## APPLICATION

# `Rphylopars`: fast multivariate phylogenetic comparative methods for missing data and within-species variation

**Eric W. Goolsby[1,2,*], Jorn Bruggeman[3] and Cécile Ané[4,5]**

[1]*Arnold Arboretum, Harvard University, Boston, MA 02131, USA;* [2]*Department of Plant Biology, Interdisciplinary Toxicology Program, University of Georgia, Athens, GA 30602, USA;* [3]*Plymouth Marine Laboratory, Prospect Place, The Hoe, Plymouth PL1 3DH, UK;* [4]*Department of Statistics, University of Wisconsin-Madison, 1300 University Avenue, Madison, WI 53706, USA; and* [5]*Department of Botany, University of Wisconsin-Madison, Wisconsin-Madison, WI 53706, USA*

### Summary

**1.** Over the past several years, phylogenetic comparative studies have increasingly approached trait evolution in a multivariate context, with a number of taxa that continues to rise dramatically. Recent methods for phylogenetic comparative studies have provided ways to incorporate measurement error and to address computational challenges. However, missing data remain a particularly common problem, in which data are unavailable for some but not all traits of interest for a given species (or individual), leaving researchers with the choice between omitting observations or utilizing imputation-based approaches.

**2.** Here, we introduce an R implementation of `PhyloPars`, a tool for phylogenetic imputation of missing data and estimation of trait covariance across species (phylogenetic covariance) and within species (phenotypic covariance). `Rphylopars` provides expanded capabilities over the original `PhyloPars` interface including a fast linear-time algorithm, thus allowing for extremely large data sets (which were previously computationally infeasible) to be analysed in seconds or minutes rather than hours.

**3.** In addition to providing fast and computationally efficient implementations, we introduce in `Rphylopars` methods to estimate macroevolutionary parameters under alternative evolutionary models (e.g. Early-Burst, multivariate Ornstein-Uhlenbeck).

**4.** By providing fast and computationally efficient methods with flexible options for various phylogenetic comparative approaches, `Rphylopars` expands the possibilities for researchers to analyse large and complex data with missing observations, within-species variation and deviations from Brownian motion.

**Key-words:** fast methods, linear-time algorithm, missing data, multivariate Ornstein-Uhlenbeck, phylogenetic comparative method, phylogenetic generalized least squares, phylogenetic imputation

## Introduction

Phylogenetic comparative methods provide tools for studying trait evolutionary history and trait covariance while accounting for non-independence of data collected across species. Since the introduction of phylogenetically independent contrasts (PICs; Felsenstein 1985), thousands of studies have incorporated phylogeny into statistical analyses. Since then, a number of theoretical advances have been made, including a flexible generalization of PICs with phylogenetic generalized least squares (PGLS; Martins & Hansen 1997), and subsequent generalizations in the form of phylogenetic mixed models (PMMs; Housworth, Martins & Lynch 2004; Ives, Midford & Garland 2007; Hadfield & Nakagawa 2010). These developments provide a unified framework for conducting comparative analyses, allowing for flexible incorporation of alternative evolutionary models (Hansen 1997), within-species variation

(Ives, Midford & Garland 2007; Hansen & Bartoszek 2012) and high-dimensional multivariate extensions (Adams 2014a, b; Denton & Adams 2015). Felsenstein (2008) introduced an extension of PICs incorporating within-species variation in an expectation-maximization algorithm that simultaneously estimates across-species (phylogenetic) and within-species (phenotypic) trait covariation. This model is conceptually and statistically similar to PMMs as well as within-species methods proposed by Ives, Midford & Garland (2007), with the main difference being that Ives, Midford & Garland (2007) use summary statistics (species trait means and standard errors) whereas Felsenstein's model utilizes raw observations from individuals (Felsenstein 2008).

