

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 27, 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 of
6 input values. We find that λ behaves as a Bernoulli random variable,
7 where estimates are increasingly skewed at larger and smaller input
8 levels of phylogenetic signal. Further, the precision of λ varies with
9 input signal. Another measure, Blomberg's K , is more consistent
10 across a range of tree sizes, and exhibits a positive relationship with
11 input levels of phylogenetic signal. However, that relationship is de-
12 cidedly nonlinear. Thus, neither λ nor K are suitable as effect sizes
13 for measuring the strength of phylogenetic signal, and comparing
14 that signal across datasets. As an alternative, we propose a stan-
15 dardized effect size based on K , (Z_K), which measures the strength
16 of phylogenetic signal more reliably than does λ , and places that sig-
17 nal on a common scale for statistical comparison. We develop tests
18 based on Z_K to provide a mechanism for formally comparing the
19 strength of phylogenetic signal across datasets, in much the same
20 manner as effect sizes may be used to summarize patterns in quan-
21 titative meta-analysis. Our approach extends the phylogenetic com-
22 parative toolkit to address hypotheses that compare the strength of
23 phylogenetic signal between various phenotypic traits, even when
24 those traits are found in different evolutionary lineages or have dif-
25 ferent units or scales.

phylogenetic signal | macroevolution | lambda | kappa

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

15 Several analytical tools have been developed to quantify
16 phylogenetic signal in phenotypic datasets (11, 12, 14–17),
17 and their statistical properties – namely type I error rates and
18 statistical power – have been investigated to determine under
19 what conditions phylogenetic signal can be detected (13, 16,
20 18–23). One of the most widely used methods for character-
21 izing phylogenetic signal is Pagel's λ (11), which transforms
22 the lengths of the internal branches of the phylogeny to im-
23 prove the fit of data to the phylogeny via maximum likelihood

(11, 24). When incorporated in PGLS, λ serves as a tuning
24 parameter which is optimized via log-likelihood profiling while
25 evaluating the covariation between the dependent and indepen-
26 dent 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 is Blomberg's K (12), which
32 characterizes phylogenetic signal as the ratio of observed trait
33 variation to the amount of variation expected under Brownian
34 motion. Blomberg's K can be treated as a test statistic by
35 employing a permutation test to generate its sampling distribu-
36 tion (12, 16) for determining whether significant phylogenetic
37 signal is present in data. Both λ and K seem intuitive to inter-
38 pret, as a value of 0 for both corresponds to no phylogenetic
39 signal, while a value of 1 corresponds to the amount of phylo-
40 genetic signal expected under Brownian motion. Thus, it is
41 tempting to regard both λ and K as descriptive statistics that
42 measure the relative strength of phylogenetic signal, providing
43 an estimate of its magnitude for comparison.

44 The appeal of Pagel's λ and Blomberg's K as descriptive
45 statistics is that they provide a basis for interpreting “weak”
46 versus “strong” phylogenetic signal; i.e., small versus large
47

Significance Statement

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

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

The authors declare no conflict of interest.

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>

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

values of $\hat{\lambda}$ or K , 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 $\hat{\lambda}$ is a random variable that follows the distribution of a Bernoulli probability parameter (31); 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 $\hat{\lambda}$ would be more platykurtic at intermediate values of λ , becoming more leptokurtic toward 0 and 1 (31). 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 K , which is positively unbounded, we might expect that for any λ used to generate data, estimates of K might be a random variable that follows a normal distribution, with values distributed symmetrically (31). 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 K 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 K , 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 K , 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 K are more consistent across tree sizes, and are normally distributed across the range of input levels of λ , making K a more reliable statistic. We then propose an effect size based on K , (Z_K), which provides consistent estimates of the strength of phylogenetic signal across tree sizes and

