1    **Permutation tests for comparative data**

2    James G. Saulsbury[1]

3    [1]University of Michigan, Ann Arbor

4

5    **Abstract**

6    The analysis of patterns in comparative data has come to be dominated by least-squares

7    regression, mainly as implemented in phylogenetic generalized least-squares (PGLS). This

8    approach has two main drawbacks: it makes relatively restrictive assumptions about distributions

9    and can only address questions about the conditional mean of one variable as a function of other

10   variables. Here I introduce two new non-parametric constructs for the analysis of a broader range

11   of comparative questions: phylogenetic permutation tests, based on cyclic permutations and

12   permutations conserving phylogenetic signal. The cyclic permutation test, an extension of the

13   restricted permutation test that performs exchanges by rotating nodes on the phylogeny, performs

14   well within and outside the bounds where PGLS is applicable but can only be used for balanced

15   trees. The signal-based permutation test has identical statistical properties and works with all

16   trees. The statistical performance of these tests compares favorably with independent contrasts

17   and surpasses that of a previously developed permutation test that exchanges closely related pairs

18   of observations more frequently. Three case studies illustrate the use of phylogenetic

19   permutations for quantile regression with non-normal and heteroscedastic data, testing

20   hypotheses about morphospace occupation, and comparative problems in which the data points

21   are not tips in the phylogeny.

22

23   **Introduction**

24   For a biologist interested in the role of natural selection in evolution, questions about

25   relative trait values are easier to address than questions about absolute trait values. For example,

26   "do bears from colder climates have longer fur" is far more analytically tractable than "is long

27   hair an adaptation for cold climates," even if the latter is the original question of interest (Sober

28   and Orzack 2003). Comparative or cross-species data are a fruitful source of insights into how

29   natural selection works in populations, and also into broad-scale phenomena that are interesting

30   in themselves, but their analysis is non-trivial. Comparative data often carry a detectable signal

31   of the phylogenies on which they evolved, and covariation between the trait values of close

32    relatives can cause serious problems for a statistical analysis, most conspicuously in the form of

33    inflated false positive rates (Felsenstein 1985). The dominant paradigm for the past several

34    decades of comparative research was established by Felsenstein (1985), who showed that the

35    independent values in a comparative analysis are not the trait states at the tips of a phylogeny but

36    their divergences (or contrasts) at phylogenetic splits. Unlike the raw trait values, these

37    "phylogenetically independent contrasts" (PICs) can be safely analyzed with least-squares

38    regression. Phylogenetic generalized least squares (PGLS; Grafen 1989) was developed as a

39    more general comparative framework that can accommodate non-linear relationships via link

40    functions (for example, phylogenetic logistic regression; Ives and Garland 2010), trees with

41    polytomies, and a variety of evolutionary models. PGLS is a kind of generalized least squares

42    regression that uses a phylogeny as the variance-covariance matrix, and it returns identical

43    results as the PIC approach in its simplest form.

44    PICs/PGLS have enjoyed immense success as a framework for understanding

45    relationships among traits in comparative data while accounting for phylogenetic autocorrelation

46    (Symonds and Blomberg 2014), but they have two chief limitations. First, as regression tests

47    they are assumption-rich: their reliability depends on, among other things, the residuals being

48    normally distributed and homoscedastic (equal variance across the values of the predictors)

49    (Mundry 2014). The other limitation is that least-squares regression is a rather specific analytical

50    framework: questions about the relationship between one or more variables and the conditional

51    mean of another variable occupy only a small corner of the universe of biologically interesting

52    comparative problems. This has pernicious implications for the use of phylogenetic regression as

53    the "go-to" method among comparative biologists. In the last section of this paper I highlight

54    three examples of comparative problems that are off-limits to PGLS: quantile regression,

55    morphospace occupation, and ecogeographic rules. PICs/PGLS is quite powerful within rather

56    circumscribed bounds (Orzack and Sober 2001), but methods that promote the creative

57    exploration of questions and datasets outside those bounds can only improve comparative

58    biology. Here I develop phylogenetically informed permutation tests, validate them with toy

59    scenarios, and illustrate their use with empirical case studies. The first of these tests generates a

60    set of nulls using cyclic permutations and is a conceptually straightforward extension of an

61    existing test, but it can only be used with balanced trees. The second conserves the phylogenetic

62  signal in the data, has identical statistical properties to the first, and can be used with any
63  phylogeny.

64

## Phylogenetic permutation tests

66  Permutation testing is widespread for some questions in evolutionary biology – for
67  example, in tests of phylogenetic signal (Blomberg et al. 2003) – but strangely has not permeated
68  to comparative testing (with one important exception, discussed below). The gist of a
69  permutation test is to take as a null distribution the set of test statistics associated with every
70  unique rearrangement (or permutation) of the data and to compare the empirical test statistic with
71  this distribution (Good 2000). In each of these permutations the "labels" on at least one variable
72  in the dataset are randomly rearranged, breaking the empirical association between variables
73  without changing their distributions. The proportion of permutations in which the test statistic is
74  at least as extreme as the empirical one is taken as the probability of obtaining the results under
75  the null hypothesis: the p-value (Perezgonzalez 2015). For example, the subject in Fisher's
76  famous "Lady Tasting Tea" experiment correctly guessed the method of preparation for 8 cups
77  of tea, and Fisher used permutations of the order of guesses to determine that guessing randomly
78  would have achieved this result with a low probability of $p = 1/70$. This example was simple
79  enough for every possible permutation to be enumerated analytically, but for bigger datasets this
80  can be computationally infeasible, so in practice the distribution is typically estimated by
81  permuting randomly many times (Good 2000). Many "flavors" of permutation test have been
82  developed, differing mainly in the null distribution they generate (Anderson 2001). The primary
83  virtue of the permutation test is its elegance: unlike parametric tests, it does not rely on
84  theoretical probability distributions (the population of interest is the empirical one), and it can be
85  used with a broader range of test statistics (Good 2000).

86  Despite its strengths, the ordinary permutation test cannot be applied to data that evolved
87  on phylogenies. The test is, in a sense, distribution-free, but not assumption-free: permutation
88  tests assume among other things that the observations being shuffled are *exchangeable*, meaning
89  that rearrangements of those observations have the same joint probability distribution (Anderson
90  2001). This is quite close to the assumption in least-squares regression that variables are
91  independent and identically distributed, and these assumptions are violated by the complex
92  covariance structure of comparative data. In other words, because the traits of closely related

93   taxa tend to covary due to their shared evolutionary history, comparative data are not

94   exchangeable. However, a modified test that uses phylogenetic information to preserve

95   exchangeability can be used for sound non-parametric hypothesis-testing with comparative data.

*Lapointe-Garland phylogenetic permutations*

97   Lapointe and Garland (2001) proposed a permutation test for comparative data in which

98   pairs of values at the tips are exchanged with probability proportional to their phylogenetic

99   proximity, such that the most probable exchange is between a trait value and itself (Fig. 1A).

