

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

Dean C. Adams<sup>a,1</sup>, Erica K. Baken<sup>a,b</sup>, and Michael L. Collyer<sup>b</sup>

<sup>a</sup>Department of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames, Iowa, 50010. USA.; <sup>b</sup>Department of Science, Chatham University, Pittsburgh, Pennsylvania, 15232. USA.

This manuscript was compiled on July 17, 2020

Macroevolutionary studies frequently characterize the phylogenetic signal in phenotypes, however, analytical tools for comparing the strength of that signal across traits remain largely underdeveloped. Here we evaluate the efficacy of Pagel's  $\lambda$  to correctly estimate the strength of phylogenetic signal in phenotypic traits across a range of input values. We find that  $\lambda$  behaves as a Bernoulli random variable, where estimates are increasingly skewed at larger and smaller input levels of phylogenetic signal. Further, the precision of  $\lambda$  varies with input signal. Another measure, Blomberg's  $\kappa$ , is more consistent across a range of tree sizes, and exhibits a positive relationship with input levels of phylogenetic signal. However, that relationship is decidedly nonlinear. Thus, neither  $\lambda$  nor  $\kappa$  are suitable as effect sizes for measuring the strength of phylogenetic signal, and comparing that signal across datasets. As an alternative, we propose a standardized effect size based on  $\kappa$ , ( $Z_\kappa$ ), which measures the strength of phylogenetic signal more reliably than does  $\lambda$ , and places that signal on a common scale for statistical comparison. We develop tests based on  $Z_\kappa$  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 (3–10). 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, 11, 12). 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, 12, 13).

Several analytical tools have been developed to quantify phylogenetic signal in phenotypic datasets (11, 12, 14–17), 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 (13, 16, 18–23). One of the most widely used methods for characterizing phylogenetic signal is Pagel's  $\lambda$  (11), which transforms the lengths of the internal branches of the phylogeny to im-

prove the fit of data to the phylogeny via maximum likelihood (11, 24). When incorporated in PGLS,  $\lambda$  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 (11, 24). To infer whether phylogenetic signal differs from no signal or a Brownian motion model of evolutionary divergence, the observed model fit using  $\hat{\lambda}$  may be statistically compared to that using  $\lambda = 0$  or  $\lambda = 1$  via likelihood ratio tests (24–26) or confidence limits (27).

Another widely used measure of phylogenetic signal is Blomberg's  $\kappa$  (12), which characterizes phylogenetic signal as the ratio of observed trait variation to the amount of variation expected under Brownian motion. Blomberg's  $\kappa$  can be treated as a test statistic by employing a permutation test to generate its sampling distribution (12, 16) for determining whether significant phylogenetic signal is present in data. Both  $\lambda$  and  $\kappa$  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  $\lambda$  and  $\kappa$  as descriptive statistics that measure the relative strength of phylogenetic signal, providing an estimate of its magnitude for comparison.

The appeal of Pagel's  $\lambda$  and Blomberg's  $\kappa$  as descriptive

## Significance Statement

Evolutionary biologists wish to quantify and compare the strength of phylogenetic signal across datasets, but analytical tools for these comparisons are generally lacking. Here we develop a standardized effect size,  $Z_\kappa$ , 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  $\lambda$  and Blomberg's  $\kappa$ ), 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.

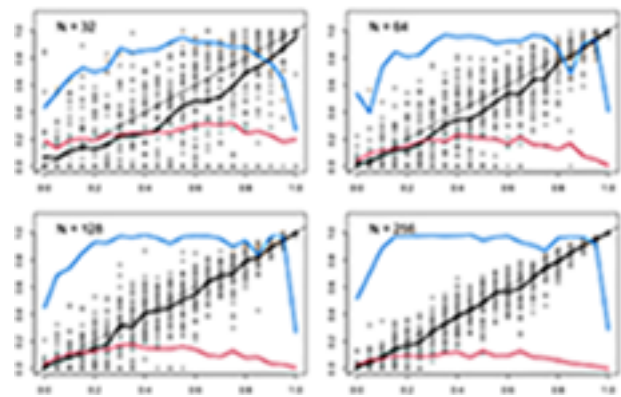
Data deposition: Data for the empirical example may be found on DRYAD: doi:10.5061/dryad.b554m44 and doi:10.5061/dryad.59zw3r23m. R-scripts for simulation tests are found on Github: XXX. Computer code for implementing the two-sample comparison of effect sizes is found in geomorph: <https://cran.r-project.org/web/packages/geomorph/index.html>

