comparisons of the strength of signal across datasets.
289 Thus when viewed from this perspective, the approach devel-
290 oped here aligns well with other statistical approaches such as
291 meta-analysis (31, 41, 42), where summary statistics across
292 datasets are converted to standardized effect sizes for subse-
293 quent “higher order” statistical summaries or comparisons. As
294 such, our approach enables evolutionary biologists to quanti-
295 tatively examine the relative strength of phylogenetic signal
296 across a wide range of phenotypic traits, and thus opens the
297 door for future discoveries that inform on how phenotypic
298 diversity accumulates in macroevolutionary time across the
299 tree of life.

300 One important advantage of the approach advocated here
301 is that the resulting effect sizes (Z_κ) are dimensionless, as the
302 units of measurement cancel out during the calculation of Z
303 (43). Thus, Z_κ represents the strength of phylogenetic signal
304 on a common and comparable scale – measured in standard
305 deviations – regardless of the initial units and original scale
306 of the phenotypic variables under investigation. This means
307 that the strength of phylogenetic signal may be compared
308 across datasets for continuous phenotypic traits measured in
309 different units and scale, because those units have been stan-
310 dardized through their conversion to Z_κ . For example, our
311 approach could be utilized to determine whether the strength
312 of phylogenetic signal (say, in response to ecological differ-
313 entiation) is stronger in morphological traits (linear traits:
314 mm), physiological traits (metabolic rate: $\frac{O^2}{min}$), or behavioral
315 traits (aggression: $\frac{\# \text{displays}}{\text{second}}$). In fact, our empirical example
316 provided just such a comparison, as SA:V is represented in
317 mm^{-1} while relative body size is a unitless ratio ($\frac{BW}{SVL}$). Ad-
318 ditionally, our method is capable of comparing the strength
319 of phylogenetic signal in traits of different dimensionality, as
320 estimates of phylogenetic signal using κ have been generalized
321 for multivariate data (16). Furthermore, tests based on \hat{Z}_{12}
322 may be utilized for comparing the strength of phylogenetic
323 signal among datasets containing a different number of species,
324 and even for phenotypes obtained from species in different
325 lineages, because their phylogenetic non-independence and
326 observed variation are taken into account in the generation of
327 the empirical sampling distribution via permutation.

328 This study is not the first to compare λ and κ for their
329 ability as statistics to measure phylogenetic signal. Our re-
330 sults for λ and κ values are consistent with those found in
331 the simulations performed by Münkemüller et al. (18), but
332 that study investigated type I error rates and statistical power,
333 finding that λ performed better in both regards, irrespec-

334 of species number in trees. Although not the central focus of
335 their study, the same tendency for variable λ and consistent
336 κ at intermediate phylogenetic signal strengths was observed
337 (Fig. 2 of ref. (18)). Recent work by Molina-Venegas and
338 Rodríguez (21) found that κ but not λ tended to inflate the
339 estimate of phylogenetic signal, leading to moderate type I
340 and type II biases, if polytomous chronograms were used. Their
341 work more thoroughly addressed previous observations of in-
342 flated κ for incompletely resolved phylogenetic trees (18, 44).
343 An interesting question is whether an inflated κ value leads
344 to an inflated Z_κ or does a tendency of a particular tree to
345 inflate estimates of κ also inflate the values in random permuta-
346 tions of a test, in which case Z_κ is robust to polytomies? We
347 repeated the analyses in Figs. 1 & 2, adjusting trees to have
348 50% collapsed nodes, per the technique of Molina-Venegas and
349 Rodríguez (21), and found results were consistent (Supporting
350 Information). This confirms that any tendency of incompletely
351 resolved trees to inflate κ as a descriptive statistic does not
352 inflate Z_κ as an effect size. Furthermore, because comparison
353 of effect sizes in a test is a comparison of locations of observed
354 values in their sampling distributions, which would shift con-
355 comitantly because of this tendency, the Z_{12} test statistic in
356 equation 4 appears to be robust in spite of unresolved trees.