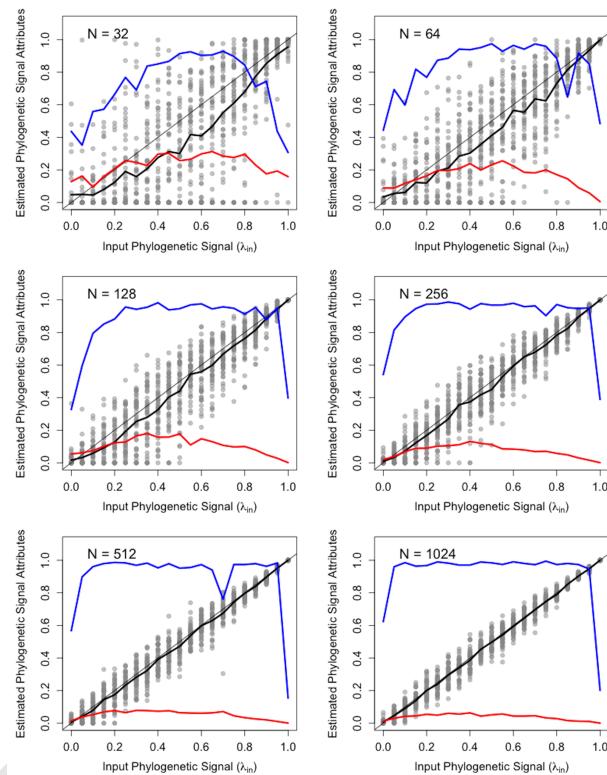


Fig. 1. Response of Pagel’s λ to increasing strength of Brownian motion. Gray line signifies the 1:1 line where the input value matches the estimate $\hat{\lambda}$. At each input level, the dark black line represents the empirically derived expected value (mean) of $\hat{\lambda}$, the red line is the 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).

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. Second, the precision in estimating λ varies across the range of input values, as the standard deviation of $\hat{\lambda}$ is largest at intermediate values of λ and smallest at extreme values (Fig. 1 red line). Third, the distributions of $\hat{\lambda}$ tend toward normal distributions at intermediate levels of λ 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 $\hat{\lambda}$. However, the mean value of $\hat{\lambda}$ is negatively-biased (particularly for small tree sizes but also consistently across most of its range; Fig. 1 black line) and standard deviations of $\hat{\lambda}$ are negatively associated with tree size. For trees of 128 species or less, $\hat{\lambda}$ are quite variable, except for cases when λ is near or equal to 1. Taken together these results reveal that $\hat{\lambda}$ is a biased statistic that inconsistently estimates phylogenetic signal, both across tree sizes and across the range of input values. Additional simulations (Supporting 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, bias in coefficients). Thus, it is reasonable to

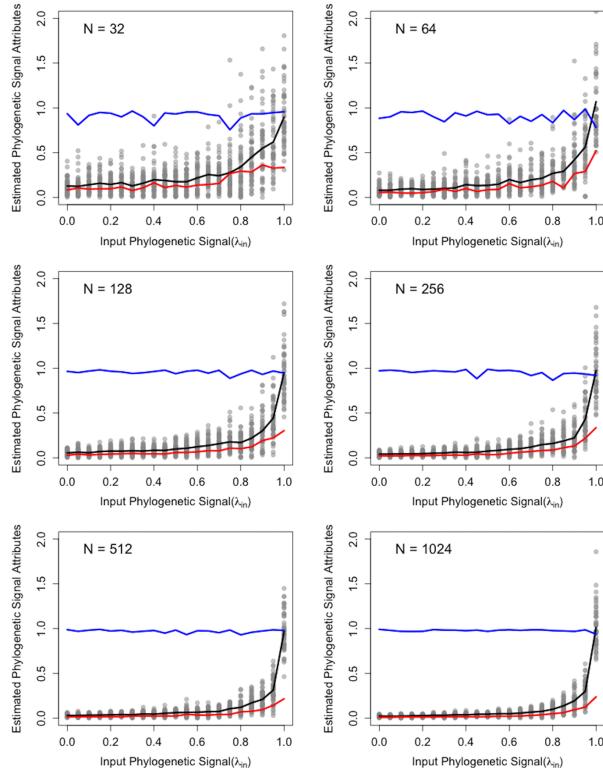


Fig. 2. Response of Blomberg's K^* to increasing strength of Brownian motion. At each input level, the black line represents the empirically derived expected value (mean) of K^* , the red line is the standard deviation of K^* , and the blue line is Shapiro Wilks statistic of K^* ($W = 1.0$ signifies normality, $W < 1.0$ represent skewed distributions).

for both λ and K . 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 (32–34). 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 (35–38) have shown that $E(\theta)$ and σ_θ may also be obtained from an empirical sampling distribution of θ simulated from permutation procedures.

