

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

Dean C. Adams^{a,1}, 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 July 16, 2020

Macroevolutionary studies frequently characterize the phylogenetic signal in phenotypes, and wish to compare the strength of that signal across traits. However, analytical tools for such comparisons have largely remained underdeveloped. Here we evaluate the efficacy of one commonly used parameter (Pagel's λ) to estimate the strength of phylogenetic signal in phenotypic traits, and evaluate the degree to which λ correctly identifies known levels of phylogenetic signal. We find that λ behaves as a Bernoulli random variable, and that estimates are increasingly skewed at larger and smaller input levels of phylogenetic signal. Further, the precision of λ in estimating actual levels of phylogenetic signal is often inaccurate, and biological interpretations of the strength of phylogenetic signal based on λ are therefore compromised. As an alternative, we propose a standardized effect size based on κ , (Z_κ), which measures the strength of phylogenetic signal more reliably than does λ , and places that signal on a common scale for statistical comparison. We develop tests based on Z_κ to provide a mechanism for formally comparing the strength of phylogenetic signal across datasets, in much the same manner as effect sizes may be used to summarize patterns in quantitative meta-analysis. Our approach extends the phylogenetic comparative toolkit to 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.

phylogenetic signal | macroevolution | lambda | kappa

Investigating macroevolutionary patterns of trait variation requires a phylogenetic perspective, because the shared ancestry among species violates the assumption of independence among trait values that is common for statistical tests (1, 2). Accounting for this evolutionary non-independence is the purview of *phylogenetic comparative methods* (PCMs): a suite of analytical tools that condition trends in the data on the phylogenetic relatedness of observations (3–12). These methods are predicated on the notion that phylogenetic signal – the tendency for closely related species to display similar trait values – is present in cross-species datasets (1, 13, 14). Indeed, under numerous evolutionary models, phylogenetic signal is to be expected, as stochastic character change along the hierarchical structure of the tree of life generates trait covariation among related taxa (1, 14, 15).

Several analytical tools have been developed to quantify phylogenetic signal in phenotypic datasets (13, 14, 16–19), and their statistical properties – namely type I error rates and statistical power – have been investigated to determine under what conditions phylogenetic signal can be detected (15, 18, 20–25). One of the most widely used methods for characterizing phylogenetic signal is Pagel's λ (13), which transforms the lengths of the internal branches of the phylogeny to improve the fit of data to the phylogeny via maximum likelihood

(13, 26). When incorporated in PGLS, λ serves as a tuning parameter which is optimized via log-likelihood profiling while evaluating the covariation between the dependent and independent variables, given the phylogeny (13, 26). To infer whether phylogenetic signal differs from no signal or a Brownian motion model of evolutionary divergence, the observed model fit using λ may be statistically compared to that using $\lambda = 0$ or $\lambda = 1$ via likelihood ratio tests (26–28) or confidence limits (29). Another widely used measure is Blomberg's κ (14), which measures the ratio of observed trait variation to the amount expected under Brownian motion. Blomberg's κ can be treated as a test statistic by employing a permutation test to generate its sampling distribution (14, 18) for determining whether significant phylogenetic signal is present in data. Both 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 Brownian motion. Thus, it is tempting to regard both λ and κ as descriptive statistics that measure the relative strength of phylogenetic signal, providing an estimate of its magnitude for comparison.

The appeal of Pagel's λ and Blomberg's κ as descriptive statistics is that they provide a basis for interpreting “weak” versus “strong” phylogenetic signal; i.e., small versus large values of λ or κ , respectively, in a comparative sense (30–

Significance Statement

Evolutionary biologists wish to quantify and compare the strength of phylogenetic signal across traits, but analytical tools for these comparisons are generally lacking. Here we develop a standardized effect size based on κ , (Z_κ), which measures the strength of phylogenetic signal on a common statistical scale, and provides a mechanism for formally comparing the strength of phylogenetic signal across datasets. Additionally, we find that a commonly used parameter (Pagel's λ) is unsuitable for this purpose. Our 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>

¹ To whom correspondence should be addressed. E-mail: dcadams@iastate.edu

32). Nonetheless, an important question that has yet to be considered is 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 treesize; 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 (20) 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. (20)). 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 (20), interpretation of these patterns is challenging.