100   This approach uses a relatedness matrix which can be "flattened out" using a parameter k; for

101   values of k higher than one the test approaches an ordinary permutation test. This was the first

102   and apparently the only previous attempt at developing a comparative permutation test.

103   Although the test is less vulnerable than the ordinary permutation test to phylogenetically

104   induced false positives, it has several undesirable properties. One of these is the high rate of

105   exchanges between an observation and itself (Fig. 1A), or auto-exchanges, which results in a set

106   of permutations that is tightly constrained around the empirical statistic (Appendix 1). This

107   should reduce statistical power because the permutations look so much like the empirical

108   arrangement. In this respect the Lapointe-Garland (LG) test also strays from one of the essential

109   features of permutation tests: enumerating each unique rearrangement of the data. The high rate

110   of auto-exchanges has the effect of up-weighting some possible rearrangements over others; it is

111   not clear why it would be desirable to give more weight to rearrangements that look more like

112   the empirical one. The rate of auto-exchanges can be dampened by increasing the value of k,

113   which makes the exchange matrix flatter, but that defeats the point of incorporating phylogeny.

114   Moreover, there is no principled way to choose a value of k above one, nor is there a clear use-

115   case for permutations that are only partially informed by phylogeny. Another theoretical problem

116   with LG permutations is that they do not conserve phylogenetic signal, the key feature that

117   makes interpreting comparative datasets difficult. For the leftmost tree in Fig. 1C, the

118   phylogenetic signal of LG permutations varies by a factor of 2.3.

119   Simulations show these features of the LG approach have consequences for its statistical

120   performance. For instance, one desirable feature of a significance test is that the rate of false

121   positives should be exactly equal to the significance level: thus, 5% of cases in which the null

122   hypothesis is true should have $p < 0.05$. I tested the false positive rate of the LG permutation test

123   by simulating uncorrelated Brownian Motion evolution of two continuous traits on a rooted 8-

124    taxon tree with two polytomies of four tips each. I computed the absolute correlation coefficient

125    between the two simulated traits and tested it with an LG permutation test (500 permutations) for

126    each of 1000 pairs of simulated traits (Fig. 2). I ran the same test with independent contrasts,

127    after making the tree amenable to PICs by making the two 4-taxon polytomies in the tree into

128    pectinate subtrees with added branches having length 0. Unlike independent contrasts, LG

129    permutation tests yielded non-uniformly distributed p-values, returning intermediate values most

130    frequently (Kolmogorov-Smirnov test against a uniform distribution, $p = 0.011$). In other words,

131    the p-values from this test do not tell the user what they are supposed to: the probability of

132    observing a statistic at least as extreme as the empirical one under the null hypothesis. I also

133    evaluated the false negative rate with the same procedure, except the two simulated traits were

134    correlated with an evolutionary covariance of 0.75. LG permutation tests return false negatives at

135    higher rates than independent contrasts: $p < 0.05$ for 469 of 1000 simulations of truly correlated

136    evolution, compared with 564. The LG phylogenetic permutation test is a valuable and

137    interesting non-parametric approach to comparative data, but, motivated by the conceptual and

138    statistical problems outlined here, I develop two new phylogenetically informed permutation

139    tests. The cyclic permutation test is a conceptually straightforward extension of an existing class

140    of permutation tests that can only be used with balanced trees; the signal-based permutation test

141    has identical statistical properties to the first and can be used with any phylogeny.

*Cyclic permutations*

143    An elegant solution to the problem of relatedness in comparative data can be found in the

144    restricted permutation test, in which rearrangements are restricted to only occur between

145    exchangeable data points or sets of data points (Anderson 2001). As a non-phylogenetic

146    example, an investigator testing the significance of a correlation between environmental

147    variables sampled in different regions might consider permuting only within regions and not

148    across them, especially if those variables were spatially autocorrelated by region. The resulting

149    restricted permutations would retain the same kind of spatial autocorrelation as the empirical

150    data. In a comparative dataset, the exchangeable units are not the values at the tips but the

151    descendants of each node in the tree. This is similar to Felsenstein's (1985) observation that

152    contrasts at nodes rather than tip values are independent of one another, and similar also to the

153    "radiation principle" that motivates Grafen's (1989) PGLS. A set of phylogenetically informed

154    permutations can therefore be generated with cyclic permutations of the values at the tips; that is,

155    by randomly rotating the descendants of each internal node for at least one variable in the dataset

156    (Fig. 1B). A statistic calculated for the set of these permuted datasets can be compared with the

157    empirical one in what is here called a cyclic permutation test. Because the units being permuted

158    can either be tip values (for the shallowest internal nodes) or sets of tip values (for deeper nodes),

159    this is a form of hierarchical restricted permutation.

160    The cyclic permutation test performs at least as well as independent contrasts and lacks

161    the statistical issues of the LG permutation test. In the test for false positives (Fig. 2), the set of

162    p-values is indistinguishable from a uniform distribution (Kolmogorov-Smirnov test, $p = 0.370$),

163    which is ideal. In the test for false negatives with cyclic permutations, p was below 0.05 for

164    702/1000 simulations, corresponding to a false negative rate of around 30% (Fig. 2).

165    Interestingly, this is a better rate than what independent contrasts recovered (p below 0.05 for

166    564/1000 simulations). This indicates that the cyclic permutation test and PICs have at least

167    comparable statistical power, even though the former is a non-parametric test.

168    Cyclic permutations will change an unbalanced tree's two-dimensional projection, so the

169    cyclic permutation test can only be used with a topologically balanced tree. If a trait is permuted

170    cyclically on an unbalanced tree, it will no longer share the same evolutionary history as other

171    traits in the dataset, and it defeats the point of the test – namely, to ask what kinds of patterns can

172    result from the independent evolution of different traits on the same phylogeny. Because of the

173    restriction to balanced trees, the cyclic permutation test cannot be used with most empirical

174    datasets. However, because it is so conceptually straightforward and because it works (Fig. 2), it

175    is a useful yardstick against which to measure another new approach in which permutations

176    conserve the amount of phylogenetic signal in the data.

177    *Signal-based permutations*

178    The following permutation test can be used with real phylogenies: compare an empirical

179    test statistic with the set of permutations in which phylogenetic signal is equal or sufficiently

180    close to the empirical signal (Fig. 1C). The logic here is that the only rearrangements that can be

181    meaningfully compared with empirical data are those in which trait values are just as conserved

182    on, or structured by, the phylogeny. Phylogenetic signal is quantified here with Moran's I rather

183    than another metric like Pagel's λ or Blomberg's K in the non-parametric spirit of the

184    permutation test: whereas those other metrics explicitly model the evolutionary process that

185    generated a given trait, Moran's I  simply quantifies the degree to which the trait values of

186    closely-related species covary (Gittleman and Kot 1990; Appendix 2). I implement signal-based

187    permutation with a simple hill-climbing algorithm in which first the values of a trait are shuffled,

188    then randomly-selected pairs of observations are swapped if doing so brings the phylogenetic

189    signal closer to the empirical signal, and the procedure stops when the permuted phylogenetic

190    signal is within some specified tolerance of the empirical value. The test could be implemented

191    without this hill-climbing procedure, but it would make the test extraordinarily time-consuming

192    for some datasets. Because phylogenetic signal depends on the values at the tips, the

193    rearrangements that are included in the set of signal-based permutations depend on the values of

194    the trait being permuted, unlike cyclic permutations and the LG permutation test. For applicable

195    trees, the set of signal-based permutations is always at least as inclusive as the set of cyclic

196    permutations: every cyclic permutation has identical phylogenetic signal, but non-cyclic

197    permutations of a dataset can too (Fig 1B, rightmost permutation), and additional rearrangements

198    can be accepted if the specified tolerance is large enough. For example, there are 2^(number of

199    internal nodes) = 128 possible cyclic permutations of the 8-taxon tree in Figure 1B but 256

200    permutations with identical phylogenetic signal. The positions of clades (C,D) and (G,H) are

201    switched in the 128 non-cyclic permutations.