### Dealing with missing data in comparative studies

A frequently encountered problem in comparative studies is the difficulty (or impossibility) of obtaining observations for every trait from each species in a study. For example, a realistic

*Correspondence author. E-mail: eric.goolsby.evolution@gmail.com

scenario might be 10 variables and a 10% missing rate. Most often, researchers handle missing data by omitting individuals for which all observations are not available. In the previous example, this would lead to excluding 65% of individuals $(1 - 0.9^{10})$, a drastic reduction in the data size, if the values are missing at random. Alternatively, researchers sometimes rely on pairwise observations to estimate pairwise trait covariance, which may result in a non-invertible covariance matrix because each covariance element is calculated from a different subset of observations (Arbuckle, Marcoulides & Schumacker 1996). This is undesirable, as a singular covariance matrix corresponds to an undefined log-likelihood, rendering likelihood-based parameter estimation, model diagnostics and model selection procedures (e.g. likelihood ratio tests, AIC, BIC) impossible. Both approaches are problematic, resulting in unnecessary loss of statistical power and risking substantial bias in parameter estimates (Pakeman 2014). Accordingly, comparative studies in which data are assumed to be missing at random should utilize methods incorporating all available observations.

Although Felsenstein's model requires that all observations are completely available for each trait, individual and species, the possibility of developing a likelihood-based modification of the algorithm to estimate trait covariance in the presence of missing data was suggested (Felsenstein 2008). This idea was implemented in an online web interface called `PhyloPars` (Bruggeman, Heringa & Brandt 2009), which is a statistical framework for estimating phylogenetic trait covariance while accounting for both within-species variation and missing data. `PhyloPars` can also be used to phylogenetically impute missing species data, perform ancestral state reconstruction and test hypotheses of correlated trait evolution, among others.

## Computational feasibility of large-scale comparative studies

Until recently, most comparative studies involved at most a few hundred species. However, in recent years, comparative studies have begun expanding in size to include several thousand species (e.g. Smith *et al.* 2011; FitzJohn *et al.* 2014). Because PGLS-based analyses require inversion of the phylogenetic covariance matrix corresponding to the tree topology and the evolutionary model [most often Brownian motion (BM)], large-scale comparative studies can become prohibitively time-consuming – with single inversions taking several days or weeks to complete or, in some cases, failing entirely (Ho & Ané 2014), as the computational time required to invert square matrices grows faster than the square of the number of rows and columns. Additionally, many comparative algorithms, including `PhyloPars`, require thousands of matrix inversions for likelihood-based parameter estimation. The inclusion of multiple traits and multiple within-species observations dramatically exacerbates this problem: for a study with *s* species, *m* traits and *k* within-species observations, a matrix with *smk* rows and columns must be inverted. For example, a study with 5000 species, 4 traits and 5 within-species

observations per trait results in a matrix with 100 000 rows and columns (10 billion individual cells total). A study of this size relying on repeated direct matrix inversions is simply not feasible.

To meet the demand for computationally feasible comparative methods, multiple algorithms have been developed to run in linear rather than polynomial time (Felsenstein 1973; FitzJohn 2012; Freckleton 2012; Ho & Ané 2014). For example, PIC-based calculations, which are statistically equivalent to PGLS (assuming BM evolution and a single observation per species), are linear in time with the number of species included. This property of PICs was exploited by Freckleton (2012) to develop fast methods for comparative likelihood calculations in a variety of applications. More recently, Ho & Ané (2014) developed a versatile linear-time algorithm which can be used to calculate many different quantities required for different phylogenetic comparative applications.

Here, we modify the methods presented in Ho & Ané (2014) to develop a linear-time implementation of `PhyloPars` (Bruggeman, Heringa & Brandt 2009) in ʀ called `Rphylopars` (Appendix S1, Supporting Information). Our algorithm avoids large matrix inversions and allows for extremely large problems to be analysed using modest computational resources (e.g. a standard laptop) while avoiding excessive memory burdens typically associated with large-scale comparative analyses. We further extend the `Rphylopars` model to allow for within-species trait correlations (the original `PhyloPars` implementation assumed zero within-species correlations). We also implement methods for incorporating alternative evolutionary models, such as Early-Burst (EB; Harmon *et al.* 2010) and Ornstein-Uhlenbeck (OU; Hansen 1997), for multivariate data with missing observations and within-species variation.

## Rphylopars description

Models can be fit in `Rphylopars` using the *phylopars* function. This function estimates the specified evolutionary model using restricted maximum likelihood, performs ancestral state reconstruction, imputes values for missing data and provides prediction variances for ancestral states and imputed data.

### WITHIN-SPECIES VARIATION AND MISSING DATA

If multiple within-species observations are available, `Rphylopars` automatically estimates within-species (phenotypic) trait covariance in addition to among-species (phylogenetic) covariance. As in Felsenstein (2008), phenotypic covariance is assumed to be equivalent among species. The estimation of phenotypic covariance can be suppressed by setting the `pheno_error` option to `FALSE`, and the algorithm instead uses species means to estimate phylogenetic covariance assuming no within-species covariance. Alternatively, within-species variance may be estimated without estimating within-species correlation (i.e. a diagonal phenotypic

covariance matrix) by setting the `pheno_correlated` option to `FALSE`.