Formalizing the suggestion of Adams and Collyer (39), an effect size for K may be found as:

$$Z_K = \frac{K_{obs} - \hat{\mu}_K}{\hat{\sigma}_K}, \quad [2]$$

where K_{obs} is the observed phylogenetic signal, and $\hat{\mu}_K$ and $\hat{\sigma}_K$ are the mean and standard deviation of the empirical sampling distribution of K obtained via permutation. The empirical sampling distribution of K can be first transformed via a Box-Cox transformation 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 $\lambda = 0$ (40). 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.

Simulations reveal that both Z_λ and Z_K 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_K is much more consistent across tree sizes, and increases more linearly with increasing levels of phylogenetic signal. Additionally, Z_K exhibits a much stronger association with phylogenetic signal strength as compared to tree size (Fig. 3), and its standard deviation is more consistent, implying similar levels of precision across the range of input signal (Supporting Information). Thus, Z_K 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 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:

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 (K) estimates of phylogenetic signal are more reliable. Simulation results demonstrate that K displays better statistical properties. First, as expected, mean values of K increase with increasing signal (λ) irrespective of tree size, albeit non-linearly (Fig. 2 black line). Second, the standard deviation of K is consistent across tree sizes (Fig. 2 red line), and while it increases with λ , it is always less than the mean (low coefficient of variation). This finding is perhaps unsurprising, as K is lower-bounded by 0, and is never large for small values of λ . Importantly, K is normally distributed across the range of input λ ; a consistent pattern regardless of tree size (Fig. 2 blue line). This differs from results of (18), where the skewing appears to be due to combining random values generated independently, rather than being a property of K itself. Overall, these findings reveal that while K 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.

Effect sizes from K (Z_K) better characterize phylogenetic signal. To measure the strength of phylogenetic signal on a common scale, we propose effect sizes (Z-scores)

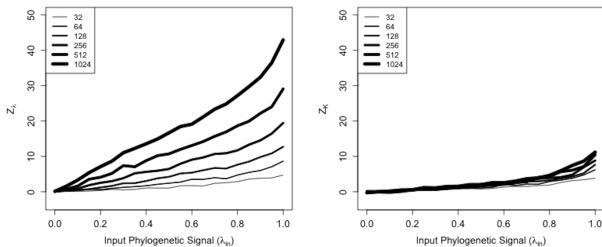


Fig. 3. Response of effect sizes Z_λ and Z_K to increasing strength of Brownian motion. Means from simulation runs are shown for comparative ease. Individual values from each simulation run are available in Supporting Information.

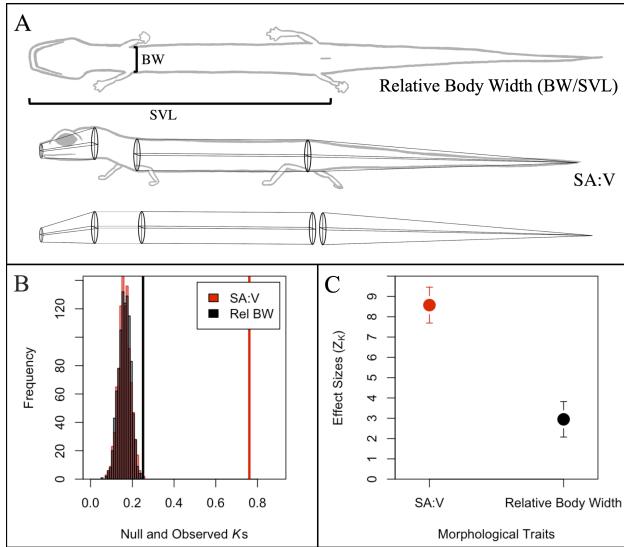


Fig. 4. (A) Linear measures for relative body size, and regions of the body used to estimate surface area to volume (SA:V) ratios. (B) Permutation distributions of phylogenetic signal for SA:V and $\frac{BW}{SVL}$, with observed values shown as vertical bars. (C) Effect sizes (Z_K) for SA:V and $\frac{BW}{SVL}$,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

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

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. 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 (36, 38).