By contrast, Blomberg’s κ , which is positively unbounded, might be expected to follow a normal distribution (add ref). Thus we might expect that κ is symmetrically distributed across different strengths of phylogenetic signal, for any λ used to generate data. This attribute seemed less reasonable based on the simulations performed by Münkemüller et al. (20), 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. Furthermore, these are statistics that appear to have expected values that vary with tree size (20), making comparisons across studies challenging. Therefore, transformation of these statistics into Z-scores in the same simulation experiments 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 clades.

Here we used 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 of the strength of phylogenetic signal across tree sizes and

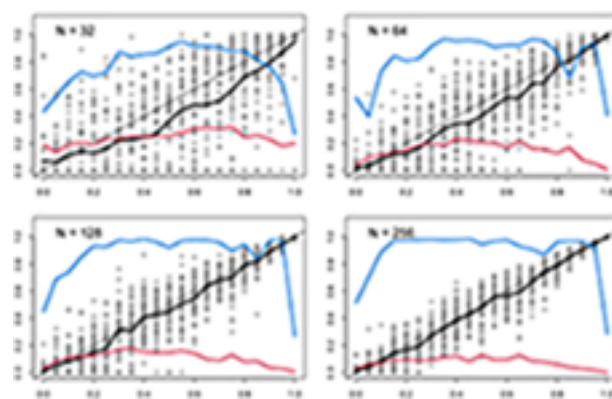


Fig. 1. Patterns in λ .

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 revealed that for $\hat{\lambda}$, the distributional expectations of a Bernoulli variable were mostly upheld. First, the mean value of $\hat{\lambda}$ did increase as λ increased, but it was negatively-biased (particularly for small tree sizes), and was consistently less than the input λ value across most of its range (Fig. 1 black line). Second, the standard deviation of $\hat{\lambda}$ was largest at intermediate values of λ and smallest at extreme values, implying that the precision in estimating λ varied across the range of input values (Fig. 1 red line). Additionally, standard deviations of $\hat{\lambda}$ were negatively associated with tree size, and for trees of 128 species or less, $\hat{\lambda}$ was quite variable, except for cases when λ was near or equal to 1. Third, the distributions of $\hat{\lambda}$ were not normal across its range, but became increasingly skewed at more extreme values of λ (Fig. 1 blue line). For small tree sizes, it was also clear that distributions were more platykurtic at intermediate values of $\hat{\lambda}$. Taken together these results reveal that $\hat{\lambda}$ inconsistently estimated phylogenetic signal, both across tree sizes and across the range of input values. Additional simulations (Supplemental Information) revealed that incorporating λ in PGLS anova and regression did not adversely affect the statistical properties of model fitting (type I error, power, parameter estimation); thus, it is reasonable to include λ in PGLS as a parameter for tuning the degree of phylogenetic signal in the dependent variables during the analysis. However, the statistical properties revealed in Fig. 1 demonstrate that λ is unsuitable as an effect size for measuring the strength of phylogenetic signal in data.

Kappa (κ) estimates of phylogenetic signal are more stable. Simulation results for κ demonstrated that this measure displayed better statistical properties. First, as expected, mean values of κ increased consistently with increasing signal (λ) irrespective of tree size (Fig. 2 black line). Additionally, the standard deviation of κ was consistent across tree sizes (Fig. 2 red line), and while it increased with λ , it was always less than the corresponding mean. This finding is perhaps unsurprising, as κ is bound by 0, and was never large for small values of λ . Importantly, κ was normally distributed across the range of input λ , and remained consistent in this

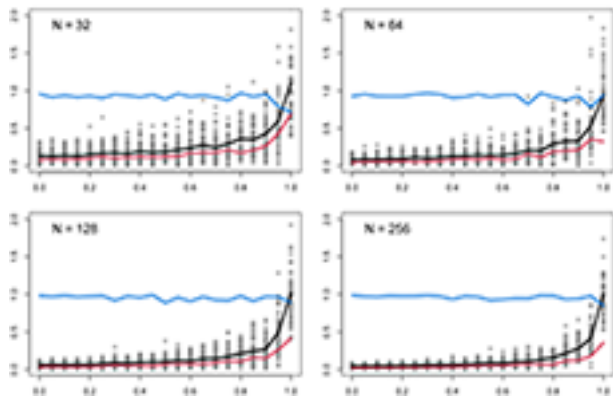


Fig. 2. Patterns in κ .