202         Despite these striking differences from the cyclic permutation test, simulations show that

203    signal-based and cyclic permutations have apparently identical statistical properties. Like the

204    other test, the signal-based test correctly returns a uniform distribution of p values for 1000

205    simulations of uncorrelated evolution (Fig. 2, "Signal-based permutation"; Kolmogorov-Smirnov

206    test against a uniform distribution, $p = 0.413$). The false negative rate is also comparable with

207    that for the cyclic permutation test (Fig. 2; 716/1000 p-values below 0.05), and higher than that

208    for PICs. Thus, the cyclic and signal-based permutation tests do not have the problems with

209    statistical power and size that characterize the Lapointe-Garland test.

210         As a visual illustration of the two new phylogenetic permutation tests, consider their

211    application to Felsenstein's (1985) "worst case scenario" in which uncorrelated Brownian

212    Motion evolution of two traits on a rooted tree of two polytomies with 20 tips each (all branch

213    lengths equal) generates a spurious correlation among traits (Fig. 3A). An ordinary permutation

214    test yields a distribution of mainly low absolute correlation coefficients (Fig. 3B) and a very high

215    level of significance ($p < 0.001$). The investigator who makes the mistake of treating all tip

216    values as exchangeable incorrectly rejects the null hypothesis of independent evolution.

217   Conversely, the distribution of correlation coefficients for 1000 cyclic permutations is centered

218   close to the empirical correlation coefficient (Fig. 3B), yielding a p-value of 0.31. Because of the

219   clustering of trait values within subclades, every cyclic permutation preserves a relatively strong

220   correlation coefficient: 95% of the cyclic permutations have |r| between 0.34 and 0.65. The null

221   distribution generated by signal-based permutations depends on the tolerance: a set of 1000

222   signal-based permutations with the broadest possible tolerance (2) is statistically

223   indistinguishable from an ordinary permutation test (Kolmogorov-Smirnov test, p = 0.7226)

224   because the phylogenetic signal of every possible permutation is within its tolerance (Fig. 3B).

225   For smaller tolerances, the distribution of test statistics for signal-based    permutations    more

226   closely approximates the set of cyclic permutations (Fig. 3B), such that with a margin of 0.01

227   (Moran's I of permuted variable Y between 0.512 and 0.532) they are statistically

228   indistinguishable (p=0.536). Thus, signal-based permutations converge on the statistical

229   properties of cyclic permutations.

230   Interestingly, phylogenetic permutation tests succeed in a case where PICs and PGLS

231   both fail: a second "worst case" constructed by Uyeda et al. (2018). In this scenario, simulated

232   traits evolve in the same way and on the same phylogeny as in Felsenstein's worst case, but with

233   one modification: a single extreme shift in both traits near the root generates a contrast that is a

234   strong enough outlier to make the two traits appear significantly associated, even when

235   "correcting for phylogeny." PICs/PGLS incorrectly recover significant relationships between

236   traits because these methods are parametric, and their assumptions are violated by the dramatic

237   outlier. The cyclic permutation test is unburdened by these assumptions: the extreme outlier is

238   incorporated into every permutation, and the test correctly yields a non-significant result

239   (Appendix 3). Likewise, the only rearrangements of the data that conserve phylogenetic signal

240   are those in which exchanges only occur within clades and not between them, so a signal-based

241   permutation test succeeds in the same way. Cyclic and signal-based permutations both represent

242   reasonable null models against which to compare empirical patterns, but only the latter is

243   applicable to real trees, so I use signal-based permutation tests to explore the following case

244   studies.

245

246                                          **Case studies**

247    The preceding sections established that phylogenetic permutations perform at least as
248  favorably as PICs in "toy scenarios" in which the truth is known. These scenarios involved
249  modeled BM evolution of traits with normal distributions, and the only test statistic considered
250  was the correlation between two traits. PICs/PGLS perform comfortably within these bounds. In
251  the following case studies, I use phylogenetic permutation to explore scientific questions that are
252  effectively off-limits to PGLS-type methods because they involve strange distributions and test
253  statistics beyond the least-squares regression framework. The first case study involves quantile
254  regression on a heteroskedastic dataset with a non-normal response variable. The second
255  explores the statistical significance of patterns of morphospace occupation. The third tests an
256  ecogeographic rule: the data points are not tips in the phylogeny but the aggregate property of all
257  the tips that occur in each geographic area.

258    The statistical significance of some of these findings could potentially be tested by
259  comparing empirical statistics with null simulations rather than permutations, like what Mahler
260  et al. (2013) used to demonstrate exceptional convergence in anoles. However, this requires
261  assumptions about distributions and the evolutionary processes that generated a dataset which an
262  investigator may not want or be able to make. If a dataset exhibits a more extreme test statistic
263  than a set of simulations, is it because there was a mechanistic association between those traits,
264  or because the simulations were unrealistic? Such questions may be hard to answer and are
265  avoided by taking the non-parametric approach.

266    *Quantile regression and peculiar distributions: arm number in feather stars*

267    Saulsbury and Baumiller (2020) investigated a wedge-shaped relationship between
268  absolute latitude and arm number among feather stars, a group of suspension-feeding marine
269  echinoderms: species near the poles typically have around 10 arms, whereas those around the
270  equator have between 5 and 150. Arm number varies widely within many families, but across the
271  dataset it has a strange distribution, probably due to the unique and complex ontogeny of feather
272  star arms (Shibata and Oji 2003): about half the species in the dataset have exactly 10 arms, and
273  the rest of the distribution is markedly right skewed. More importantly, this non-normality also
274  characterizes the residuals in a PGLS regression of arm number, and log(arm number), on
275  absolute latitude. Another aspect of the dataset that poses obvious problems for least-squares
276  regression is also the dataset's most biologically interesting feature: arm number is
277  heteroskedastic across absolute latitude. Beyond these more technical challenges, questions

278    about the spread of a response variable as a function of a predictor cannot be readily addressed

279    with least-squares regression. Instead they are the purview of quantile regression, which

280    estimates quantiles (for example, the median, or the $10^{th}$ percentile) conditional on predictors.