357 Phylogenetic signal can be thought of as both an attribute
358 to be measured in the data and a parameter that can be tuned
359 to account for the phylogenetic non-independence among ob-
360 servations, for analysis of the data. As such, λ is appealing,
361 as a statistic that potentially fulfills both roles. However,
362 the inability to estimate phylogenetic signal with λ for data
363 simulated with known phylogenetic signal is troublesome, and
364 we recommend evolutionary biologists refrain from viewing it
365 as a statistic to describe the amount of phylogenetic signal in
366 the data. Interestingly, κ – when standardized to an effect size
367 Z_κ – is a better statistic for measuring the amount of phylo-
368 genetic signal in data simulated with respect to known levels
369 of λ . Although λ might be viewed as an important parameter
370 for modifying the conditional estimation of linear model
371 coefficients with respect to phylogeny, it is neither a statistic
372 that has meaningful comparative value as a measure of phylo-
373 genetic signal nor a statistic that lends itself well to reliable
374 calculation of a test statistic. By contrast, κ has been shown
375 here to be a reliable statistic, but only when standardized by
376 the mean and standard deviation of its empirical sampling
377 distribution (i.e., when converted to the effect size, Z_κ). Be-
378 cause one has control over the number of permutations used
379 in analysis, one can be assured with many permutations that
380 the empirical sampling distribution is representative of true
381 probability distributions (10). Given the greater consistency in
382 estimates of Z_κ across tree sizes and input signal, it is difficult
383 to imagine a hypothesis test that can improve equation 4 for
384 efficiently comparing phylogenetic signal for different traits,
385 different trees, or a combination of both.

3. Methods

386 **Simulations.** Simulations were conducted by generating
387 pure-birth phylogenies at each of six different tree sizes
388 ($n = 2^5, 2^6, \dots, 2^{10}$), and with differing levels of phylogenetic
389 signal ($\lambda = 0.0, 0.5, \dots, 1.0$). We generated 100 random trees
390 for each intersection of tree size and λ . For each λ within
391 each tree size, continuous traits were then simulated on each
392 phylogeny under a BM model of evolution. For each set of 100
393