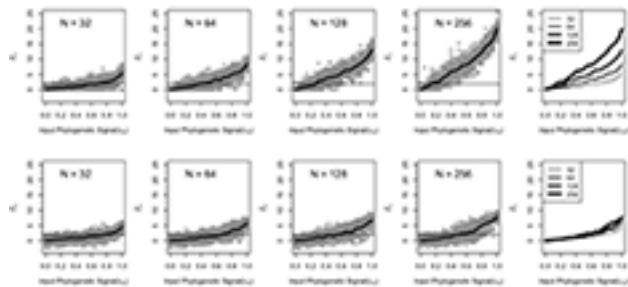


Fig. 3. Patterns in Z .

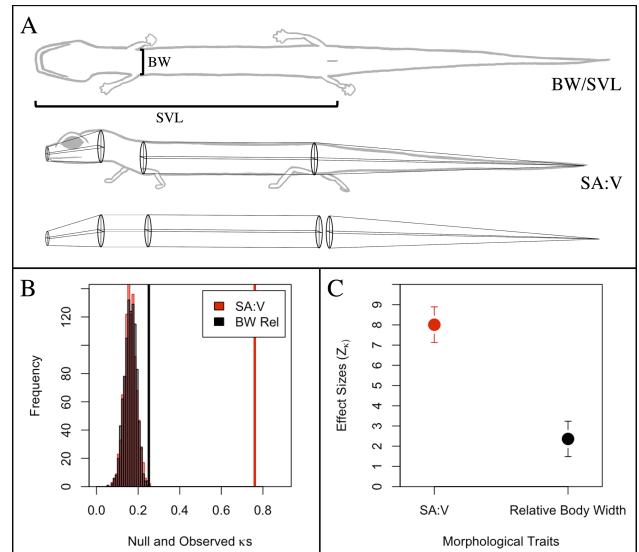


Fig. 4. (A) Linear measures for relative body size, and regions of the body used to estimate surface area to volume (SA:V) ratios. (B) Permutation distributions of phylogenetic signal for SA:V and $\frac{BW}{SVL}$, with observed values shown as vertical bars. (C) Effect sizes (Z_k) for SA:V and $\frac{BW}{SVL}$,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

cantly stronger in SA:V (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 (34). Thus, greater phylogenetic signal in SA:V is to be expected.

2. Discussion

To be edited. 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.

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 (13, 14, 16–18), methods that explicitly measure the strength of phylogenetic signal, or facilitate comparisons among datasets, have remained underdeveloped. In this study, 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 the precision of λ increased with increasing sample sizes; a pattern noted previously (25), and one that conformed with parametric statistical theory (35). However, we also found that vastly different λ estimates could be obtained from data containing the same level of phylogenetic signal, and that similar λ estimates may be obtained from data containing differing levels of phylogenetic signal. Further, the precision of λ varied with the strength of phylogenetic signal, where lower precision was observed when in data whose phylogenetic signal was of intermediate strength. From these findings we conclude that λ is not a reliable indicator of the observed strength of phylogenetic signal in phenotypic datasets, and

that biological interpretations of the strength of signal based on this parameter may inaccurately characterize such effects.

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

Based on the effect size Z_κ , we then proposed a two-sample test, which provides means of determining whether the strength of phylogenetic signal is greater in one phenotypic trait as compared to another, via a hypothesis test. Prior studies have summarized patterns of variation in phylogenetic signal across datasets using summary test values, such as κ (14). However, κ does not scale linearly with input levels of phylogenetic signal, and its variance increases (i.e., precision decreases) with increasing strength of phylogenetic signal (20, 22). Thus, κ should not be considered an effect size that measures the strength of phylogenetic signal on a common scale. By contrast, standardizing κ (Z_κ , via equation 2) alleviates these concerns, and facilitates formal statistical comparisons of the strength of signal across datasets. Thus when viewed from this perspective, the approach developed here aligns well with other statistical approaches such as meta-analysis (36–38), where summary statistics across datasets are converted to standardized effect sizes for subsequent “higher order” statistical summaries or comparisons. As such, our approach enables evolutionary biologists to quantitatively examine the relative strength of phylogenetic signal across a wide range of phenotypic traits, and thus opens the door for future discoveries that inform on how phenotypic diversity accumulates in macroevolutionary time across the tree of life.