281    There is not currently an equivalent to quantile regression in the PGLS framework. As such, the

282    authors used signal-based phylogenetic permutation tests to consider whether the latitudinal

283    gradient in arm number could have plausibly emerged through independent evolution on feather

284    star phylogeny.

285    Although both absolute latitude and arm number exhibit phylogenetic signal, and thus

286    might be prone to spurious associations, the empirical relationships between the two are more

287    extreme than almost all phylogenetic permutations. The $90^{th}$ and $95^{th}$ conditional percentiles,

288    which characterize how maximum arm number relates to latitude, were significantly negative (p

289    = 0.017 and 0.009, respectively), as was Spearman's rank-correlation coefficient (p < 0.001).

290    Concluding that the pattern could not be explained away as the result of random evolution, the

291    authors drew on ecological and functional morphological evidence to argue that a latitudinal

292    gradient in the intensity of predation represented the most plausible explanation for their

293    findings. This simple case study illustrates the value of a comparative method that makes

294    minimal assumptions about the distribution of the data. It also hints at the extent of the patterns

295    that can be evaluated with phylogenetic permutation, although that is more fully illustrated by

296    the following examples.

297    *Morphospace occupation: Triassic ammonoids*

298    Why are some theoretically possible morphologies not realized in nature, and why are

299    some realized more frequently than others? These questions are the domain of theoretical

300    morphology, a subdiscipline catapulted to the forefront of evolutionary biology for a time by

301    David Raup. He found (1966) that the breadth of shell morphologies realized by mollusks and

302    brachiopods was surprisingly well-summarized by a model in which a generating curve or whorl

303    increases in size as it revolves around an axis. Shell geometry is controlled by three parameters:

304    whorl expansion rate, translation of successive whorls along the axis, and the distance of

305    successive whorls from the shell axis. Interestingly, most theoretically possible combinations of

306    parameter values are not realized in nature; Raup cautiously submitted that either these

307    unrealized forms were physiologically impossible, or shell-building invertebrates simply had not

308    had time to reach those parts of morphospace yet. A companion paper (Raup 1967) focused on

ammonoids, an extinct group of mostly "planispiral" mollusks in which typically no whorl translation occurs and variation is constrained along two axes of theoretical shell morphospace: the distance of successive whorls from the axis (D), and the whorl expansion rate (W) (Fig. 4A). Again, much of the rectangle defined by ammonoid occupation in D-W space is unoccupied – for example, almost no ammonoids fall above the line $W = 1/D$ (Fig. 4A). Shells above this curve are open-coiled, making them, among other things, weaker and easier for a predator to crush. For shells under this curve, each whorl can incorporate part of the previous whorl in its construction, so open-coiled shells ($W > 1/D$) also waste the building materials they otherwise would have saved. Thus, the patterns in theoretical morphospace occupation are interesting because of the underlying fitness surface they suggest.

The problem with inferences of selective forces from the pattern of morphospace occupation is that they rely on the equilibrium assumption (Lauder 1982): namely, that the phenotypes under study are at equilibrium with the selective forces that act on them. The alternate explanation for un- or under-occupied regions of morphospace is that, by chance, ammonoids simply have not had time to reach those regions yet – in other words, the system is historical and not at equilibrium. Raup (1967) raised this possibility, but admitted that in order to make headway he had to "assume that the observed morphology has had, in evolution, a selective advantage over other possible morphologies." Subsequent studies have made the same assumption: for example, Tendler et al. (2015) tested whether ammonoids fill out a triangle in D-W space as a demonstration of Pareto optimality theory, which predicts that functional "archetypes" should form the vertices of polygons in trait space (Fig. 4A). They tested whether the ammonoid data are more triangular than the set of ordinary permutations, but this procedure incorrectly assumes that the data are exchangeable, or in other words that each data point obtained its morphology independently – a problem pointed out by Edelaar (2013) for another study of Pareto optimality. The equilibrium assumption leaves comparative studies vulnerable to the kinds of false positives discussed by Felsenstein (1985) in which a pattern apparently supported by a high number of replicates actually only represents a few evolutionary events. Theoretical morphology has not been incorporated with phylogeny in the way other comparative subdisciplines have in recent decades. However, it is not amenable to PGLS because it is not a regression problem: the question is not about the conditional mean of a response variable but about why certain combinations of traits are unrealized.

340  The phylogenetic permutation approach is a promising way forward for theoretical
341  morphology because it can be used to ask what kinds of patterns in morphospace occupation can
342  emerge without any dependence between traits. If empirical patterns fall outside the range of
343  phylogenetic permutations, more interesting evolutionary explanations for the pattern in
344  morphospace occupation can be explored – for example, certain morphologies could be
345  unrealized because they are less fit. No broad-scale phylogeny of ammonoids is available, so I
346  used taxonomy as a polytomy-rich phylogeny to explore morphospace occupation in the database
347  of 322 Triassic ammonoid genera from McGowan (2004). These genera belong to 79 families in
348  18 superfamilies. This is a very coarse way to approximate phylogeny, so this exploration should
349  be taken as a proof of concept and a hint at the role of contingency in ammonoid evolution.

350  I used ordinary and signal-based phylogenetic permutations to test the significance of two
351  test statistics: the number of genera over the W=1/D line (6/322 genera), and the triangularity of
352  the dataset in D-W space, defined as the ratio of the area of the convex hull to the area of the
353  smallest triangle that encloses all the data (triangularity = 0.8535; Fig. 4A; Appendix 4). Both D
354  and W have low signal on the "phylogeny", with values of Moran's I of 0.072 and 0.049,
355  respectively. So, inasmuch as ammonoid taxonomy approximates phylogeny, the various
356  ammonoid clades appear to have independently explored a lot of D-W space: for example, there
357  are five superfamilies that each occupy more than half the area of the total convex hull. In a
358  system characterized by this much exploration of morphospace, it seems unlikely that
359  particularly strong patterns could emerge from random chance alone. The phylogenetic
360  permutation test quantifies this preliminary impression: $p < 0.001$ for both test statistics for both
361  ordinary and phylogenetic permutations (Fig. 4B-C). In other words, all phylogenetic
362  permutations of the data have more open-coiled genera and are less triangular than the empirical
363  dataset. Considering phylogeny (that is, going from ordinary to phylogenetic permutations) does
364  not visibly affect the null distribution for the first test statistic; it does slightly for triangularity,
365  shifting it to the right. So, because of phylogeny there is a slight tendency for permuted datasets
366  to look more triangular, but not enough to make a difference for the p-value.

367  Thus, the independent evolution of shell growth parameters D and W constitutes a poor
368  explanation for both the triangularity of the dataset and the paucity of open-coiled genera. One
369  could easily imagine a hypothetical phylogenetic history for which a non-significant result would
370  be obtained – for example, if every ammonoid with $D > 0.3$ were part of the same clade, it would

371     be easier to explain away the pattern of morphospace occupation as a historical accident. The

372     phylogenetic permutation test is well-suited for this problem because characterizing the

373     biologically interesting features of morphospace occupation often requires the use of creative or

374     novel statistics. Note that this is not the only way to evaluate the evolutionary "significance" of

375     morphospace occupation: Tendler et al. (2015) showed that ammonoids refilled roughly the same

376     region of morphospace several times after mass extinctions, representing semi-independent

377     replicates. In the final case study, I explore a comparative problem in which the data points are

378     not tips in the phylogeny.