To demonstrate the utility of \hat{Z}_{12} , we compared Z_K for two ecologically-relevant traits in plethodontid salamanders (Fig. 4): surface area to volume ratios (SA:V) and relative body width ($\frac{BW}{SVL}$) (41, 42). 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 the tropical species – which form a monophyletic group within plethodontids – display greater variation in SA:V, which covaries with disparity in their climatic niches (42). Thus, greater phylogenetic signal in SA:V is to be expected.

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 K , and found that Z_K was a better estimate of the strength of phylogenetic signal in phenotypic data. First, values of Z_K more accurately tracked known changes in the magnitude of phylogenetic signal, as demonstrated by the near linear relationship between Z_K and input signal. Additionally, the precision of Z_K was more consistent across the range of input levels of phylogenetic signal (Fig S1; Supporting Information). Thus, Z_K is a more reliable measure of the relative strength of phylogenetic signal, and places that effect on a common and comparable scale. We therefore recommend that future studies interested in evaluating the strength of phylogenetic signal incorporate Z_K as a statistical measure of this effect.

Next we proposed a two-sample test (\hat{Z}_{12}), which provides a formal statistical procedure for determining whether the strength of phylogenetic signal is greater in one phenotypic trait as compared to another. Prior studies have summarized patterns of variation in phylogenetic signal across datasets using summary test values, such as K (12). However, because K does not scale linearly with input levels of phylogenetic signal (Fig. 2), and its variance increases with increasing strength of phylogenetic signal (18, 20), it should not be considered an effect size that measures the strength of phylogenetic signal on a common scale. By contrast, standardizing K to Z_K 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 (32, 43, 44), 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

291 diversity accumulates in macroevolutionary time across the
292 tree of life.

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

321 This study is not the first to compare λ and K for their
322 ability as statistics to measure phylogenetic signal. Our re-
323 sults for λ and K values are consistent with those found in
324 the simulations performed by Münkemüller et al. (18), but
325 that study investigated type I error rates and statistical power,
326 finding that λ performed better in both regards, irrespective
327 of species number in trees. Although not the central focus of
328 their study, the same tendency for variable λ and consistent
329 K at intermediate phylogenetic signal strengths was observed
330 (Fig. 2 of ref. (18)). Recent work by Molina-Venegas and
331 Rodríguez (21) found that K but not λ tended to inflate the
332 estimate of phylogenetic signal, leading to moderate type I and
333 type II biases, if polytomic chronograms were used. Their work
334 more thoroughly addressed previous observations of inflated
335 K for incompletely resolved phylogenetic trees (18, 46). An
336 interesting question is whether an inflated K value leads to an
337 inflated Z_K or does a tendency of a particular tree to inflate
338 estimates of K also inflate the values in random permutations
339 of a test, in which case Z_K is robust to polytomies? We re-
340 peated the analyses in Figs. 1 & 2, adjusting trees to have
341 20% collapsed nodes, per the technique of Molina-Venegas and
342 Rodríguez (21), and found results were consistent (Supporting
343 Information). This confirms that any tendency of incompletely
344 resolved trees to inflate K as a descriptive statistic does not
345 inflate Z_K as an effect size. Furthermore, because comparison
346 of effect sizes in a test is a comparison of locations of observed
347 values in their sampling distributions, which would shift con-
348 comitantly because of this tendency, the Z_{12} test statistic in
349 equation 4 appears to be robust in spite of unresolved trees.

350 Phylogenetic signal can be thought of as both an attribute
351 to be measured in the data and a parameter that can be tuned