One important advantage of the approach advocated here is that the resulting effect sizes (Z_κ) are dimensionless, as the units of measurement cancel out during the calculation of Z (39). Thus, Z_κ represents the strength of phylogenetic signal on a common and comparable scale – measured in standard deviations – regardless of the initial units and original scale of the phenotypic variables under investigation. This means that the strength of phylogenetic signal may be compared across datasets for continuous phenotypic traits measured in different units and scale, because those units have been standardized through their conversion to Z_κ . For example, our approach could be utilized to determine whether the strength of phylogenetic signal (say, in response to ecological differentiation) is stronger in morphological traits (linear traits: mm), physiological traits (metabolic rate: $\frac{O^2}{min}$), or behavioral traits (aggression: $\frac{\#displays}{second}$). In fact, our empirical example provided such a comparison, as SA:V is represented in mm^{-1} while relative body size is a unitless ratio ($\frac{BW}{SVL}$). Ad-

ditionally, our method is capable of comparing the strength of phylogenetic signal in traits of different dimensionality, as estimates of phylogenetic signal using κ have been generalized for multivariate data (18). Furthermore, tests based on \hat{Z}_{12} may be utilized for comparing the strength of phylogenetic signal among datasets containing a different number of species, and even for phenotypes obtained from species in different lineages, because their phylogenetic non-independence and observed variation are taken into account in the generation of the empirical sampling distribution via permutation.

This study is not the first to compare λ and κ for their ability as statistics to measure phylogenetic signal. Our results for λ and κ values are consistent with those found in the simulations performed by Münkemüller et al. (20), but that study investigated type I error rates and statistical power, finding that λ performed better in both regards, irrespective of species number in trees. Although not the central focus of their study, the same tendency for variable λ and consistent κ at intermediate phylogenetic signal strengths was observed (20). Recent work by Molina-Venegas and Rodríguez (23) found that κ but not λ tended to inflate the estimate of phylogenetic signal, leading to moderate type I and type II biases, if polytomic chronograms were used. Their work more thoroughly addressed previous observations of inflated κ for incompletely resolved phylogenetic trees (20, 40). An interesting question is whether an inflated κ value leads to an inflated Z_κ or does a tendency of a particular tree to inflate estimates of κ also inflate the values in random permutations of a test, in which case Z_κ is robust to polytomies? We repeated the analyses in Figure 4, adjusting trees to have 50% collapsed nodes, per the technique of Molina-Venegas and Rodríguez (23), and found results were consistent (see Supporting Information). This confirms that any tendency of incompletely resolved trees to inflate κ as a descriptive statistic does not inflate Z_κ as an effect size. Furthermore, because comparison of effect sizes in a test is a comparison of locations of observed values in their sampling distributions, which would shift concomitantly because of this tendency, the Z_{12} test statistic in equation 4 appears to be robust in spite of unresolved trees.

Phylogenetic signal can be thought of as both an attribute to be measured in the data and a parameter that can be tuned to account for the phylogenetic non-independence among observations, for analysis of the data. As such, λ is appealing, as a statistic that potentially fulfills both roles. However, the inability to estimate phylogenetic signal with λ for data simulated with known phylogenetic signal is troublesome, and we recommend evolutionary biologists refrain from viewing it as a statistic to describe the amount of phylogenetic signal in the data. Interestingly, κ – when standardized to an effect size Z_κ – is a better statistic for measuring the amount of phylogenetic signal in data simulated with respect to known levels of λ . Although λ might be viewed as an important parameter for modifying the the conditional estimation of linear model coefficients with respect to phylogeny, it is neither a statistic that has meaningful comparative value as a measure of phylogenetic signal nor a statistic that lends itself well to reliable calculation of a test statistic. By contrast, κ has been shown here to be a reliable statistic, but only when standardized by the mean and standard deviation of its empirical sampling distribution (i.e., when converted to

the effect size, Z_κ). Because one has control over the number of permutations used in analysis, one can be assured with many permutations that the empirical sampling distribution is representative of true probability distributions (12). With low coefficients of variation for Z_κ (Figure 4), it is difficult to imagine that a hypothesis test can improve equation 4 for efficiently comparing phylogenetic signal for different traits, different trees, or a combination of both.