- 394 trees we measured the mean values of $\hat{\lambda}$ and κ , their standard
 395 deviation, and calculated the Shapiro-Wilk W statistic as a
 396 departure from normality (symmetry). For the latter, a value
 397 of 1.0 indicates normally distributed values, while departures
 398 from 1.0 indicate skewness. Simulations were then repeated for
 399 both balanced and pectinate trees, which yielded qualitatively
 400 similar results (see Supporting Information). Trees containing
 401 polytomies, and an evaluation of $\hat{\lambda}$ from models of linear
 402 regression and phylogenetic ANOVA, were also investigated,
 403 and results were qualitatively similar to those reported above
 404 (see Supporting Information).
- 405 **Empirical Data.** Surface area to volume ratios (SA:V)
 406 and relative body width ($\frac{BW}{SVL}$) measures were obtained from
 407 individuals of 305 species, from which species means were
 408 obtained (39, 40). A time-dated molecular phylogeny for the
 409 group (45) was pruned to match the species in the pheno-
 410 typic dataset. The phylogenetic signal in each trait was then
 411 characterized using κ , which was converted to its effect size
 412 (Z_κ) using geomorph 3.3.1 (46, 47), and routines by the au-
 413 thors (to be incorporated in geomorph upon manuscript
 414 acceptance).
- 415 **ACKNOWLEDGMENTS.** We thank E. Glynne and B. Juarez
 416 for comments on early drafts of the manuscript. This work was
 417 supported in part by NSF grant DBI-1902511 (to D.C.A.) and
 418 DBI-1902694 (to M.L.C.).
- 419
1. Felsenstein J (1985) Phylogenies and the comparative method. *American Naturalist* 125(1):1–15.
 2. Harvey PH, Pagel MD (1991) *The comparative method in evolutionary biology* (Oxford University Press, Oxford).
 3. Grafen A (1989) The phylogenetic regression. *Philosophical Transactions of the Royal Society of London B, Biological Sciences* 326:119–157.
 4. Garland TJ, Ives AR (2000) Using the past to predict the present: Confidence intervals for regression equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
 5. Rohlf FJ (2001) Comparative methods for the analysis of continuous variables: Geometric interpretations. *Evolution* 55:2143–2160.
 6. Martins EP, Hansen TF (1997) Phylogenies and the comparative method: A general approach to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist* 149:646–667.
 7. O’Meara BC, Ane C, Sanderson MJ, Wainwright PC (2006) Testing for different rates of continuous trait evolution using likelihood. *Evolution* 60:922–933.
 8. Beaulieu JM, Jhuang DC, Boettiger C, O’Meara BC (2012) Modeling stabilizing selection: Expanding the ornstein-uhlenbeck model of adaptive evolution. *Evolution* 66:2369–2383.
 9. Adams DC (2014) A method for assessing phylogenetic least squares models for shape and other high-dimensional multivariate data. *Evolution* 68:2675–2688.
 10. Adams DC, Collyer ML (2018) Phylogenetic anova: Group-clade aggregation, biological challenges, and a refined permutation procedure. *Evolution* 72(6):1204–1215.
 11. Pagel MD (1999) Inferring the historical patterns of biological evolution. *Nature* 401:877–884.
 12. Blomberg SP, Garland T, Ives AR (2003) Testing for phylogenetic signal in comparative data: Behavioral traits are more labile. *Evolution* 57:717–745.
 13. Revell LJ, Harmon LJ, Collar DC (2008) Phylogenetic signal, evolutionary process, and rate. *Systematic Biology* 57:591–601.
 14. Abouheif E (1999) A method for testing the assumption of phylogenetic independence in comparative data. *Evolutionary Ecology Research* 1:895–909.
 15. Gittleman JL, Kot M (1990) Adaptation: Statistics and a null model for estimating phylogenetic effects. *Systematic Zoology* 39(3):227–241.
 16. Adams DC (2014) A generalized Kappa statistic for estimating phylogenetic signal from shape and other high-dimensional multivariate data. *Systematic Biology* 63:685–697.
 17. Klingenberg CP, Gidaszewski NA (2010) Testing and quantifying phylogenetic signals and homoplasy in morphometric data. *Systematic biology* 59(3):245–261.
 18. Münkemüller T, et al. (2012) How to measure and test phylogenetic signal. *Methods in Ecology and Evolution* 3:743–756.
 19. Pavoine S, Ricotta C (2012) Testing for phylogenetic signal in biological traits: The ubiquity of cross-product statistics. *Evolution: International Journal of Organic Evolution* 67(3):828–840.
 20. Diniz-Filho JAF, Santos T, Rangel TF, Bini LM (2012) A comparison of metrics for estimating phylogenetic signal under alternative evolutionary models. *Genetics and Molecular Biology* 35(3):673–679.
 21. Molina-Venegas R, Rodríguez MA (2017) Revisiting phylogenetic signal; strong or negligible impacts of polytomies and branch length information? *BMC evolutionary biology* 17(1):53.
 22. Revell LJ (2010) Phylogenetic signal and linear regression on species data. *Methods in Ecology and Evolution* 1:319–329.
 23. Boettiger C, Coop G, Ralph P (2012) Is your phylogeny informative? Measuring the power of comparative methods. *Evolution* 67:2240–2251.
 24. Freckleton RP, Harvey PH, Pagel M (2002) Phylogenetic analysis and comparative data: A test and review of evidence. *American Naturalist* 160:712–726.
 25. Cooper N, Jetz W, Freckleton RP (2010) Phylogenetic comparative approaches for studying niche conservatism. *Journal of Evolutionary Biology* 23(12):2529–2539.
 26. Bose R, Ramesh BR, Pélassier R, Munoz F (2019) Phylogenetic diversity in the western ghats biodiversity hotspot reflects environmental filtering and past niche diversification of trees. *Journal of Biogeography* 46(1):145–157.
 27. Vandeloek F, et al. (2019) Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124(2):269–279.
 28. De Meester G, Huyghe K, Van Damme R (2019) Brain size, ecology and sociality: A reptilian perspective. *Biological Journal of the Linnean Society* 126(3):381–391.
 29. Pintanel P, Tejedo M, Ron SR, Llorente GA, Merino-Viteri A (2019) Elevational and microclimatic drivers of thermal tolerance in andean pristimantis frogs. *Journal of Biogeography* 46(8):1664–1675.
 30. Su G, Villéger S, Brosse S (2019) Morphological diversity of freshwater fishes differs between realms, but morphologically extreme species are widespread. *Global ecology and biogeography* 28(2):211–221.
 31. Glass GV (1976) Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
 32. Cohen J (1988) *Statistical power analysis for the behavioral sciences* (Routledge).
 33. Rosenthal R (1994) The handbook of research synthesis. ed Cooper LV H Hedges (Russell Sage Foundation), pp 231–244.
 34. Collyer ML, Sekora DJ, Adams DC (2015) A method for analysis of phenotypic change for phenotypes described by high-dimensional data. *Heredity* 115:357–365.
 35. Adams DC, Collyer ML (2016) On the comparison of the strength of morphological integration across morphometric datasets. *Evolution* 70:2623–2631.
 36. Collyer ML, Adams DC (2018) RRPP: An r package for fitting linear models to high-dimensional data using residual randomization. *Methods in Ecology and Evolution* 9:1772–1779.
 37. Adams DC, Collyer ML (2019) Comparing the strength of modular signal, and evaluating alternative modular hypotheses, using covariance ratio effect sizes with morphometric data. *Evolution* 73(12):2352–2367.
 38. Adams DC, Collyer ML (2019) Phylogenetic comparative methods and the evolution of multivariate phenotypes. *Annual Review of Ecology, Evolution, and Systematics* 50:405–425.
 39. Baken EK, Adams DC (2019) Macroevolution of arboreality in salamanders. *Ecology and Evolution* 9(12):7005–7016.
 40. Baken EK, Mellenthin LE, Adams DC (2020) Macroevolution of desiccation-related morphology in plethodontid salamanders as inferred from a novel surface area to volume ratio estimation approach. *Evolution* 74:476–486.