352 to account for the phylogenetic non-independence among ob-
353 servations, for analysis of the data. As such, λ is appealing,
354 as a statistic that potentially fulfills both roles. However,
355 the inability to estimate phylogenetic signal with λ for data
356 simulated with known phylogenetic signal is troublesome, and
357 we recommend evolutionary biologists refrain from viewing it
358 as a statistic to describe the amount of phylogenetic signal in
359 the data. Interestingly, K – when standardized to an effect
360 size Z_K – is a better statistic for measuring the amount of
361 phylogenetic signal in data simulated with respect to known
362 levels of λ . Although λ might be viewed as an important
363 parameter for modifying the the conditional estimation of
364 linear model coefficients with respect to phylogeny, it is nei-
365 ther a statistic that has meaningful comparative value as a
366 measure of phylogenetic signal nor a statistic that lends itself
367 well to reliable calculation of a test statistic. By contrast,
368 K has been shown here to be a reliable statistic, but only
369 when standardized by the mean and standard deviation of its
370 empirical sampling distribution (i.e., when converted to the
371 effect size, Z_K). Because one has control over the number
372 of permutations used in analysis, one can be assured with
373 many permutations that the empirical sampling distribution is
374 representative of true probability distributions (10). Given the
375 greater consistency in estimates of Z_K across tree sizes and
376 input signal, it is difficult to imagine a hypothesis test that
377 can improve equation 4 for efficiently comparing phylogenetic
378 signal for different traits, different trees, or a combination of
379 both.

3. Methods

380 **Simulations.** Simulations were conducted by generating
381 pure-birth phylogenies at each of six different tree sizes
382 ($n = 2^5, 2^6, \dots, 2^{10}$), and with differing levels of phylogenetic
383 signal ($\lambda = 0.0, 0.5, \dots, 1.0$). We generated 50 random trees
384 for each intersection of tree size and λ . For each λ within
385 each tree size, continuous traits were then simulated on each
386 phylogeny under a BM model of evolution. For each set of 50
387 trees we measured the mean values of $\hat{\lambda}$ and K , their standard
388 deviation, and calculated the Shapiro-Wilk W statistic as a
389 departure from normality (symmetry). For the latter, a value
390 of 1.0 indicates normally distributed values, while departures
391 from 1.0 indicate skewness. Simulations were then repeated for
392 both balanced and pectinate trees, which yielded qualitatively
393 similar results (see Supporting Information). Trees contain-
394 ing polytomies, and an evaluation of $\hat{\lambda}$ from models of linear
395 regression and phylogenetic ANOVA, were also investigated,
396 and results were qualitatively similar to those reported above
397 (see Supporting Information).

398 **Empirical Data.** Surface area to volume ratios (SA:V)
399 and relative body width ($\frac{BW}{SVL}$) measures were obtained from
400 individuals of 305 species, from which species means were
401 obtained (41, 42). A time-dated molecular phylogeny for the
402 group (47) was pruned to match the species in the phenotypic
403 dataset. The phylogenetic signal in each trait was then char-
404 acterized using K , which was converted to its effect size (Z_K)
405 using geomorph 3.3.1 (48, 49), and routines by the authors (**to**
406 **be incorporated in geomorph upon acceptance**).

407 **ACKNOWLEDGMENTS.** We thank E. Glynne and B. Juarez
408 for comments on early drafts of the manuscript. This work was
409 supported in part by NSF grant DBI-1902511 (to D.C.A.) and
410