3. Methods

Derivation of effect sizes. 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 (35, 37, 41). 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 (42–45) 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 (46), 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 thus confidence intervals) may be obtained. 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:

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

where, $\hat{\lambda}$ is the maximized likelihood value of λ and χ_λ^2 is the likelihood ratio statistic for the value. Note that alternative formulations could be envisioned: $Z_\lambda = d\sqrt{\chi_\lambda^2}$, where d is a binary value ($-1, 1$) to indicate a direction based on whether $\hat{\lambda}$ is below or above a critical value of λ , for a quantile from a χ^2 distribution at a probability of 0.5. However, preliminary investigations found that mapping χ_λ^2 values in this manner did not produce effect sizes that were symmetrical about $Z = 0$, as the mapping was not linear and the log-likelihood profiles can be rather flat for small trees. **Is alternative needed here?**

Derivation of two-sample test statistic. To compare the strength of phylogenetic signal across datasets, 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 κ_1 , κ_2 , $\hat{\mu}_{\kappa_1}$, $\hat{\mu}_{\kappa_2}$, $\hat{\sigma}_{\kappa_1}$, and $\hat{\sigma}_{\kappa_2}$ are as defined above for equation 2. The right side of the equation 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 (43, 45).

Simulations. Simulations were conducted by generating pure-birth phylogenies at each of six different tree sizes ($n = 2^5, 2^6, \dots, 2^{10}$), and with differing levels of phylogenetic signal ($\lambda = 0.0, 0.5, \dots, 1.0$). We generated 100 random trees for each intersection of tree size and λ . For each λ within each tree size, continuous traits were then simulated on each phylogeny under a BM model of evolution. For each set of 100 trees we measured the mean values of $\hat{\lambda}$ and κ , their standard deviation, and calculated the Shapiro-Wilk W statistic as a departure from normality (symmetry). For the latter, a value of 1.0 indicates normally distributed values, while departures from 1.0 indicate skewness. Simulations were then repeated for both balanced and pectinate trees, which yielded qualitatively similar results (see Supporting Information). Trees containing polytomies, and an evaluation of $\hat{\lambda}$ from models of linear regression and phylogenetic ANOVA, were also investigated, and results were qualitatively similar to those reported above (see Supporting Information).

Empirical Data. Surface area to volume ratios (SA:V) and relative body width ($\frac{BW}{SVL}$) measures were obtained from as species means from individuals of 305 species, from which species means were obtained (33, 34). A time-dated molecular phylogeny for the group (47) was pruned to match the species in the phenotypic dataset. The phylogenetic signal in each trait was then characterized using κ , which was converted to its effect size (Z_κ) using **geomorph** 3.3.1 (48, 49), and routines by the authors **(to be incorporated in geomorph upon manuscript acceptance)**.

ACKNOWLEDGMENTS. We thank E. Glynne and B. Juarez for comments on early drafts of the manuscript. This work was supported in part by NSF grant DBI-1902511 (to D.C.A.) and DBI-1902694 (to M.L.C.).