379               *Ecogeographic rules: Thorson's rule in muricid gastropods*

380        Some of the most productive hypotheses in biology predict the way some biological

381     feature changes across space. Well-known examples of "ecogeographic rules" like these include

382     Bergmann's rule, the tendency for endotherms to be larger toward the poles (Olalla-Tárraga

383     2011), and Rapoport's rule, the putative tendency for species' latitudinal ranges to be smaller in

384     the tropics (Stevens 1989). The analysis of ecogeographic rules entails an interesting and under-

385     researched problem: species typically exist in more than one place, rather than at a single point

386     as in other kinds of comparative studies. It might seem that a straightforward comparative

387     analysis could address this by using a summary statistic of the range of each species, such as the

388     range midpoint, and indeed many studies take this shortcut. However, such an approach removes

389     biological information and is susceptible to false positives, especially if the trait in question

390     corresponds with range size in some way (Saulsbury and Baumiller 2020; Colwell and Hurtt

391     1994). In the most well-known and straightforward example, a test for a relationship between

392     absolute latitudinal midpoint and range size tends to recover strong negative relationships even

393     none really exists: geometrically, large ranges cannot be centered at high latitude, so these taxa

394     have their latitudinal midpoints "pulled" toward the equator (Colwell and Hurtt 1994). An

395     alternative approach is to consider the ecogeographical data as such – that is, as a set of places

396     and the aggregate properties of all the species in each place (Stevens 1989) – but this is

397     analytically fraught as well. Such data are beset not only by the phylogenetic autocorrelation that

398     complicates other comparative studies, but also by spatial autocorrelation to the degree that

399     species occur in multiple places (Rohde et al. 1993). Here I show that both the phylogenetic

400     permutation test can circumvent both sources of autocorrelation using a case study of larval

401     development across latitude in muricid gastropods.

402    Thorson's rule predicts that the larvae of marine invertebrates near the equator are more
403    likely to be planktotrophs – feeding larvae that persist in the water column for a long time –
404    whereas toward the poles there should be a predominance of non-feeding larvae, including direct
405    developers and lecithotrophs (yolk-supplied larvae) (Thorson 1950). Thorson proposed that
406    vulnerable planktotrophic larvae would not be able to cope with the extreme conditions and
407    variable food supply at high latitudes, but this mechanism and the latitudinal pattern have
408    subsequently received mixed empirical support (Marshall et al. 2012). Yet the idea persists:
409    Pappalardo et al. (2014) claimed support for Thorson's rule in a dataset of 44 muricid gastropod
410    species (Fig 5A). A logistic PGLS regression of larval development (planktotrophic vs. non-
411    feeding) on sea surface temperature [taken either from a single confirmed occurrence (69%) or
412    from the latitudinal midpoint of each species (31%)] recovered marginally significant
413    relationships: $p = 0.087$ for the regression of feeding (planktotrophic) vs. non-feeding mode on
414    temperature, and $p = 0.045$ for the regression of pelagic (planktotrophic and lecithotrophic) vs.
415    non-pelagic mode on temperature. Analyzing the same dataset, I found a similar degree of
416    support in a PGLS logistic regression of larval development on latitudinal midpoints (Appendix
417    5). Other recent studies of Thorson's rule use latitudinal or environmental midpoints as well
418    (Ibáñez et al. 2018; Ewers-Saucedo and Pappalardo 2019), presumably because the PGLS
419    framework requires it. Notably this seems to be a recent development, as Thorson and others
420    who worked on this problem since were mostly considering the proportion of planktotrophic
421    species at each latitude (Thorson 1950; Mileikovsky 1971; Jablonski and Lutz 1983; Collin
422    2003). Importantly, the use of midpoints can be vulnerable to complications involving range
423    size: for example, if planktotrophic species have larger ranges, it would artificially strengthen the
424    relationship between latitude and development by dragging the latitudinal midpoints of wide-
425    ranging species toward the equator (Colwell and Hurtt 1994). In fact, the broad geographic
426    ranges of feeding larvae are famous among invertebrate zoologists (Jablonski 1986), and the
427    median latitudinal range of planktotrophic species in the muricid dataset is 3.25 times that of
428    non-planktotrophs (Fig 5A). Using a latitude or temperature value selected randomly from the
429    range might not be biased like the midpoint method is, but is not an ideal solution because it
430    removes information and adds noise.

431    If instead the ecogeographic data are considered as such – for example, with a plot of the
432    percentage of species with planktotrophic larvae in each 1° latitudinal bin – the trend is still

433 apparent, with a strong correlation of r = 0.927 (Fig. 5B). A non-phylogenetic significance test

434 that nevertheless accounts for spatial autocorrelation can be performed by permuting modes of

435 larval development randomly across the tips of the phylogeny and re-computing the correlation

436 coefficient (Fig. 5C). The resulting absolute correlation coefficients are spread evenly between 0

437 and 1, yielding marginal statistical significance with a p-value of 0.033. We can take both spatial

438 and phylogenetic autocorrelation into account with signal-based phylogenetic permutations of

439 mode of larval development. Phylogenetic signal of planktotrophic vs. non-planktotrophic

440 development is high (Moran's I = 0.60), and indeed larval development only appears to have

441 transitioned on the phylogeny a few times, providing an investigator with very low sample size:

442 the most parsimonious history of development involves only three transitions to or away from

443 planktotrophy. Accordingly, almost all phylogenetic permutations have high absolute correlation

444 coefficients, from which the empirical correlation is statistically indistinguishable (p = 0.387).

445 Thus, the muricid dataset cannot provide strong evidence against the completely independent

446 evolution of latitude and larval development. Notably, the authors focused on temperature not

447 latitude; it is unclear if a similarly non-significant result would be obtained for the correlation

448 between temperature and larval development, but temperature and latitude are closely correlated,

449 and the same analytical problem applies because species occupy a range of temperatures.

450 Ecogeographic data present an interesting challenge to the comparative biologist because

451 the data points, cast most directly, do not represent tips in the phylogeny but the aggregate

452 properties of all the tips in the phylogeny that occur in each place. It might be possible to

453 consider such data in a PGLS framework, but it would require the specification of a rather

454 complex variance-covariance matrix. Crucially, this phylogenetic permutation test does not

455 provide evidence *against* Thorson's rule in this group. At the pattern level, the group is a clear

456 example of the rule, with a strong negative correlation between latitude and the proportion of

457 species with planktotrophic larvae. This trend probably has important implications for their

458 modern ecology and future evolution, because it predicts, for example, that low-latitude

459 planktotrophic species should be buffered against extinction by their broad ranges (Jablonski and

460 Lutz 1983; Jablonski 1986). However, the key point is that, with phylogenetic and spatial

461 autocorrelation this strong, such a trend could have easily arisen without any mechanistic

462 relationship between latitude and larval development. In fact, given the distribution of

463 phylogenetic permutations (Fig. 5C), it would be much more surprising to find no latitudinal

464 trend in larval development. This might explain why so many groups appear to follow the rule
465 (Ibáñez et al. 2018), especially since mode of larval development appears to evolve infrequently
466 among marine invertebrates (Collin 2004). It would require an exceptionally strong trend to
467 support a mechanistic Thorson's rule in a dataset like this one – or more plausibly, a different
468 kind of data. This might mean a clade in which larval development transitions more frequently,
469 or it might mean a different kind and scale of evolutionary repetition.