- 537 41. Hedges L. V., Olkin I (1985) *Statistical methods for meta-*
538 *analysis* (Elsevier).
- 539 42. Arnqvist G., Wooster D (1995) Meta-analysis: Synthesizing
540 research findings in ecology and evolution. *Trends in Ecology and*
541 *Evolution* 10:236–240.
- 542 43. Sokal R. R., Rohlf FJ (2012) *Biometry* (W.H. Freeman &
543 Co., San Francisco). 4th Ed.
- 544 44. Davies TJ, Kraft NJ, Salamin N, Wolkovich EM (2012)
545 Incompletely resolved phylogenetic trees inflate estimates of phylo-
546 genetic conservatism. *Ecology* 93(2):242–247.
- 547 45. Bonett RM, Blair AL (2017) Evidence for complex life cycle
548 constraints on salamander body form diversification. *Proceedings*
549 *of the National Academy of Sciences, USA* 114:9936–9941.
- 550 46. Adams DC, Otárola-Castillo E (2013) Geomorph: An r
551 package for the collection and analysis of geometric morphometric
552 shape data. *Methods in Ecology and Evolution* 4:393–399.
- 553 47. Adams DC, Collyer ML, Kallontzopoulou A (2020) Geo-
554 morph: Software for geometric morphometric analyses. R pack-
555 age version 3.3.1. Available at: <https://cran.r-project.org/package=geomorph>.
- 556

DRAFT