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. Butler MA, King AA (2004) Phylogenetic comparative analysis: A modeling approach for adaptive evolution. *American Naturalist* 164:683–695.
7. 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.
8. O'Meara BC, Ane C, Sanderson MJ, Wainwright PC (2006) Testing for different rates of continuous trait evolution using likelihood. *Evolution* 60:922–933.
9. Revell LJ, Harmon LJ (2008) Testing quantitative genetic hypotheses about the evolutionary rate matrix for continuous characters. *Evolutionary Ecology Research* 10:311–331.
10. Beaulieu JM, Jhwueng DC, Boettiger C, O'Meara BC (2012) Modeling stabilizing selection: Expanding the ornstein-uhlenbeck model of adaptive evolution. *Evolution* 66:2369–2383.
11. Adams DC (2014) A method for assessing phylogenetic least squares models for shape and other high-dimensional multivariate data. *Evolution* 68:2675–2688.
12. Adams DC, Collyer ML (2018) Phylogenetic anova: Group-clade aggregation, biological challenges, and a refined permutation procedure. *Evolution* 72(6):1204–1215.
13. Pagel MD (1999) Inferring the historical patterns of biological evolution. *Nature* 401:877–884.
14. Blomberg SP, Garland T, Ives AR (2003) Testing for phylogenetic signal in comparative data: Behavioral traits are more labile. *Evolution* 57:717–745.
15. Revell LJ, Harmon LJ, Collar DC (2008) Phylogenetic signal, evolutionary process, and rate. *Systematic Biology* 57:591–601.
16. Abouheif E (1999) A method for testing the assumption of phylogenetic independence in comparative data. *Evolutionary Ecology Research* 1:895–909.
17. Gittleman JL, Kot M (1990) Adaptation: Statistics and a null model for estimating phylogenetic effects. *Systematic Zoology* 39(3):227–241.
18. Adams DC (2014) A generalized Kappa statistic for estimating phylogenetic signal from shape and other high-dimensional multivariate data. *Systematic Biology* 63:685–697.
19. Klingenberg CP, Gidaszewski NA (2010) Testing and quantifying phylogenetic signals and homoplasy in morphometric data. *Systematic biology* 59(3):245–261.
20. Münkemüller T, et al. (2012) How to measure and test phylogenetic signal. *Methods in Ecology and Evolution* 3:743–756.
21. 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.
22. 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.
23. 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.
24. Revell LJ (2010) Phylogenetic signal and linear regression on species data. *Methods in Ecology and Evolution* 1:319–329.
25. Boettiger C, Coop G, Ralph P (2012) Is your phylogeny informative? Measuring the power of comparative methods. *Evolution* 67:2240–2251.
26. Freckleton RP, Harvey PH, Pagel M (2002) Phylogenetic analysis and comparative data: A test and review of evidence. *American Naturalist* 160:712–726.
27. Cooper N, Jetz W, Freckleton RP (2010) Phylogenetic comparative approaches for studying niche conservatism. *Journal of Evolutionary Biology* 23(12):2529–2539.
28. Bose R, Ramesh BR, Pélissier 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.
29. Vandeloof F, et al. (2019) Nectar traits differ between pollination syndromes in balsaminaceae. *Annals of Botany* 124(2):269–279.
30. 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.
31. 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.
32. 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.
33. Baken EK, Adams DC (2019) Macroevolution of arboreality in salamanders. *Ecology and Evolution* 9(12):7005–7016.
34. 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.
35. Cohen J (1988) *Statistical power analysis for the behavioral sciences* (Routledge).
36. Hedges L. V., Olkin I (1985) *Statistical methods for meta-analysis* (Elsevier).
37. Glass GV (1976) Primary, secondary, and meta-analysis of research. *Educational Researcher* 5:3–8.
38. Arnqvist G., Wooster D (1995) Meta-analysis: Synthesizing research findings in ecology and evolution. *Trends in Ecology and Evolution* 10:236–240.
39. Sokal R. R., Rohlf FJ (2012) *Biometry* (W.H. Freeman & Co., San Francisco). 4th Ed.
40. Davies TJ, Kraft NJ, Salamin N, Wolkovich EM (2012) Incompletely resolved phylogenetic trees inflate estimates of phylogenetic conservatism. *Ecology* 93(2):242–247.
41. Rosenthal R (1994) The handbook of research synthesis. ed Cooper LV H Hedges (Russell Sage Foundation), pp 231–244.
42. 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.
43. Adams DC, Collyer ML (2016) On the comparison of the strength of morphological integration across morphometric datasets. *Evolution* 70:2623–2631.
44. 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.
45. 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.
46. Adams DC, Collyer ML (2019) Phylogenetic comparative methods and the evolution of multivariate phenotypes. *Annual Review of Ecology, Evolution, and Systematics* 50:405–425.
47. Bonett RM, Blair AL (2017) Evidence for complex life cycle constraints on salamander body form diversification. *Proceedings of the National Academy of Sciences, USA* 114:9936–9941.
48. Adams DC, Otárola-Castillo E (2013) Geomorph: An r package for the collection and analysis of geometric morphometric shape data. *Methods in Ecology and Evolution* 4:393–399.
49. Adams DC, Collyer ML, Kaliontzopoulou A (2020) Geomorph: Software for geometric morphometric analyses. R package version 3.3.1. Available at: <https://cran.r-project.org/package=geomorph>.