470

471 **Conclusions**

472 The regression-based approach to comparative biology has been hugely successful, but it
473 is also inflexible: it fails for strangely distributed response variables, but more importantly, the
474 range of questions it can address is limited. Permutation tests represent a powerful alternative
475 that performs well both within and outside the bounds where PGLS is applicable. Case studies
476 illustrate the use of phylogenetic permutations for pushing comparative methods to new places.

477 Rather than being a purely technical matter, the distinction between PGLS and
478 permutation-based approaches is underlain by a more substantive difference in attitude toward
479 comparative data. PGLS is typically described as a way to "correct for phylogeny" (Symonds
480 and Blomberg 2014). Other comparative methods take an even more direct approach by
481 transforming the data to "remove phylogenetic effects" (Stearns 1984; Cheverud et al. 1985;
482 Felsenstein 1985; Gittleman and Kot 1990). The implication is that comparative data have been
483 contaminated or affected by an agent called phylogeny, and that this contamination needs to be
484 isolated and removed before the real relationships in the data can be studied. It is a drawback of
485 these methods that they put the user at a remove from the raw data. Patterns in phylogenetically
486 autocorrelated data are also no less real than those in transformed data: biological phenomena
487 that could have arisen purely by chance, like Thorson's rule in some taxa, can nevertheless have
488 real and important consequences. Transformations and corrections also remove information and
489 limit the kinds of statistics and questions that can be applied to a dataset.

490 The phylogenetic permutation test is mostly unique among comparative methods in that it
491 treats the raw data as such. The test is subject to some of the same criticisms to which all
492 frequentist tests are subject, including that statistical significance tells an investigator nothing
493 about effect size (a reaction to the widespread conflation of "significance" with importance;
494 Dushoff et al. 2019). This is true, but for many biological phenomena including the case studies

495  discussed here, the most relevant effect size is arguably the empirical test statistic. Only six of

496  322 Triassic ammonoid genera have open-coiled shells; it is hard to imagine a more meaningful

497  phylogenetic transformation of this test statistic. The phylogenetic permutation framework,

498  which considers whether raw data look typical for cases of independent evolution, is in a way the

499  reverse of the reigning paradigm of transforming comparative data or their expected covariances

500  to fit into a regression analysis. Hopefully, these new approaches can help facilitate scientific

501  creativity among comparative biologists.

502                                    **Figures**



503

504  **Figure 1.** The three kinds of phylogenetic permutations discussed in this paper, with examples of

505  each kind on a balanced, rooted tree of 8 taxa A-H. In a test, each kind of permutation is applied

506  to at least one variable in the dataset, and a population of many such permuted datasets is

507  compared with the empirical arrangement. **1A.** The phylogenetic permutation approach

508  developed by Lapointe and Garland (2001). Note that many trait values do not change position

509  across permutations because the highest probability of exchange is between a trait and itself. **1B.**

510  Cyclic permutations: the set of permutations that can be generated by rotating about nodes in the

511 tree (double-sided arrows). **1C.** Signal-based permutations. These are more inclusive than cyclic

512 permutations: they include all possible cyclic permutations because the latter always conserves

513 phylogenetic signal, but also non-cyclic permutations that retain the same or nearly the same

514 signal (rightmost rearrangement). This is the only test in which the set of permutations depends
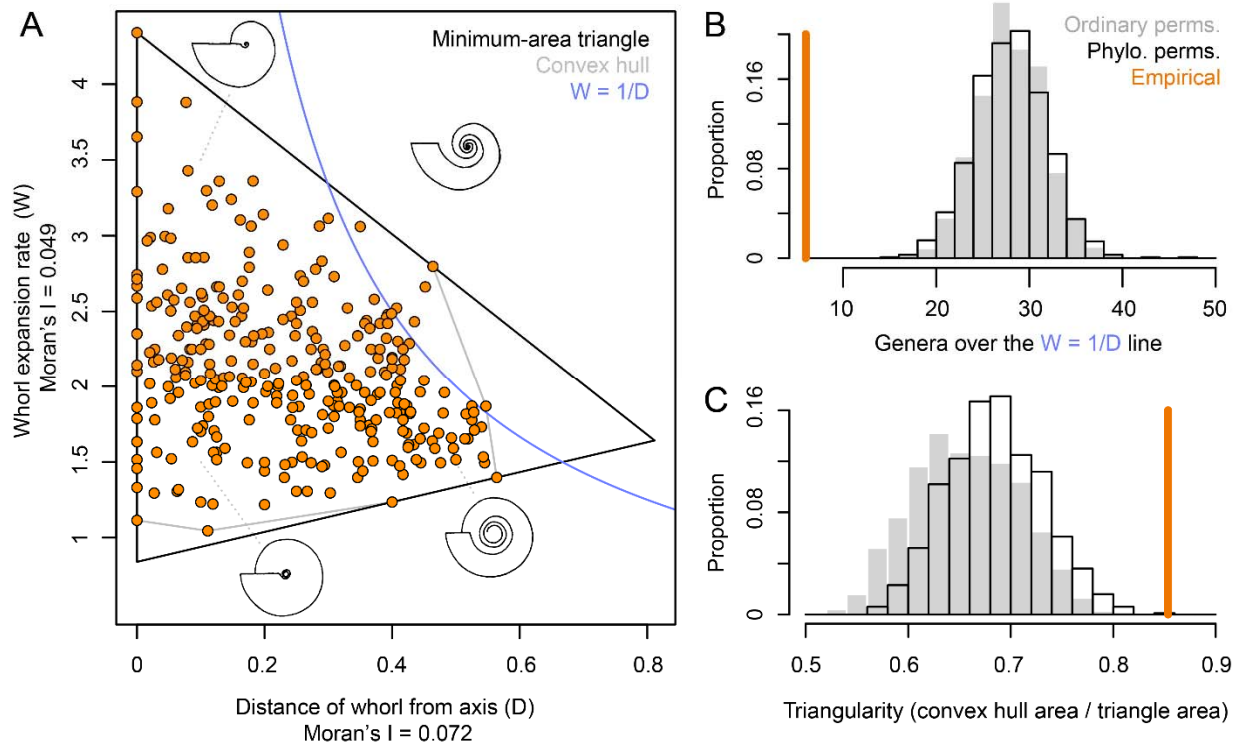
515 on the values at the tips.