`Rphylopars` also readily incorporates missing observations by maximizing the log-likelihood of the covariance parameters using all available data (Bruggeman, Heringa & Brandt 2009). Using the estimated evolutionary model, missing data and ancestral states (which can also be viewed as missing data) are phylogenetically imputed as the best linear unbiased predictions, which is mathematically equivalent to universal kriging in spatial statistics (Bruggeman, Heringa & Brandt 2009; Ho & Ané 2014; Cressie 2015). This method may also be used to predict phenotypic values in completely unobserved species. By modifying the methods of Ho & Ané (2014), `Rphylopars` is able to compute these quantities in linear time, providing the maximum likelihood ancestral reconstructions and prediction covariances for each internal node, as well as predicted means and covariances for missing values at the tips of the tree (Appendix S1).

### ALTERNATIVE EVOLUTIONARY MODELS

The OU model is a popular alternative to BM, which fits parameters for both $\alpha$ (adaptation rate towards an optimal trait value) and $\Sigma$ (rate of drift variance accumulation; Hansen 1997; Bartoszek *et al.* 2012). Because the OU model can fit a variety of evolutionary patterns in addition to adaptation, $\alpha$ should be more broadly interpreted as the overall strength of phylogenetic correlation, where low values of $\alpha$ correspond to high phylogenetic correlation (Hansen 1997; Harmon *et al.* 2010). Ho & Ané (2014) discussed a fast linear-time algorithm for the univariate OU model, but concluded that a fast algorithm for the multivariate OU model could not be adapted because different traits (and their covariances) operate under different three-point structured matrices, each matrix representing a BM-like process on a transformed tree. Estimation of the $\alpha$ and $\Sigma$ matrices requires numerical optimization, and as such poses a large computational burden when large numbers of species and traits are present. In the `mvMORPH` package, Clavel, Escarguel & Merceron (2015) use Cholesky decomposition to speed up log-likelihood computations, but this and related approaches are still nonlinear in complexity and remain infeasible for extremely large problems. Additionally, the complexity of filling in the large covariance matrix for log-likelihood computations is itself nonlinear [up to $O(n^2)$]. Our linear-time algorithm overcomes both of these problems, as it allows for different three-point structured matrices among traits while simultaneously avoiding the need to explicitly build large matrices (Appendix S1). The multivariate OU model may be specified in the *phylopars* function by setting the option `model='mvOU'`. By default, $\alpha$ is a full positive-definite matrix (`full_alpha=TRUE`), in which case $\alpha$ influences both the phylogenetic correlation between species and correlated evolution between traits. Alternatively, `full_alpha` may be set to `FALSE` to estimate the model with a diagonal $\alpha$ to represent adaptation acting on traits separately, in which case $\alpha$ influences the phylogenetic correlation between species but not correlations between traits.

A special case of the multivariate OU model assumes a diagonal $\alpha$ matrix with equal values along the diagonal. This model can be estimated by applying a branch length transformation (OU $\alpha$) such that a BM-like process applies on the transformed tree. Other evolutionary models, including EB, $\lambda$, $\kappa$ and $\delta$, may be fit as tree transformations in a similar manner (Pagel 1997, 1999; Harmon *et al.* 2010; Ho & Ané 2014). Additionally, BM covariance parameters may be estimated assuming a star phylogeny (equivalent to setting $\lambda = 0$). These models may be specified in the *phylopars* function by setting the option `model` to `'OU'`, `'EB'`, `'lambda'`, `'kappa'`, `'delta'` or `'star'`, respectively.