- 411 DBI-1902694 (to M.L.C.).
 412
 413 1. Felsenstein J (1985) Phylogenies and the comparative method.
 414 *American Naturalist* 125(1):1–15.
 415 2. Harvey PH, Pagel MD (1991) *The comparative method in*
 416 *evolutionary biology* (Oxford University Press, Oxford).
 417 3. Grafen A (1989) The phylogenetic regression. *Philosophical*
 418 *Transactions of the Royal Society of London B, Biological Sciences*
 419 326:119–157.
 420 4. Garland TJ, Ives AR (2000) Using the past to predict the
 421 present: Confidence intervals for regression equations in phylogenetic
 422 comparative methods. *American Naturalist* 155:346–364.
 423 5. Rohlf FJ (2001) Comparative methods for the analysis of
 424 continuous variables: Geometric interpretations. *Evolution* 55:2143–
 425 2160.
 426 6. Martins EP, Hansen TF (1997) Phylogenies and the comparative
 427 method: A general approach to incorporating phylogenetic
 428 information into the analysis of interspecific data. *American Natu-*
 429 *ralist* 149:646–667.
 430 7. O’Meara BC, Ane C, Sanderson MJ, Wainwright PC (2006)
 431 Testing for different rates of continuous trait evolution using likeli-
 432 hood. *Evolution* 60:922–933.
 433 8. Beaulieu JM, Jhuang DC, Boettiger C, O’Meara BC (2012)
 434 Modeling stabilizing selection: Expanding the ornstein-uhlenbeck
 435 model of adaptive evolution. *Evolution* 66:2369–2383.
 436 9. Adams DC (2014) A method for assessing phylogenetic least
 437 squares models for shape and other high-dimensional multivariate
 438 data. *Evolution* 68:2675–2688.
 439 10. Adams DC, Collyer ML (2018) Phylogenetic anova: Group-
 440 clade aggregation, biological challenges, and a refined permutation
 441 procedure. *Evolution* 72(6):1204–1215.
 442 11. Pagel MD (1999) Inferring the historical patterns of biological
 443 evolution. *Nature* 401:877–884.
 444 12. Blomberg SP, Garland T, Ives AR (2003) Testing for phy-
 445 logenetic signal in comparative data: Behavioral traits are more
 446 labile. *Evolution* 57:717–745.
 447 13. Revell LJ, Harmon LJ, Collar DC (2008) Phylogenetic signal,
 448 evolutionary process, and rate. *Systematic Biology* 57:591–601.
 449 14. Abouheif E (1999) A method for testing the assumption
 450 of phylogenetic independence in comparative data. *Evolutionary*
 451 *Ecology Research* 1:895–909.
 452 15. Gittleman JL, Kot M (1990) Adaptation: Statistics and a
 453 null model for estimating phylogenetic effects. *Systematic Zoology*
 454 39(3):227–241.
 455 16. Adams DC (2014) A generalized Kappa statistic for esti-
 456 mating phylogenetic signal from shape and other high-dimensional
 457 dultivariate data. *Systematic Biology* 63:685–697.
 458 17. Klingenberg CP, Gidaszewski NA (2010) Testing and quan-
 459 tifying phylogenetic signals and homoplasy in morphometric data.
 460 *Systematic biology* 59(3):245–261.
 461 18. Münkemüller T, et al. (2012) How to measure and test
 462 phylogenetic signal. *Methods in Ecology and Evolution* 3:743–756.
 463 19. Pavoine S, Ricotta C (2012) Testing for phylogenetic signal in
 464 biological traits: The ubiquity of cross-product statistics. *Evolution:*
 465 *International Journal of Organic Evolution* 67(3):828–840.
 466 20. Diniz-Filho JAF, Santos T, Rangel TF, Bini LM (2012)
 467 A comparison of metrics for estimating phylogenetic signal under
 468 alternative evolutionary models. *Genetics and Molecular Biology*
 469 35(3):673–679.
 470 21. Molina-Venegas R, Rodríguez MA (2017) Revisiting phylo-
 471 genetic signal; strong or negligible impacts of polytomies and branch
 472 length information? *BMC evolutionary biology* 17(1):53.
 473 22. Revell LJ (2010) Phylogenetic signal and linear regression
 474 on species data. *Methods in Ecology and Evolution* 1:319–329.
 475 23. Boettiger C, Coop G, Ralph P (2012) Is your phylogeny infor-
 476 mative? Measuring the power of comparative methods. *Evolution*
 477 67:2240–2251.
 478 24. Freckleton RP, Harvey PH, Pagel M (2002) Phylogenetic
 479 analysis and comparative data: A test and review of evidence.
 480 *American Naturalist* 160:712–726.
 481 25. Cooper N, Jetz W, Freckleton RP (2010) Phylogenetic
 482 comparative approaches for studying niche conservatism. *Journal of*