516

517 **Figure 2.** p-values for three phylogenetic permutation tests of correlation and one test with

518 independent contrasts, applied to uncorrelated evolution of X and Y (above) and correlated

519 evolution with an evolutionary covariation of 0.75 (below). Traits simulated on a rooted 8-taxon

520 tree containing two polytomies with 4 taxa each, all branch lengths equal. p-values should ideally

521 be uniformly distributed for uncorrelated evolution and as low as possible for correlated

522 evolution. Red horizontal line indicates a uniform distribution; only Lapointe-Garland

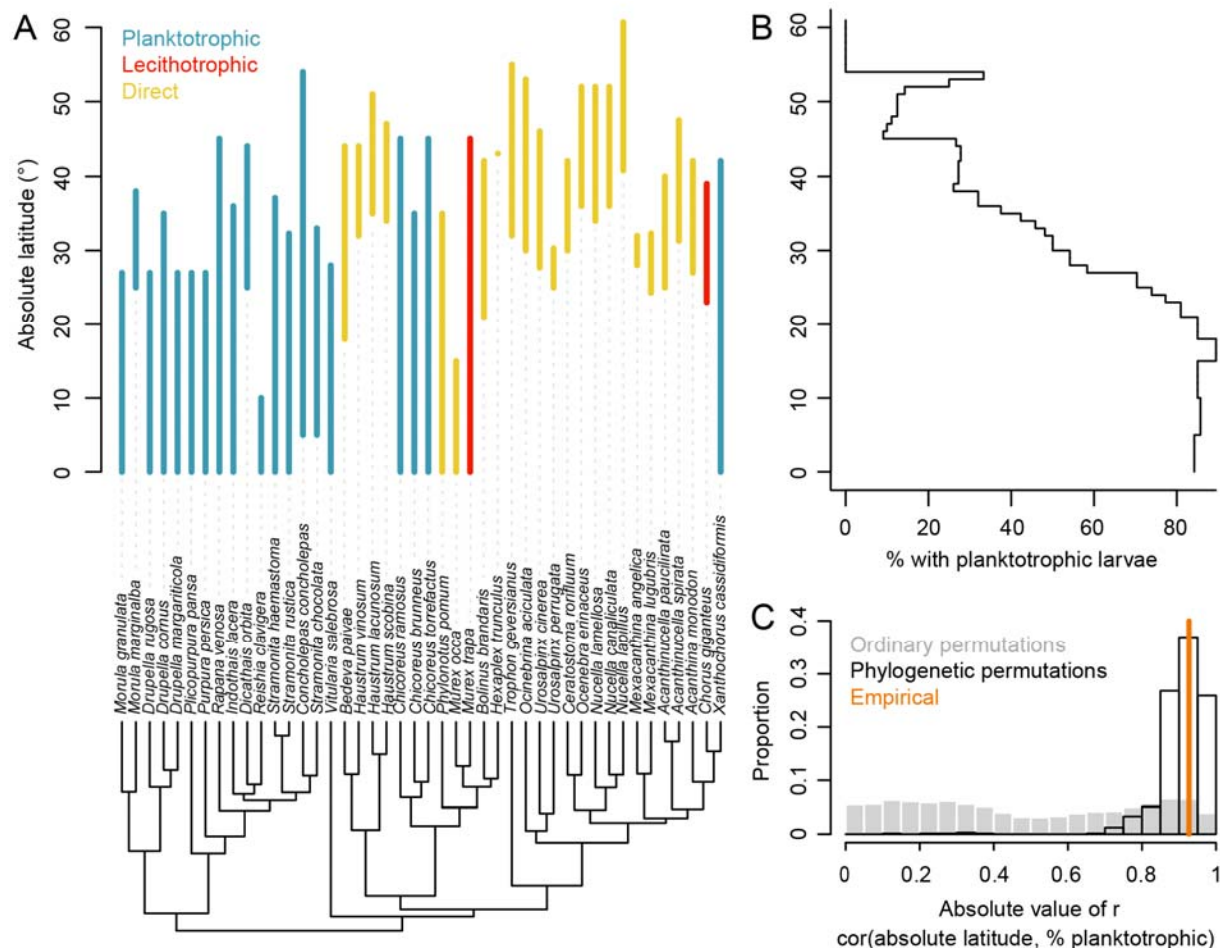523 permutations differ significantly from this distribution. All bins have width 0.05.

**Figure 3.** Phylogenetic permutations applied to Felsenstein's "worst case" in which two traits evolve independently on a tree whose shape tends to induce spurious correlations. **3A.** Scatterplot of traits X and Y with lines connecting values at the tips to ancestral state reconstructions. **3B.** Histograms showing the correlation between X and Y for sets of 1000

529 permutations of the variable Y with different approaches: ordinary permutations, signal-based
530 permutations with progressively smaller signal tolerances, and cyclic permutations. With a
531 tolerance of 0.01, signal-based permutations are statistically indistinguishable from cyclic
532 permutations. Vertical orange bars indicate the empirical correlation coefficient.



533

**Figure 4.** Phylogenetic permutation tests applied to theoretical morphospace occupation in
535 Triassic ammonoids. **4A.** Two parameters controlling shell geometry in 322 Triassic ammonoid
536 genera from McGowan (2004). Four theoretical ammonoid shells redrawn from Raup (1967)
537 illustrate how different shell geometries correspond to different combinations of these
538 parameters. Also plotted are the convex hull around the points, the smallest possible triangle
539 around the points, and the line $W = 1/D$, above which shells are open-coiled. **4B.** The empirical
540 number of genera with $W > 1/D$ compared with the same statistic for 1000 ordinary and 1000
541 signal-based phylogenetic permutations. This statistic was taken by Raup (1967) as evidence for
542 the reduced fitness of open-coiled forms. **4C.** The empirical ratio of the area of the convex hull
543 around the data to the area of the smallest triangle that fits around the data, compared with the
544 same statistic for 1000 ordinary and 1000 phylogenetic permutations. This metric of triangularity
545 was interpreted by Tendler et al. (2015) in light of Pareto optimality theory.

**Figure 5.** Phylogenetic permutation applied to Thorson's rule in muricid gastropods. **5A.** Mode of larval development (color-coded), absolute latitudinal range, and phylogeny for the 44 species from Pappalardo et al. (2014). **5B.** Thorson's rule plotted "as such": the percentage of species with planktotrophic larval development in each 1° bin of absolute latitude. **5C.** The correlation between absolute latitude and the percentage of planktotrophic species in each 1° latitudinal bin, shown for the empirical data, 1000 ordinary permutations of mode of larval development, and 1000 phylogenetic permutations.

**References**

556 Anderson M.J. 2001. Permutation tests for univariate or multivariate analysis of variance and
557 regression. Can. J. Fish. Aquat. Sci. 58:626–639.

559 Blomberg S.P., Theodore Garland J., Ives A.R. 2003. Testing for phylogenetic signal in
560 comparative data: behavioral traits are more labile. Evolution (N. Y). 57:717–745.