### PARAMETER ESTIMATION

Aside from simple BM with no missing data, the evolutionary model parameters described here lack closed-form solutions. `Rphylopars` uses Broyden–Fletcher–Goldfarb–Shannon (BFGS) optimization for parameter estimation using a modified log-Cholesky parametrization (Pinheiro & Bates 1996). To provide reasonable starting parameters for optimization, two Expectation-Maximization algorithms are implemented: $EM_{Felsenstein}$ and $EM_{missing}$ (Dempster, Laird & Rubin 1977; Felsenstein 2008). The $EM_{Felsenstein}$ algorithm is used to simultaneously estimate phylogenetic and phenotypic covariance matrices, as described in Felsenstein (2008). The $EM_{missing}$ algorithm is used to estimate phylogenetic trait covariance in the presence of missing data. If multiple within-species observations are present with missing data, $EM_{Felsenstein}$ is nested within each iteration of $EM_{missing}$ using current parameter estimates to impute missing observations (Appendix S1). EM algorithms are not necessarily guaranteed to converge on the maximum likelihood solution, but generally provide reasonable starting parameters for BFGS optimization for BM models and BM-like tree transformations. For the multivariate OU model, maximum likelihood BM covariance parameters and the identity matrix for $\alpha$ are used by default to initialize numerical optimization.

For high-dimensional optimization problems, it is well known that optimization routines may converge on local optima rather than the maximum likelihood solution. To increase confidence in estimated parameters, multiple starting parameters may be tried by overriding EM-generated starting parameters. User-defined starting parameters may be supplied for phylogenetic covariance (`phylocov_start`), phenotypic covariance (`phenocov_start`) or for alternative evolutionary model parameters (`model_par_start`). Similarly, any of these parameters may be fixed during optimization by supplying the arguments `phylocov_fixed`, `phenocov_fixed` or `model_par_fixed`.
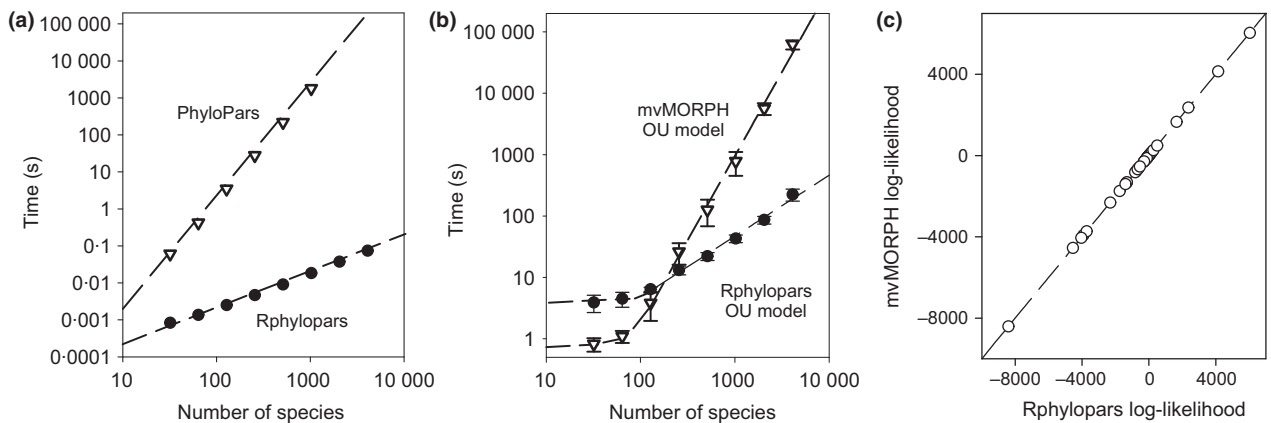
### ALTERNATIVES TO RPHYLOPARS

The original online `PhyloPars` interface can fit a BM model with missing data and multiple within-species observations (Bruggeman, Heringa & Brandt 2009). Additionally, the R package `mvMORPH` (Clavel, Escarguel &

Merceron 2015) can fit most of the models described here, including multivariate evolutionary models incorporating missing data, within-species variation and alternatives to BM evolution. Three main differences exist between `Rphylopars` and `mvMORPH` regarding the implementation of these models: (i) `mvMORPH` utilizes summary statistics (means and standard errors) to accommodate within-species variation, whereas `Rphylopars` directly utilizes raw data and can incorporate intraspecific correlations; (ii) `Rphylopars` provides reconstructed ancestral states and (co)variances at each node of the tree (as well as predicted missing species values), whereas `mvMORPH` solely estimates the root ancestral state; and (iii) `mvMORPH` relies on Cholesky factorization or similar methods to speed up log-likelihood calculations for the multivariate OU model (which exhibits nonlinear polynomial increases in computation time as the number of species or traits increases), whereas `Rphylopars` implements a fast linear-time algorithm (note: the original `PhyloPars` interface also relies on Cholesky factorization for log-likelihood calculations). Of these three differences, the third (computational cost) is perhaps the most striking: for large data sets, `Rphylopars'` linear-time algorithm reduces computation time with several orders of magnitude for large data sets, as discussed in detail in the next section. In addition to the models' specifications shared by `Rphylopars` and `mvMORPH`, `mvMORPH` also allows for many other model specifications, including shifts in evolutionary rates and multi-optima multivariate OU models. Multivariate OU models may also be fit in the OUCH (Butler & King 2004) and MVS-LOUCH (Bartoszek *et al.* 2012) packages.