<sup>1</sup> To 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  $\kappa$ , respectively, in a comparative sense (28–30). 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  $\lambda$  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  $\lambda$ ; the dispersion of  $\hat{\lambda}$  would be largest at intermediate values of  $\lambda$ ,  $\hat{\lambda}$  would be predictable over the range of  $\lambda$  with respect to tree size; the distribution of  $\hat{\lambda}$  would be symmetric at intermediate values of  $\lambda$  and more skewed toward values of 0 or 1; and that the distribution of  $\hat{\lambda}$  will be more platykurtic at intermediate values of  $\lambda$ , 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  $\lambda$  (18), interpretation of these patterns is challenging.

By contrast, for Blomberg’s  $\kappa$ , which is positively unbounded, we might expect that for any  $\lambda$  used to generate data, estimates of  $\kappa$  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  $\kappa$  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  $\lambda$  or Blomberg’s  $\kappa$ , evaluation of statistical moments across a range of  $\lambda$  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 (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  $\kappa$ , 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  $\lambda$ , and that the precision of  $\hat{\lambda}$  is not constant across its range. By contrast, estimates of  $\kappa$  are more consistent across tree sizes, and are normally distributed across the range of input levels of  $\lambda$ , making  $\kappa$  a more reliable statistic. We then propose an effect size based on  $\kappa$ , ( $Z_\kappa$ ), which provides consistent estimates



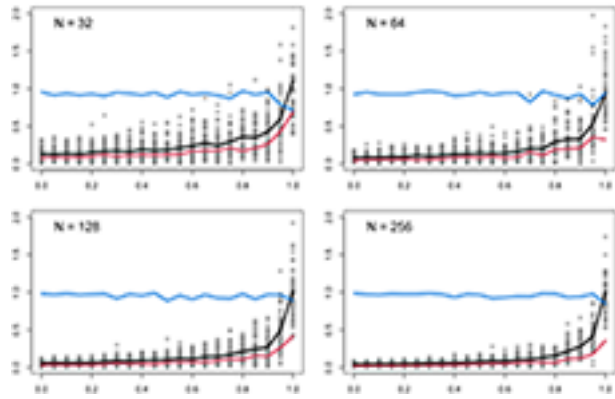
**Fig. 1.** Response of Pagel’s  $\lambda$  to increasing strength of Brownian motion. Gray line signifies the 1:1 line where the input value matches the estimate  $\hat{\lambda}$ . At each input level, the dark black line represents the empirically derived expected value (mean) of  $\hat{\lambda}$ , 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 ( $\lambda$ ) 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  $\lambda$  increases, but it is negatively-biased (particularly for small tree sizes), and is consistently less than the input  $\lambda$  value across most of its range (Fig. 1 black line). Second, the standard deviation of  $\hat{\lambda}$  is largest at intermediate values of  $\lambda$  and smallest at extreme values, implying that the precision in estimating  $\lambda$  varies across the range of input values (Fig. 1 red line). 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  $\lambda$  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  $\lambda$  (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, 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  $\lambda$  is unsuitable as an effect size for measuring the strength of phylogenetic signal in data, or for comparing that signal across datasets.