 483 *Evolutionary Biology* 23(12):2529–2539.
 484 26. Bose R, Ramesh BR, Pélišsier R, Munoz F (2019) Phylo-
 485 genetic diversity in the western ghats biodiversity hotspot reflects
 486 environmental filtering and past niche diversification of trees. *Jour-*
 487 *nal of Biogeography* 46(1):145–157.
 488 27. Vandeloek F, et al. (2019) Nectar traits differ between polli-
 489 nation syndromes in balsaminaceae. *Annals of Botany* 124(2):269–
 490 279.
 491 28. De Meester G, Huyghe K, Van Damme R (2019) Brain size,
 492 ecology and sociality: A reptilian perspective. *Biological Journal of*
 493 *the Linnean Society* 126(3):381–391.
 494 29. Pintanel P, Tejedo M, Ron SR, Llorente GA, Merino-Viteri
 495 A (2019) Elevational and microclimatic drivers of thermal tolerance
 496 in andean pristimantis frogs. *Journal of Biogeography* 46(8):1664–
 497 1675.
 498 30. Su G, Villéger S, Brosse S (2019) Morphological diversity
 499 of freshwater fishes differs between realms, but morphologically
 500 extreme species are widespread. *Global ecology and biogeography*
 501 28(2):211–221.
 502 31. Forbes C, Evans M, Hastings N, Peacock B (2011) *Statistical*
 503 *distributions* (John Wiley & Sons).
 504 32. Glass GV (1976) Primary, secondary, and meta-analysis of
 505 research. *Educational Researcher* 5:3–8.
 506 33. Cohen J (1988) *Statistical power analysis for the behavioral*
 507 *sciences* (Routledge).
 508 34. Rosenthal R (1994) The handbook of research synthesis. ed
 509 Cooper LV H Hedges (Russell Sage Foundation), pp 231–244.
 510 35. Collyer ML, Sekora DJ, Adams DC (2015) A method for
 511 analysis of phenotypic change for phenotypes described by high-
 512 dimensional data. *Heredity* 115:357–365.
 513 36. Adams DC, Collyer ML (2016) On the comparison of the
 514 strength of morphological integration across morphometric datasets.
 515 *Evolution* 70:2623–2631.
 516 37. Collyer ML, Adams DC (2018) RRPP: An r package for
 517 fitting linear models to high-dimensional data using residual ran-
 518 domization. *Methods in Ecology and Evolution* 9:1772–1779.
 519 38. Adams DC, Collyer ML (2019) Comparing the strength of
 520 modular signal, and evaluating alternative modular hypotheses, us-
 521 ing covariance ratio effect sizes with morphometric data. *Evolution*
 522 73(12):2352–2367.
 523 39. Adams DC, Collyer ML (2019) Phylogenetic comparative
 524 methods and the evolution of multivariate phenotypes. *Annual*
 525 *Review of Ecology, Evolution, and Systematics* 50:405–425.
 526 40. Orme D, et al. (2013) CAPER: Comparative analyses of
 527 phylogenetics and evolution in r. *Methods in Ecology and Evolution*
 528 3:145–151.
 529 41. Baken EK, Adams DC (2019) Macroevolution of arboreality
 530 in salamanders. *Ecology and Evolution* 9(12):7005–7016.
 531 42. Baken EK, Mellenthin LE, Adams DC (2020) Macroevolu-
 532 tion of desiccation-related morphology in plethodontid salamanders
 533 as inferred from a novel surface area to volume ratio estimation
 534 approach. *Evolution* 74:476–486.
 535 43. Hedges L. V., Olkin I (1985) *Statistical methods for meta-*
 536 *analysis* (Elsevier).
 537 44. Arnqvist G., Wooster D (1995) Meta-analysis: Synthesizing
 538 research findings in ecology and evolution. *Trends in Ecology and*
 539 *Evolution* 10:236–240.
 540 45. Sokal R. R., Rohlf FJ (2012) *Biometry* (W.H. Freeman &
 541 Co., San Francisco). 4th Ed.
 542 46. Davies TJ, Kraft NJ, Salamin N, Wolkovich EM (2012)
 543 Incompletely resolved phylogenetic trees inflate estimates of phylo-
 544 genetic conservatism. *Ecology* 93(2):242–247.
 545 47. Bonett RM, Blair AL (2017) Evidence for complex life cycle
 546 constraints on salamander body form diversification. *Proceedings*
 547 *of the National Academy of Sciences, USA* 114:9936–9941.
 548 48. Adams DC, Otárola-Castillo E (2013) Geomorph: An r
 549 package for the collection and analysis of geometric morphometric
 550 shape data. *Methods in Ecology and Evolution* 4:393–399.
 551 49. Adams DC, Collyer ML, Kaliontzopoulou A (2020) Geo-
 552 morph: Software for geometric morphometric analyses. R pack-
 553 age version 3.3.1. Available at: <https://cran.r-project.org/package=geomorph>.
 554