COMPUTATIONAL PERFORMANCE

Simulations were performed to compare the computation times of `Rphylopars` and alternative implementations (Fig. 1). First, the speed of Cholesky decomposition (the rate-limiting step for log-likelihood calculations in the online `PhyloPars` interface) was compared to log-likelihood calculation speed in `Rphylopars` (Fig. 1a). Simulations were



**Fig. 1.** Computation time and log-likelihood comparisons between `Rphylopars` and alternative implementations for simulated data sets. For (a), four-trait data sets with five within-species replicates per species per trait were simulated on 32-, 64-, 128-, 256-, 512-, 1024-, 2048- and 4096-species pure-birth phylogenies, analysed with a BM model, and the computation time for a single log-likelihood calculation is presented. (a) `PhyloPars` (open triangles) uses a Cholesky decomposition, whose time increases cubically as the number of species increases. Cholesky decomposition failed entirely for 2048 and 4096-species data sets (a 4096-species data set with four traits and five within-species replicates corresponds to an 81 920 × 81 920 species-trait covariance matrix, or 6·71 billion matrix cells). Computation time for `Rphylopars` (solid circles) increased linearly, with the computation time for a 4096-species data set completing in < 0·1 s. For (b), the multivariate OU was fit on bivariate data sets simulated on 32-, 64-, 128-, 256-, 512-, 1024-, 2048- and 4096-species pure-birth phylogenies. Five data sets were simulated for each number of species, and the mean time to convergence in `mvMORPH` and `Rphylopars` is compared (error bars correspond to standard deviation). (c) Demonstration of the equivalence between log-likelihoods of converged parameters using `mvMORPH` and `Rphylopars` for the models fit in (b).

**Table 1.** Mean (±SD) parameter estimates and differences in parameter estimates (elements of $\Sigma$ and $\alpha$) between MVMORPH and `Rphylopars` corresponding to simulations described in the Computational performance section (Fig. 1b,c), along with the $R^2$, slopes and intercepts of linear regression between parameter estimates fit in both packages

| Parameter | Mean estimate | Avg. $\Delta$ | $R^2$ | Slope | Intercept |
|---|---|---|---|---|---|
| $\Sigma_{1,1}$ | 1·073 ± 2·456 | −0·001 ± 0·006 | 1·000 | 1·000 | 0·000 |
| $\Sigma_{1,2}$ | 0·208 ± 1·311 | −0·003 ± 0·022 | 1·000 | 1·000 | 0·003 |
| $\Sigma_{2,2}$ | 1·755 ± 1·616 | −0·014 ± 0·065 | 0·999 | 1·017 | −0·016 |
| $\alpha_{1,1}$ | 1·639 ± 2·660 | 0·0430 ± 0·273 | 0·989 | 0·999 | −0·042 |
| $\alpha_{1,2}$ | 0·116 ± 1·352 | −0·002 ± 0·010 | 1·000 | 1·000 | 0·002 |
| $\alpha_{2,2}$ | 2·272 ± 2·424 | −0·002 ± 0·016 | 1·000 | 0·999 | 0·003 |