**Kappa ( $\kappa$ ) 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 ( $\lambda$ ) irrespective of tree size, though the increase does not scale linearly with input levels of phylogenetic signal (Fig. 2 black line). Additionally,



**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).

the standard deviation of  $\hat{\kappa}$  is consistent across tree sizes (Fig. 2 red line), and while it increases with  $\lambda$ , 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  $\lambda$ . Importantly,  $\hat{\kappa}$  is normally distributed across the range of input  $\lambda$ , and remains consistent in this pattern regardless of tree size (Fig. 2 blue line). In fact, distributions of  $\hat{\kappa}$  display slight skewing only at small tree sizes and when  $\lambda$  is near or equal to 1. This result differs from those of (18), where skewing appears to be from combining random values generated independently. These findings reveal that while  $\kappa$  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  $\kappa$  ( $Z_\kappa$ ) better characterize phylogenetic signal.** To measure the strength of phylogenetic signal on a common scale, here we propose effect sizes ( $Z$ -scores) for both  $\lambda$  and  $\kappa$ . Statistically, a standardized effect size may be found as:

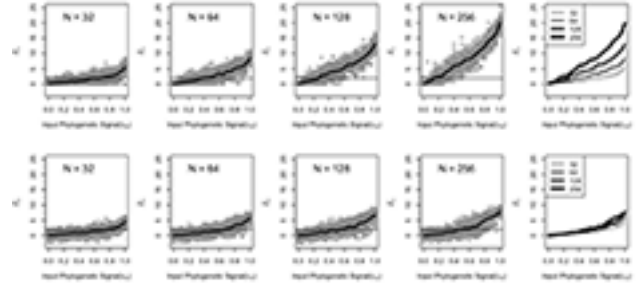
$$Z_\theta = \frac{\theta_{obs} - E(\theta)}{\sigma_\theta} \quad [1]$$

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

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

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

where  $\kappa_{obs}$  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  $\kappa$  obtained via permutation. The



**Fig. 3.** Response of effect sizes  $Z_\lambda$  and  $Z_\kappa$  to increasing strength of Brownian motion.

empirical sampling distribution of  $\kappa$  is first transformed via Box-Cox to better adhere to the assumption of normality.

For  $\lambda$ , deriving an effect size is more challenging, as  $\lambda$  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  $\lambda$  that intersect the log-likelihood profile for corresponding percentiles of the  $\chi^2$  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  $\lambda$  may be found as:

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

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

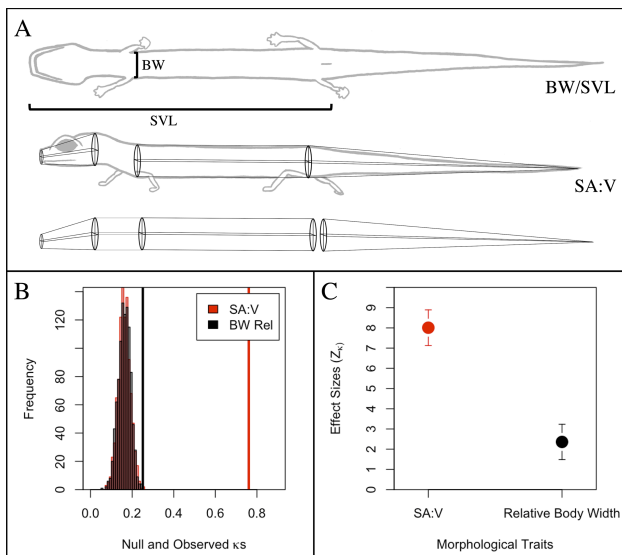
Here we evaluate the ability of  $Z_\lambda$  and  $Z_\kappa$  to characterize known levels of phylogenetic signal. Both  $Z_\lambda$  and  $Z_\kappa$  are associated with input phylogenetic signal ( $\lambda$ ), indicating that both statistics capture the observed signal (Fig. 3). However, effect sizes 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,  $Z_\kappa$  is much more consistent across tree sizes, and increases more linearly with increasing levels of phylogenetic signal. Additionally,  $Z_\kappa$  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  $Z_\kappa$  across the range of input values. Thus between the two statistics,  $Z_\kappa$  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:

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

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  $Z_\kappa$  has already been calculated for two sampling distributions as in equation 2, the sampling distributions have unit variance for





**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  $\sqrt{n}$ .

each of the  $Z_k$  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  $Z_k$  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.