561 Cheverud J.M., Dow M.M., Leutenegger W. 1985. The Quantitative Assessment of Phylogenetic
562 Constraints: Sexual Dimorphism in Body Weight Among Primates. Evolution (N. Y).
563 39:1335–1351.

564 Collin R. 2003. Worldwide patterns in mode of development in calyptraeid gastropods. Mar.
565 Ecol. Prog. Ser. 247:103–122.

566 Collin R. 2004. Phylogenetic effects, the loss of complex characters, and the evolution of
567 development in calyptraeid gastropods. Evolution (N. Y). 58:1488–1502.

568 Colwell R.K., Hurtt G.C. 1994. Nonbiological gradients in species richness and a spurious
569 Rapoport effect. Am. Nat. 144:570–595.

570 Dushoff J., Kain M.P., Bolker B.M. 2019. I can see clearly now: Reinterpreting statistical
571 significance. Methods Ecol. Evol. 2019:3–6.

572 Edelaar P. 2013. Comment on "Evolutionary trade-offs, Pareto optimality, and the geometry of
573 phenotype space." Science (80-. ). 339:1–3.

574 Ewers-Saucedo C., Pappalardo P. 2019. Testing adaptive hypotheses on the evolution of larval
575 life history in acorn and stalked barnacles. Ecol. Evol. 9:11434–11447.

576 Felsenstein J. 1985. Phylogenies and the comparative method. Am. Nat. 125:1–15.

577 Gittleman J.L., Kot M. 1990. Adaptation: Statistics and a null model for estimating phylogenetic
578 effects. Syst. Zool. 39:227–241.

579 Good P. 2000. Permutation tests: A practical guide to resampling methods for testing hypotheses.
580 New York: Springer Science+Business Media.

581 Grafen A. 1989. The phylogenetic regression. Philos. Trans. R. Soc. B Biol. Sci. 326:119–157.

582 Ibáñez C.M., Rezende E.L., Sepúlveda R.D., Avaria-Llautureo J., Hernández C.E., Sellanes J.,
583 Poulin E., Pardo-Gandarillas M.C. 2018. Thorson's rule, life-history evolution, and
584 diversification of benthic octopuses (Cephalopoda: Octopodoidea). Evolution (N. Y).
585 72:1829–1839.

586 Ives A.R., Garland T. 2010. Phylogenetic logistic regression for binary dependent variables.

587    Syst. Biol. 59:9–26.

588  Jablonski D. 1986. Larval ecology and macroevolution in marine invertebrates. Bull. Mar. Sci.
589    39:565–587.

590  Jablonski D., Lutz R.A. 1983. Larval ecology of marine benthic invertebrates: paleobiological
591    implications. Biol. Rev. 58:21–89.

592  Lapointe F.-J., Garland T. 2001. A generalized permutation model for the analysis of cross-
593    species data. J. Classif. 18:109–127.

594  Lauder G. V. 1982. Historical biology and the problem of design. J. Theor. Biol. 97:57–67.

595  Mahler D.L., Ingram T., Revell L.J., Losos J.B. 2013. Exceptional convergence on the
596    macroevolutionary landscape in island lizard radiations. Science (80-. ). 341:292–296.

597  Marshall D.J., Krug P.J., Kupriyanova E.K., Byrne M., Emlet R.B. 2012. The Biogeography of
598    Marine Invertebrate Life Histories. Annu. Rev. Ecol. Evol. Syst. 43:97–114.

599  McGowan A.J. 2004. The effect of the Permo-Triassic bottleneck on Triassic ammonoid
600    morphological evolution. Paleobiology. 30:369–395.

601  Mileikovsky S.A. 1971. Types of larval development in marine bottom invertebrates, their
602    distribution and ecological significance: a re-evaluation. Mar. Biol. 10:193–213.

603  Mundry R. 2014. Statistical issues and assumptions of phylogenetic generalized least squares.
604    Modern Phylogenetic Comparative Methods and their Application in Evolutionary Biology.
605    p. 131–153.

606  Olalla-Tárraga M.Á. 2011. "Nullius in Bergmann" or the pluralistic approach to ecogeographical
607    rules: A reply to Watt et al. (2010). Oikos. 120:1441–1444.

608  Orzack S.H., Sober E. 2001. Adaptation, phylogenetic inertia, and the method of controlled
609    comparisons. In: Orzack S.H., Sober E., editors. Adaptationism and optimality. Cambridge,
610    UK: Cambridge University Press. p. 45–63.

611  Pappalardo P., Rodri□uez-Serrano E., Fernańdez M. 2014. Correlated evolution between mode
612    of larval development and habitat in muricid gastropods. PLoS One. 9.

613  Perezgonzalez J.D. 2015. Fisher, Neyman-Pearson or NHST? A tutorial for teaching data testing.
614    Front. Psychol. 6:1–11.

615  Raup D.M. 1966. Geometric analysis of shell coiling: general problems. J. Paleontol. 40:1178–
616    1190.

617  Raup D.M. 1967. Geometric analysis of shell coiling: coiling in ammonoids. J. Paleontol. 41:43–

618       65.

619   Rohde K., Heap M., Heap D. 1993. Rapoport's rule does not apply to marine teleosts and cannot

620       explain latitudinal gradients in species richness. Am. Nat. 142:1–16.

621   Saulsbury J.G., Baumiller T.K. 2020. Predation as an explanation for a latitudinal gradient in arm

622       number among featherstars. J. Biogeogr.:1–14.

623   Shibata T.F., Oji T. 2003. Autotomy and arm number increase in *Oxycomanthus japonicus*

624       (Echinodermata, Crinoidea). Invertebr. Biol. 122:375–379.

625   Sober E., Orzack S.H. 2003. Common ancestry and natural selection. Br. J. Philos. Sci. 54:423–

626       437.

627   Stearns S.C. 1984. The effects of size and phylogeny on patterns of covariation in the life history

628       traits of lizards and snakes. Am. Nat. 123:56–72.

629   Stevens G.C. 1989. The latitudinal gradient in geographic range: how so many species coexist in

630       the tropics. Am. Nat. 133:240–256.

631   Symonds M.R.E., Blomberg S.P. 2014. A primer on phylogenetic generalised least squares. In:

632       Garamszegi L.Z., editor. Modern Phylogenetic Comparative Methods and Their Application

633       in Evolutionary Biology. Berlin and Heidelberg: Springer-Verlag. p. 105–130.

634   Tendler A., Mayo A., Alon U. 2015. Evolutionary tradeoffs, Pareto optimality and the

635       morphology of ammonite shells. BMC Syst. Biol. 9:1–12.

636   Thorson G. 1950. Reproductive and larval ecology of marine bottom invertebrates. Biol. Rev.

637       25:1–45.

638   Uyeda J.C., Zenil-Ferguson R., Pennell M.W. 2018. Rethinking phylogenetic comparative

639       methods. Syst. Biol. 67:1091–1109.

640

641                  **Acknowledgements**

646

647                  **Data availability statement**

648     Code, supplementary data, and appendices are available at the following GitHub repository:

649     https://github.com/jgsaulsbury/phyloperm