performed using the `simtraits` function and consisted of four-trait data sets with five within-species replicates per species per trait simulated on 32-, 64-, 128-, 256-, 512-, 1024-, 2048- and 4096-species pure-birth phylogenies (the corresponding species-trait covariance matrices for these simulated data sets are of dimension $s \times 5 \times 4$). Cholesky decomposition time increased cubically as the number of species increased ($R^2 = 1\cdot000$; $y = 0\cdot167 \times 10^{-6}x^3$) and failed entirely for 2048- and 4096-species data sets (Fig. 1a), as the memory requirements for the corresponding $40\,960 \times 40\,960$ and $81\,920 \times 81\,920$ dimension matrices (respectively) exceeded available computational resources (on a standard laptop). Conversely, log-likelihood calculations increased linearly in `Rphylopars` ($R^2 = 0\cdot999$; $y = 2 \times 10^{-5}x$), and the log-likelihood for the 4096-species data set was computed in just 0·08 s (Fig. 1a). Next, the speed and convergence of `Rphylopars` was compared to that of the R package `mvMORPH` (Clavel, Escarguel & Merceron 2015) for fitting the multivariate OU model (Fig. 1b; simulation code available in Appendix S2). Simulations were performed using the `mvSIM` function in `mv-MORPH` for bivariate traits ranging from 32 to 4006 species. Simulations were performed using randomly generated positive-definite matrices for $\Sigma$ and $\alpha$ and randomly generated root values. By default, `Rphylopars` fits parameters via REML rather than ML, so the option `REML` was set to `FALSE` in order to fit models via ML for direct comparison with `mvMORPH`. Both packages exhibited a lag effect in computation time for smaller data sets (i.e. computation time was essentially unchanged regardless of the number of species for data sets with fewer than 128 species). Following the lag, computation times for `mvMORPH` increased faster than quadratically with the number of species ($R^2 = 0\cdot985$; slope = 2·759; intercept = $-5\cdot326$ on $\log_{10}$-transformed $x$ and $y$), whereas `Rphylopars` time to convergence increased approximately linearly ($R^2 = 0\cdot977$; slope = 0·989; intercept = $-1\cdot300$ on $\log_{10}$-transformed values) (Fig. 1b). Models fit using `Rphylopars` and `mvMORPH` converged to nearly identical log-likelihoods ($R^2 = 1\cdot000$; slope = 1·000; intercept = $-0\cdot110$), with `Rphylopars` log-likelihoods on average greater than `mvMORPH` log-likelihoods by 0·099 ($\pm$ 0·745; see Fig. 1c). These differences are not due to differences in the way likelihoods are calculated between the two packages, as `Rphylopars` and `mvMORPH` return identical log-likelihoods of up to $10^{-6}$ when supplied identical parameters ($R^2 = 1\cdot000$; slope = 1·000; intercept = 0·000). Rather, differences in parameter estimates are due to the difficulty of numerical estimation and highlight the importance of trying multiple starting parameters, as convergence on the maximum likelihood solution is not guaranteed for either package. However, overall these differences appear to be negligible, as parameter estimates for $\Sigma$ and $\alpha$ were very similar (Table 1).

## Conclusion

We have implemented a fast linear-time algorithm in `Rphylopars` which allows for the estimation of phylogenetic trait covariance for data sets with missing observations and multiple within-species observations. The methods described here extend the original `PhyloPars` implementation in an R environment to allow incorporation of within-species (phenotypic) correlations and alternative evolutionary models. As comparative data sets continue to grow in number of species and traits observed, fast methods will become increasingly critical. `Rphylopars` is available on the CRAN repository (https://cran.r-project.org/web/packages/Rphylopars/), as well as on GitHub (https://github.com/ericgoolsby/Rphylopars). A tutorial with worked examples is provided in Appendix S3 for implementing the features described here, and additional information can be found on the `Rphylopars` wiki (https://github.com/ericgoolsby/Rphylopars/wiki).

## Data accessibility

This study does not use data.

## References

Adams, D.C. (2014a) A generalized K statistic for estimating phylogenetic signal from shape and other high-dimensional multivariate data. *Systematic Biology*, **63**, 685–697.

Adams, D.C. (2014b) A method for assessing phylogenetic least squares models for shape and other high-dimensional multivariate data. *Evolution*, **68**, 2675–2688.

Arbuckle, J.L., Marcoulides, G.A. & Schumacker, R.E. (1996) Full information estimation in the presence of incomplete data. *Advanced Structural Equation Modeling: Issues and Techniques*, **243**, 277.

Bartoszek, K., Pienaar, J., Mostad, P., Andersson, S. & Hansen, T.F. (2012) A phylogenetic comparative method for studying multivariate adaptation. *Journal of Theoretical Biology*, **314**, 204–215.

Bruggeman, J., Heringa, J. & Brandt, B.W. (2009) PhyloPars: estimation of missing parameter values using phylogeny. *Nucleic Acids Research*, **37**, W179–W184.

Butler, M.A. & King, A.A. (2004) Phylogenetic comparative analysis: a modeling approach for adaptive evolution. *The American Naturalist*, **164**, 683–695.