## 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  $\lambda$ , and explored its efficacy for characterizing the strength of phylogenetic signal in phenotypic data. Using computer simulations, we found that  $\lambda$  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  $\lambda$  in estimating actual levels of phylogenetic signal varies with both tree size (see also ref. (23)) and input levels of phylogenetic sig-

nal. From these findings we conclude that  $\lambda$  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  $\lambda$  and  $\kappa$ , 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 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. 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  $\kappa$  (12). However, because  $\kappa$  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  $\kappa$  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 (31, 41, 42), 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 ( $Z_k$ ) are dimensionless, as the units of measurement cancel out during the calculation of  $Z$  (43). Thus,  $Z_k$  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  $Z_k$ . 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 just such a comparison, as SA:V is represented in  $mm^{-1}$  while relative body size is a unitless ratio ( $\frac{BW}{SVL}$ ). Additionally, our method is capable of comparing the strength of phylogenetic signal in traits of different dimensionality, as estimates of phylogenetic signal using  $\kappa$  have been generalized for multivariate data (16). 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  $\lambda$  and  $\kappa$  for their ability as statistics to measure phylogenetic signal. Our results for  $\lambda$  and  $\kappa$  values are consistent with those found in the simulations performed by Münkemüller et al. (18), but that study investigated type I error rates and statistical power, finding that  $\lambda$  performed better in both regards, irrespective of species number in trees. Although not the central focus of their study, the same tendency for variable  $\lambda$  and consistent  $\kappa$  at intermediate phylogenetic signal strengths was observed (Fig. 2 of ref. (18)). Recent work by Molina-Venegas and Rodríguez (21) found that  $\kappa$  but not  $\lambda$  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  $\kappa$  for incompletely resolved phylogenetic trees (18, 44). An interesting question is whether an inflated  $\kappa$  value leads to an inflated  $Z_\kappa$  or does a tendency of a particular tree to inflate estimates of  $\kappa$  also inflate the values in random permutations of a test, in which case  $Z_\kappa$  is robust to polytomies? We repeated the analyses in Figs. 1 & 2, adjusting trees to have 50% collapsed nodes, per the technique of Molina-Venegas and Rodríguez (21), and found results were consistent (Supporting Information). This confirms that any tendency of incompletely resolved trees to inflate  $\kappa$  as a descriptive statistic does not inflate  $Z_\kappa$  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,  $\lambda$  is appealing, as a statistic that potentially fulfills both roles. However, the inability to estimate phylogenetic signal with  $\lambda$  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,  $\kappa$  – when standardized to an effect size  $Z_\kappa$  – is a better statistic for measuring the amount of phylogenetic signal in data simulated with respect to known levels of  $\lambda$ . Although  $\lambda$  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,  $\kappa$  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,  $Z_\kappa$ ). 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 (10). Given the greater consistency in estimates of  $Z_\kappa$  across tree sizes and input signal, it is difficult to imagine a hypothesis test that can improve equation 4 for efficiently comparing phylogenetic signal for different traits, different trees, or a combination of both.

### 3. Methods

**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  $\lambda$ . For each  $\lambda$  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  $\kappa$ , 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 (39, 40). A time-dated molecular phylogeny for the group (45) was pruned to match the species in the phenotypic dataset. The phylogenetic signal in each trait was then characterized using  $\kappa$ , which was converted to its effect size ( $Z_\kappa$ ) using **geomorph** 3.3.1 (46, 47), 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. 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, Jhwueng 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. Pavoiné 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, 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.
27. Vandeloof 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.
41. Hedges L. V., Olkin I (1985) *Statistical methods for meta-analysis* (Elsevier).
42. Arnqvist G., Wooster D (1995) Meta-analysis: Synthesizing research findings in ecology and evolution. *Trends in Ecology and Evolution* 10:236–240.
43. Sokal R. R., Rohlf FJ (2012) *Biometry* (W.H. Freeman & Co., San Francisco). 4th Ed.
44. Davies TJ, Kraft NJ, Salamin N, Wolkovich EM (2012) Incompletely resolved phylogenetic trees inflate estimates of phylogenetic conservatism. *Ecology* 93(2):242–247.
45. 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.
46. 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.
47. 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>.