Clavel, J., Escarguel, G. & Merceron, G. (2015) mvMORPH: an R package for fitting multivariate evolutionary models to morphometric data. *Methods in Ecology and Evolution*, **6**, 1311–1319.

Cressie, N. (2015) *Statistics for Spatial Data*. Wiley-Interscience, New York.

Dempster, A.P., Laird, N.M. & Rubin, D.B. (1977) Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, **39**, 1–38.

Denton, J.S. & Adams, D.C. (2015) A new phylogenetic test for comparing multiple high-dimensional evolutionary rates suggests interplay of evolutionary rates and modularity in lanternfishes (Myctophiformes; Myctophidae). *Evolution*, **69**, 2425–2440.

Felsenstein, J. (1973) Maximum-likelihood estimation of evolutionary trees from continuous characters. *American Journal of Human Genetics*, **25**, 471–492.

Felsenstein, J. (1985) Phylogenies and the comparative method. *The American Naturalist*, **125**, 1–15.

Felsenstein, J. (2008) Comparative methods with sampling error and within-species variation: contrasts revisited and revised. *The American Naturalist*, **171**, 713–725.

FitzJohn, R.G. (2012) Diversitree: comparative phylogenetic analyses of diversification in R. *Methods in Ecology and Evolution*, **3**, 1084–1092.

FitzJohn, R.G., Pennell, M.W., Zanne, A.E., Stevens, P.F., Tank, D.C. & Cornwell, W.K. (2014) How much of the world is woody? *Journal of Ecology*, **102**, 1266–1272.

Freckleton, R.P. (2012) Fast likelihood calculations for comparative analyses. *Methods in Ecology and Evolution*, **3**, 940–947.

Hadfield, J.D. & Nakagawa, S. (2010) General quantitative genetic methods for comparative biology: phylogenies, taxonomies and multi-trait models for continuous and categorical characters. *Journal of Evolutionary Biology*, **23**, 494–508.

Hansen, T.F. (1997) Stabilizing selection and the comparative analysis of adaptation. *Evolution*, **51**, 1341–1351.

Hansen, T.F. & Bartoszek, K. (2012) Interpreting the evolutionary regression: the interplay between observational and biological errors in phylogenetic comparative studies. *Systematic Biology*, **61**, 413–425.

Harmon, L.J., Losos, J.B., Jonathan Davies, T., Gillespie, R.G., Gittleman, J.L., Bryan Jennings, W. *et al.* (2010) Early bursts of body size and shape evolution are rare in comparative data. *Evolution*, **64**, 2385–2396.

Ho, L.S.T. & Ané, C. (2014) A linear-time algorithm for Gaussian and non-Gaussian trait evolution models. *Systematic Biology*, **63**, 397–408.

Housworth, E.A., Martins, E.P. & Lynch, M. (2004) The phylogenetic mixed model. *The American Naturalist*, **163**, 84–96.

Ives, A.R., Midford, P.E. & Garland, T. (2007) Within-species variation and measurement error in phylogenetic comparative methods. *Systematic Biology*, **56**, 252–270.

Martins, E.P. & Hansen, T.F. (1997) Phylogenies and the comparative method: a general approach to incorporating phylogenetic information into the analysis of interspecific data. *The American Naturalist*, **149**, 646–667.

Pagel, M. (1997) Inferring evolutionary processes from phylogenies. *Zoologica Scripta*, **26**, 331–348.

Pagel, M. (1999) Inferring the historical patterns of biological evolution. *Nature*, **401**, 877–884.

Pakeman, R.J. (2014) Functional trait metrics are sensitive to the completeness of the species' trait data? *Methods in Ecology and Evolution*, **5**, 9–15.

Pinheiro, J.C. & Bates, D.M. (1996) Unconstrained parametrizations for variance-covariance matrices. *Statistics and Computing*, **6**, 289–296.

Smith, S.A., Beaulieu, J.M., Stamatakis, A. & Donoghue, M.J. (2011) Understanding angiosperm diversification using small and large phylogenetic trees. *American Journal of Botany*, **98**, 404–414.

## Supporting Information

Additional Supporting Information may be found online in the supporting information tab for this article:

**Appendix S1.** Description of the linear-time algorithm.

**Appendix S2.** R code used for simulations and benchmark comparisons in Figure 1 and Table 1.

**Appendix S3.** Tutorial with worked examples.